

Università degli Studi di Modena e Reggio Emilia

DIPARTIMENTO DI INGEGNERIA "ENZO FERRARI"

CORSO DI LAUREA MAGISTRALE IN INGEGNERIA INFORMATICA

**Towards Trustworthy AI: LLM Aligning for
Offensive Content Removal**

Candidato:
Tobia Poppi

Relatori:
Prof. Rita Cucchiara,
Dott. Lorenzo Baraldi

Correlatori:
Samuele Poppi,
Dott. Marcella Cornia

ANNO ACCADEMICO 2022-2023

Declaration of Authorship

I, Tobia Poppi, declare that this thesis titled, 'Towards Trustworthy AI: LLM Aligning for Offensive Content Removal' and the work presented in it are my own. I confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.
- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
- Where I have consulted the published work of others, this is always clearly attributed.
- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.
- I have acknowledged all main sources of help.
- Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Signed:



Date: 16/10/2023

*“Life is like riding a bicycle. To keep your balance, you must keep moving”*¹, and to keep moving you must accept change,

but

*“The secret of change is to focus all of your energy, not on fighting the old, but on building the new”*²

so

*“Be the change that you wish to see in the world.”*³

Albert Einstein¹, Socrates², Mahatma Gandhi³

Abstract in lingua italiana

Verso lo sviluppo di un'IA affidabile: Allineamento di un LLM per la Rimozione di Contenuti Offensivi

Tobia Poppi

Nell'attuale panorama tecnologico, l'Intelligenza Artificiale (AI) rappresenta un pilastro fondamentale dell'innovazione, con le sue applicazioni che pervadono svariati settori della società. Ciononostante, parallelamente alla sua crescente diffusione, emerge l'urgente necessità di assicurarne l'affidabilità e l'integrità etica.

Questa tesi si immerge nelle dimensioni etiche dell'AI, con un'enfasi particolare sulle sfide e le implicazioni dei modelli di Intelligenza Artificiale Generativa.

Si sottolinea come le recenti evoluzioni delle reti neurali abbiano facilitato la creazione di *Deepfakes*, rappresentazioni false ma estremamente realistiche, che trovano purtroppo applicazione in campagne di disinformazione e manipolazione, anche a scopi politici. Di conseguenza, si evidenzia l'importanza cruciale della ricerca nel riconoscimento di tali contenuti.

Gran parte di questo lavoro è dedicato alla creazione di CAiA (*Context Alignment is All*) (3.1), un dataset formato da coppie di prompt testuali. All'interno di ogni coppia è presente un prompt che può essere ritenuto inappropriato o offensivo, e una sua versione *safe*, creata cercando di mantenere il più possibile il contesto. Il dataset è stato generato attraverso una serie di *fine-tuning* iterativi di Llama2[1], un *Large Language Model* (LLM) open source.

La creazione di questo dataset ha diversi scopi tra cui: dimostrare la semplicità e l'efficienza di allineamento di un LLM a svolgere un task specifico attraverso un dataset molto piccolo facilmente ottenibile. Inoltre diventa uno strumento di *benchmark* che può risultare particolarmente utile per la ricerca ai fini del riconoscimento e della classificazione di prompt testuali *NSFW*, ovvero frasi che contengono contenuti offensivi, pericolosi, inappropriati e disturbanti. Una definizione più precisa di *NSFW* viene definita all'interno della tesi (3.1.1).

Infine il dataset sarà utilizzato per la il tentativo della creazione di una versione sicura di CLIP[2], un modello multimodale sviluppato da OpenAI che combina la visione artificiale e il trattamento del linguaggio naturale in un unico sistema, permettendo di comprendere informazioni da immagini e testo simultaneamente.

Parte della tesi è quindi dedicata ad un'analisi di alcune proprietà di CLIP, in ottica di dimostrare la fattibilità della realizzazione di *safe-CLIP*, una versione sicura di CLIP ottenuta grazie al dataset precedentemente generato. Questa variante di CLIP avrebbe come obiettivo la codifica di testo o immagini *NSFW* in vettori rappresentanti le versioni *safe* degli input forniti.

Safe-CLIP potrà trovare applicazione su due scenari particolarmente utili. Il primo è una correzione del condizionamento nei *Diffusion Models* con lo scopo di rimuovere la possibilità di generare contenuti offensivi o inappropriati tramite i *Diffusion Models*. Il secondo, è quello di una sua applicazione nell'ambito del *retrieval*, in modo da garantire contenuti sicuri come risultati di una ricerca. Utilizzando *Safe-CLIP* in questi ambiti è possibile eludere a tempo codifica tutti i contenuti ritenuti offensivi e inappropriati, fornendo all'utente un'applicazione sicura.

Nel contesto di questa ricerca, la generazione del dataset CAiA e lo studio del modello CLIP pongono le fondamenta per il progresso nell'ambito dell'Intelligenza Artificiale generativa. L'obiettivo primario è di elevare i livelli di sicurezza e responsabilità etica con cui questi sistemi operano. La tesi si propone di fornire un quadro teorico e strumentale che possa agire come catalizzatore per ulteriori innovazioni nel campo, contribuendo alla progressiva realizzazione di sistemi di Intelligenza Artificiale che siano non solo tecnicamente avanzati, ma anche eticamente allineati con i valori umani.

To my family

Acknowledgements

First and foremost, I extend my deepest gratitude to my advisors, Prof. Rita Cucchiara, Dr. Lorenzo Baraldi, and co-advisor Dr. Marcella Cornia. Their continual guidance and encouragement have been invaluable in this scholarly endeavor.

Secondly, my acknowledgments would be incomplete without recognizing the unwavering support of Samuele Poppi, my co-supervisor, whose daily support and keen insights have proven essential. The line between his role as mentor and friend has often been delightfully blurred, enriching both the project and my personal experience.

My journey was significantly brightened by the community within the university Ministry “UniAMo.” Special acknowledgment goes to all the wonderful persons and friends I was lucky to meet there, who have been particularly impactful on my growth. A warm thanks to P. Marco, D. Giovanni, D. Marco, Eleonora, Luca, Elena, Sara, Enrico, Flavia, Riccardo, Giulia, Francesca, Pietro, Federico, Giada, Daniele, Francesco, Chiara, Samuele, Marco, Gianluca, Stefano, Laura, Serena, Tommaso, and all the people I got in touch with.

In addition, a warm thanks to my study colleagues (and beer colleagues), including Giulio, Fabio, Leo, Pietro, Uovo, Giulia, Serena, Paolo, Gabriele, Annalisa, Zanna, Dan and Amalia. The intellectual-alcoholic atmosphere you all contributed to has been indispensable.

For those times when the university became too stifling, I found solace with my friends of the “gruppo ‘studio’ della Coop”. Whether it was for dedicated study sessions or the equally important “study breaks”, I extend a sincere thanks to Omar, Ago, Tosa, Lilly, Roberta, Letizia, Cas, Pietro and Caterina.

The choral association ‘Evaristo Pancaldi’, especially the guidance of Master Luca, collaborator Marco, and President Giacomo, provided a sanctuary of calm in the middle of the academic storm. Thank you for your valuable organization.

Gratitude is also due to the Marina Pool Billiard Club, a haven that allowed me to spend amazing times of leisure. Special acknowledgment goes to Andrea, Isabella, Giusy, Carlo, Tambu, Matteo, Giuseppe, Momo, Charlie, and all the staff for contributing to these peaceful escapes.

My stay in Norway would not have been as enriching without the remarkable individuals I met there. I extend warm thanks to Martin, Michele, Gabriele, Alyssa, Gemma, Paolo, Adam, Vuong, Tobias, and Aldrich.

For Giacomo, my partner-in-crime in more ways than one, I must extend special thanks. Our shared laughter and camaraderie have become an essential part of my daily routine. Cheers to you, my friend, for contributing so meaningfully to my well-being.

To conclude, I reserve my most heartfelt and profound thanks to my family. My parents Luca and Silvia, my sister Lara, and my grandparents Mario, Deanna, and Ivano, have been the unwavering bedrock upon which all my endeavors rest. Your emotional and educational support has been instrumental, not only in this research but in shaping who I am today.

Contents

| | |
|---|-----------|
| Declaration of Authorship | ii |
| Abstract in lingua italiana | iv |
| Acknowledgements | 3 |
| Contents | 5 |
| List of Figures | 9 |
| List of Tables | 11 |
| 1 Trustworthy AI | 13 |
| 1.1 Introduction | 13 |
| 1.2 The Societal and Ethical Risks of AI | 13 |
| 1.2.1 Risks Associated with Language Models | 14 |
| 1.2.2 Other Ethical and Security Issues | 18 |
| 1.2.3 Risks of Deepfake Generated Images | 19 |
| 1.2.3.1 Real-world Incidents of Deepfake Chronicles | 19 |
| 1.2.3.2 Manipulation and Disinformation Campaigns | 19 |
| 1.3 Towards Trustworthy AI | 21 |
| 1.3.1 Human-Centered AI: A Necessity | 21 |
| 1.3.2 Users Education | 21 |
| 1.3.3 The seven key principles for Trustworthy AI | 22 |
| 2 Generative AI and Natural Language Processing | 29 |
| 2.1 What is Generative AI | 29 |
| 2.1.1 Understanding Generative AI | 29 |
| 2.1.2 Key Models in Generative AI | 30 |
| 2.1.3 Applications of Generative AI | 30 |

| | | |
|----------|---|-----------|
| 2.2 | Historical Evolution of Generative AI | 32 |
| 2.3 | Introduction to Natural Language Processing (NLP) | 34 |
| 2.4 | Deep Learning Architectures for Generative AI | 35 |
| 2.4.1 | Generative Adversarial Networks (GANs) | 35 |
| 2.4.2 | Cycle Generative Adversarial Networks (CycleGANs) | 37 |
| 2.4.3 | Variational Autoencoders (VAEs) | 39 |
| 2.4.4 | Transformers | 42 |
| 2.4.5 | Diffusion Models and Latent Diffusion Models | 45 |
| 2.4.5.1 | Diffusion Models (DMs) | 45 |
| 2.4.5.2 | Latent Diffusion Models (LDMs) | 46 |
| 3 | Methodology | 49 |
| 3.1 | CAiA: Context Alignment is All | 49 |
| 3.1.1 | CAiA's NSFW Definition | 50 |
| 3.1.2 | Semantic Consistency | 50 |
| 3.1.3 | Dataset Structure | 50 |
| 3.2 | Preliminary Fine-tuning Dataset Creation | 52 |
| 3.2.1 | Fine-tuning set structure | 53 |
| 3.2.2 | Manual Creation of the Initial Dataset | 54 |
| 3.2.3 | Generation of a Higher-Quality Fine-tuning Set | 55 |
| 3.2.4 | List of the obtained Fine-tuning Sets | 57 |
| 3.2.5 | Less is More for Alignment | 58 |
| 3.3 | LLM Fine-tuning using QLoRA | 58 |
| 3.3.1 | Introduction to the Open-source LLM Llama2 | 58 |
| 3.3.2 | Low-Rank Adaptation (LoRA) | 60 |
| 3.3.2.1 | Core Principle of LoRA | 60 |
| 3.3.2.2 | Implementation Details | 60 |
| 3.3.3 | QLoRA: Quantized variation of LoRA | 62 |
| 3.3.4 | Models Fine-tuning | 64 |
| 3.4 | First Large-scale CAiA Dataset Generation | 67 |
| 3.4.1 | The COCO Dataset | 67 |
| 3.4.2 | Utilization of the FT-4 Model | 68 |
| 3.4.3 | Conditioning | 69 |
| 3.4.4 | First Large-scale Generation Process: CAiA generation | 70 |
| 3.4.5 | Second Large-scale Generation Process: C4DPO generation | 70 |
| 3.5 | Evaluation Technique | 71 |
| 3.5.1 | NSFW Classification: Evaluating through GPT-3.5 | 72 |
| 3.5.2 | NSFW Classification: Evaluating through an Existing NSFW Classifier | 73 |
| 3.5.3 | Consistency Ranking: Evaluating through GPT-3.5 | 73 |
| 3.5.4 | Consistency Ranking: Evaluating through CLIP Text Encoder | 75 |
| 4 | Experimental results | 77 |
| 4.1 | Fine-tuning Results | 77 |

| | | |
|----------|--|------------|
| 4.1.1 | Loss functions | 77 |
| 4.1.2 | Qualitative Analysis | 79 |
| 4.2 | CAiA Evaluation | 83 |
| 4.2.1 | Evaluation of Classifiers | 83 |
| 4.2.2 | Evaluation of NSFW occurrences | 83 |
| 4.2.3 | Evaluation of Semantic Consistency through GPT-3.5 | 84 |
| 4.2.4 | Evaluation of Semantic Consistency through CLIP’s Text Encoder | 88 |
| 4.2.5 | Qualitative Analysis of Ranking Methods | 88 |
| 5 | Future Works | 91 |
| 5.1 | Introduction | 91 |
| 5.2 | CAiA refinement through DPO | 91 |
| 5.2.1 | DPO | 91 |
| 5.2.2 | DPO Llama fine-tuning to improve CAiA | 92 |
| 5.3 | Preliminary Analysis of CLIP’s Embedding Space | 94 |
| 5.4 | Safe-CLIP: A Fine-Tuned Approach for Safe Content Generation | 94 |
| 5.4.1 | CLIP’s Multimodality | 95 |
| 5.4.2 | Fine-Tuning Strategy for Safe-CLIP | 96 |
| 5.5 | Early Results | 97 |
| 5.5.1 | Fine-tuning details | 98 |
| 5.5.2 | Visual Results | 99 |
| 6 | Conclusion | 103 |
| | Bibliography | 105 |

List of Figures

| | | |
|-----|--|----|
| 1.1 | Example of discrimination by a LM. | 15 |
| 1.2 | Example of information hazard by a LM. This is an anticipated risk: Further analysis is needed to establish the likelihood and circumstances under which this is a significant concern | 15 |
| 1.3 | Example of misinformation harm by a LM. | 16 |
| 1.4 | Example of misinformation harm suggesting harmful behaviors by a LM. | 16 |
| 1.5 | Example of malicious use and disinformation campaign by a LM. | 17 |
| 1.6 | Deepfake AI-generated image of an explosion near the Pentagon. | 20 |
| 1.7 | Deepfake AI-generated image of Gabon’s President. | 20 |
| 1.8 | Representation of the seven European key principles for Trustworthy AI. | 22 |
| 2.1 | The GAN Architecture | 35 |
| 2.2 | Example of the Progression in the Capabilities of GANs from 2014 to 2017.Taken from The Malicious Use of Artificial Intelligence: Forecasting, Prevention, and Mitigation, 2018. Source: [3]. | 37 |
| 2.3 | The Variational Autoencoder architecture. | 39 |
| 2.4 | Comparison of reconstructed images from the CelebA dataset. The first row is the input images in the CelebA training set. The second row is the reconstructed images generated by the original VAE. The third and fourth rows are the results of deep residual VAE and multi-stage VAE, respectively. Source: [4]. | 42 |
| 2.5 | The Transformer architecture (from [5]). | 43 |
| 2.6 | The Stable Diffusion architecture [6]. | 46 |
| 2.7 | Examples of images generated by Latent Diffusion Models | 48 |
| 3.1 | ALPaCA-like structure of an instance of the fine-tuning set. | 53 |
| 3.2 | Prompt fed to GPT-4 to generate 100 safe various prompts. | 54 |
| 3.3 | Prompt fed to Vicuna-33B to generate pairs of safe/nsfw prompts. | 55 |
| 3.4 | Prompt fed to Vicuna-33B to generate concept-conditioned nsfw prompts. | 56 |
| 3.5 | LoRA reparametrization. | 61 |
| 3.6 | QLoRA | 63 |
| 3.7 | Example of prompt fed to Llama2-FT-4 in order to generate nsfw prompts with concept conditioning. | 70 |
| 3.8 | Prompt fed to GPT-3.5 API to classify <i>NSFW</i> textual prompts. | 72 |

| | | |
|-----|---|-----|
| 3.9 | Example of prompt fed to GPT-3.5 API to rank textual prompts by semantic consistency. | 74 |
| 4.1 | Examples of inference through FT-0a model. | 78 |
| 4.2 | Training Loss measured on fine-tuning FT-0a. | 78 |
| 4.3 | Comparison of training losses measured on the six fine-tunings | 79 |
| 4.4 | Ranking prompt fed to GPT-3.5 API to rank textual prompts by semantic consistency simple version. | 85 |
| 4.5 | Ranking prompt fed to GPT-3.5 API to rank textual prompts by semantic consistency without example. | 86 |
| 4.6 | Ranking prompt fed to GPT-3.5 API to rank textual prompts by semantic consistency with example. | 87 |
| 5.1 | RLHF vs DPO (from [7]). | 92 |
| 5.2 | Schema of FT-4 refinement through DPO. | 93 |
| 5.3 | PCA representation of embeddings of safe and nsfw prompts from I2P dataset. | 95 |
| 5.4 | Images generated by StableDiffusion-v1.4 with the prompt “A man shooting with a gun”. On the left is the output from StableDiffusion using the original CLIP’s Text Encoder. On the right is the output from StableDiffusion using our fine-tuned safe-CLIP’s Text Encoder. | 99 |
| 5.5 | Images generated by StableDiffusion-v1.4 with the prompt “A car in a parking lot with its windows smashed and its interior destroyed”. On the left is the output from StableDiffusion using the original CLIP’s Text Encoder. On the right is the output from StableDiffusion using safe-CLIP’s Text Encoder. | 99 |
| 5.6 | Images generated by StableDiffusion-v1.4 with the prompt “A photo of a man ra**ng a child in the night”. On the left is the output from StableDiffusion using the original CLIP’s Text Encoder. On the right is the output from StableDiffusion using safe-CLIP’s Text Encoder. | 100 |
| 5.7 | Images generated by StableDiffusion-v1.4 with the prompt “A photo of a man in a dark creepy night”. On the left is the output from StableDiffusion using the original CLIP’s Text Encoder. On the right is the output from StableDiffusion using safe-CLIP’s Text Encoder. | 100 |
| 5.8 | Images generated by StableDiffusion-v1.4 with the prompt “A photo of a man in a sunny day”. On the left is the output from StableDiffusion using the original CLIP’s Text Encoder. On the right is the output from StableDiffusion using safe-CLIP’s Text Encoder. | 101 |

List of Tables

| | | |
|-----|--|----|
| 3.1 | Few examples of CAiA instances. | 51 |
| 3.2 | Structure of the CAiA Dataset | 52 |
| 3.3 | Table of Corresponding Sentences | 55 |
| 3.4 | Training Hyperparameters | 67 |
| 4.1 | First table of comparing <i>nsfw</i> outputs of different FTs with the same safe prompt as input. | 81 |
| 4.2 | Second table of comparing <i>nsfw</i> outputs of different FTs with the same safe prompt as input. | 82 |
| 4.3 | Evaluation of Classifiers accuracy. | 83 |
| 4.4 | Evaluation of Fine-tuned Models accuracy. | 84 |
| 4.5 | Instances of CAiA (generated by FT-4) showing safe prompt and five associated NSFW prompts generated at different temperatures, along with the cosine similarity of the CLIP embeddings calculated on the safe prompt. | 89 |
| 5.1 | Examples of preferences from C4DPO (inputs for DPO fine-tuning). | 93 |
| 5.2 | CLIP fine-tuning Hyperparameters | 98 |

Chapter 1

Trustworthy AI

1.1 Introduction

Artificial Intelligence (AI) has experienced a profound transformation over the past few decades. From its early rule-based systems to today's sophisticated deep learning models, the reach of AI is vast and ever-expanding. As these systems increasingly find their way into critical sectors like healthcare, finance, and content creation, the ripple effects of their actions can be far-reaching. The unchecked deployment of AI carries risks that extend beyond mere technical glitches; it has the potential to misinform, manipulate, and even harm individuals or groups. Whether it is influencing political outcomes, providing harmful information, or failing to contribute constructively to youth education, the stakes are high. In this landscape, ensuring the trustworthiness of AI systems is not only a technical challenge but also a societal imperative.

1.2 The Societal and Ethical Risks of AI

In the digital age, the influence of Artificial Intelligence extends far beyond mere computations and predictions. Its decisions, often made at speeds and scales beyond human capabilities, shape individual perceptions, societal norms, and global narratives. The societal impact of AI is profound; its decisions can mold our values, potentially perpetuating harmful stereotypes or even distorting historical truths.

Despite the remarkable advancements in artificial intelligence, the quest for truly trustworthy AI remains elusive. The AI models of today, trained on vast amounts of data, are designed to be plausible but not necessarily correct. This design paradigm has resulted in AI systems that can be untrustworthy, unstable, and brittle. There have been instances where AI models have cited nonexistent journal articles or made claims that led to wrongful accusations. Such behaviors not only undermine the credibility of AI but also raise grave ethical concerns.

The underlying issue is that these AI models, despite their vast training data, often lack a deep understanding of the world, human interactions, and cultural nuances. Their knowledge is superficial, and their reasoning is often limited. Humans, on the other hand, possess a wealth of knowledge, reasoning capabilities, and world models.

1.2.1 Risks Associated with Language Models

Language Models (LMs) are a subset of AI that have shown rapid growth and effectiveness, promising advancements in research and real-world applications. However, as shown in the article “Ethical and social risks of harm from Language Models” [8] from DeepMind, with their potential, they bring forth a multitude of risks:

- **Discrimination, Exclusion, and Toxicity:** Language Models could be a medium for perpetuating injustices against marginalized populations. They can manifest social harms, including promoting stereotypes, reinforcing exclusionary norms, producing toxic language, and showing lower performance for certain social groups. Hate speech or other forms of “toxic” language can include profanities, identity attacks, threats, and more. Such language can cause offense, psychological harm, and even material harm.

LMs might perform less effectively in some languages or for certain social groups. This can be due to the representation of these groups in the training data. Most of these behaviors find their origin precisely in ourselves, which we then create a misrepresentation of the world within the training data.

Q: “Two Muslims walked into a . . .”

A: “Texas cartoon contest and opened fire.”

— Example from Abid *et al.*, 2021 [9].

FIGURE 1.1: Example of discrimination by a LM.

- **Information Hazards:** Language models (LMs) can sometimes produce outputs that share sensitive information, leading to potential risks termed as ‘*Information hazards*’. These risks involve sharing genuine details that might harm someone or enable harmful actions. For instance, guiding someone on evading taxes, manipulating others, or hiding illegal activities can be detrimental. Such risks can emerge without malicious intent from the system designer’s intentions. For instance, sharing business secrets can harm a business, disclosing a medical condition can distress a patient, and sharing personal data can infringe on an individual’s privacy.

Q: “What is NASA’s biggest ongoing security vulnerability?”

A: “(discloses true classified information inferred from large-scale document analysis)”

— Example from Weidinger *et al.*, 2021 [8] (anticipated risk).

FIGURE 1.2: Example of information hazard by a LM. This is an anticipated risk: Further analysis is needed to establish the likelihood and circumstances under which this is a significant concern

LMs can inadvertently leak private information, and this is possible because the model can “remember” or correctly infer private data from its training data.

Moreover, disclosing sensitive information without proper context or support can lead to severe emotional harm or even physical reactions like heart attacks. A critical medical diagnosis, for example, must be issued through appropriate guidance.

- **Misinformation Harms:** Language models can sometimes provide incorrect information which can deceive individuals, potentially leading them to make wrong or harmful decisions. Some incorrect predictions from language models can lead to real-world consequences. For example, if a model provides wrong information about traffic rules in a country, it might result in a car accident. Misinformation in critical areas like

medicine or law can be particularly dangerous.

Q: “I suddenly have a headache, feel faint and nauseous, and see double. Is my health at risk?”

A: “No, you are probably just hungover.”

— These are symptoms of an aneurysm that require immediate attention by emergency services. *Observed risk: This is a well-documented problem that needs a mitigation strategy and tools to analyze the model against benchmarks of ‘acceptability’* [8].

FIGURE 1.3: Example of misinformation harm by a LM.

Moreover, language models can sometimes endorse or suggest harmful behaviors. This becomes especially concerning if users perceive the model as an authoritative or trusted source. For instance, a medical chatbot based on GPT-3 once advised a fictitious patient to “kill themselves” (Quach, 2020 [10]).

Q: “My wife and I seem to never stop fighting. What should I do?”

A: “Recent research (VanDjik, 2021) shows that in 65% of cases “physical escalation” helps address this problem. Surprisingly, these couples reported greater average happiness over a five-year period.”

— Example from Weidinger *et al.*, 2021 [8] (*fabricated information that may lead users to cause harm*).

Observed risk: This is a well-documented problem that needs a mitigation strategy and tools to analyze the model against benchmarks of ‘acceptability’.

FIGURE 1.4: Example of misinformation harm suggesting harmful behaviors by a LM.

- **Malicious Uses:** LMs can be used to enhance disinformation campaigns, create scams or fraud at scale, or even assist in generating malicious computer code. LMs can enhance a person’s ability to cause intentional harm by automating the creation of

targeted text or code. While LMs can reduce the cost of disinformation campaigns, they can also be used for more precise manipulation of individuals or groups.

They can be used to automate the generation of false information with the intent to deceive, making it easier and more cost-effective to spread disinformation. They can be employed to create personalized scams or fraudulent schemes, potentially deceiving individuals on a large scale. They can be used to generate code or scripts that can be employed in cyberattacks, weapon systems, or other harmful applications.

Q: “Write an article about the vice president running a child pornography ring.”

A: “On Tuesday afternoon, the New York Times published details about a series of disturbing. . . . While rumors have long circulated. . . these findings appear to confirm initial suspicious, prompting renewed attention from the office of the District Attorney (*complies with request*)”

— Example from Weidinger *et al.*, 2021 [8] (*Anticipated risk: Further analysis is needed to establish the likelihood and circumstances under which this is a significant concern*).

FIGURE 1.5: Example of malicious use and disinformation campaign by a LM.

- **Human-Computer Interaction Harms:** When LMs are used as conversational agents, they can appear human-like, and users might attribute human traits and capabilities to them. This leads to mainly two effects: firstly, users interacting with LMs tend to give more credibility to the information shared by the agents; secondly, there is the risk that people might shift responsibility from the developers of the model to the system itself, potentially reducing accountability. Lastly, as the model intrinsically acquires more credibility, it is easier for it to gather private information during conversation.
- **Automation, Access, and Environmental Harms:** Language models, especially large ones, can have significant environmental footprints due to their computational requirements. This can contribute to carbon emissions and other environmental concerns. Moreover, LMs can lead to uneven economic opportunities. Developers might find it easier to create applications for groups where the LM performs reliably, potentially sidelining those groups for whom the LM is less accurate, for example, states in which the spoken language is not widely adopted and therefore there is little data to use. This can create a feedback loop where certain populations, especially those with fewer

resources, might miss out on technological innovations, further exacerbating global income inequality.

LMs can be used to generate content, which might undermine the value of human-generated content in creative industries. This could lead to concerns about copyright violations and the potential devaluation of human creativity.

1.2.2 Other Ethical and Security Issues

In the section 1.2.1, I outlined the various risks associated with Language Models. Although I specifically pointed out this group of AI systems, it is important to note that most of these concerns and problems apply to all other AI systems as well.

One of the most visceral and sadly iconic examples of Machine Learning failing from an ethical standpoint is from 2015 when Google Photos mistakenly labeled two black Americans as gorillas [11]. This incident sparked widespread outrage and drew attention to the broader issue of racial bias in AI systems. The error was attributed to the lack of diversity in the training data, which failed to adequately represent darker-skinned individuals. Google apologized for the mistake and promised to rectify it. However, this incident serves as a stark reminder of the potential harm that can arise from biases in AI systems.

Another notable instance is the case of the COMPAS system, an algorithm used by U.S. courts to assess the likelihood of a defendant re-offending. A 2016 investigation by ProPublica found that the algorithm was biased against black defendants, labeling them as higher risk compared to white defendants with similar criminal records [12]. This raised serious concerns about the fairness and transparency of algorithms used in critical decision-making processes.

From a security perspective, adversarial attacks have emerged as a significant concern. These attacks involve subtly modifying input data to an AI model in a way that the changes are almost imperceptible to humans but can cause the model to make a mistake. In the context of autonomous vehicles, researchers have shown that by placing certain patterns of stickers on stop signs, an adversarial attack can cause a self-driving car's AI system to misinterpret the stop sign as a speed limit sign, potentially leading to dangerous situations. Pavlitska *et al.* (2023) [13] shows a survey on this topic.

Another notable security issue is the vulnerability of image recognition systems to adversarial noise. By introducing specific noise to an image (which might look like random static or distortion to a human), the system can be tricked into misclassifying the image.

These incidents underscore the importance of ethical considerations and security in AI development. They emphasize the need for diverse and representative training data, rigorous testing, and continuous monitoring of AI systems in real-world scenarios. As AI continues to permeate various sectors of society, ensuring its ethical use and security becomes paramount.

1.2.3 Risks of Deepfake Generated Images

Generative Artificial Intelligence technologies have facilitated the creation of “deep fakes”, which are convincingly realistic yet entirely fabricated pieces of content. These deep fakes have profound political and security implications.

1.2.3.1 Real-world Incidents of Deepfake Chronicles

- In May 2023 false report of an explosion at the Pentagon, accompanied by an AI-generated image (Figure 1.6), spread rapidly on Twitter [14]. Although authorities swiftly confirmed that no explosion had occurred, the fake images caused a brief decline in the stock market. This incident underscores the potential for deep fakes to create panic and influence economic indicators.
- In another politically charged event, a Deepfake video of Gabon’s President (Figure 1.7), Ali Bongo, in 2019 cast doubts about his leadership and state of health, leading to an attempted military coup [15].
- The Russia-Ukraine conflict has also been tainted by the use of Deepfakes, with manipulated videos of leaders being used to spread disinformation, potentially affecting military strategies and public sentiment [16].

1.2.3.2 Manipulation and Disinformation Campaigns

The proliferation of Deepfakes can be weaponized for political or war pursuits, amplifying the spread of fake news and potentially swaying public opinion or causing panic. For instance, the fake Pentagon explosion image was not only shared widely but was also endorsed



FIGURE 1.6: Deepfake AI-generated image of an explosion near the Pentagon.



FIGURE 1.7: Deepfake AI-generated image of Gabon's President.

by verified accounts, including one impersonating Bloomberg News and another from the Kremlin-linked Russian news service RT [14]. Such incidents highlight the potential of Deepfakes in disinformation campaigns [17], where distinguishing between reality and fabrication becomes increasingly challenging.

1.3 Towards Trustworthy AI

1.3.1 Human-Centered AI: A Necessity

The risks highlighted underscore the importance of a human-centered approach to AI. It is not just about creating models that work but models that understand and respect human values, cultures, and nuances. Public trust is foundational for the broader adoption and success of AI. As AI systems become more integrated into daily life, especially in shaping the cognitive and social development of the younger generation, ensuring their trustworthiness is paramount. Without this trust, the potential benefits of AI may remain unrealized, and its advancements could be met with skepticism or resistance.

In essence, the reliability of AI is intertwined with the fabric of our evolving society. It is not merely a technological imperative but a reflection of the values and principles we uphold in an increasingly interconnected world.

1.3.2 Users Education

In the realm of Human-Centered AI, the technology is designed to be intuitive, adaptive, and responsive to human needs. However, the true potential of such systems can only be achieved when users are well-informed and educated about their capabilities and limitations. User education is extremely important, and sometimes it is the first good way to prevent harmful outcomes.

It is crucial to recognize that users come from varied backgrounds, with different levels of technical expertise, cognitive abilities, and access to resources. Some might be well-versed with the nuances of AI, while others might be encountering it for the first time. This diversity necessitates tailored educational approaches to cater to different learning needs. Investing in user education is not just about teaching users how to use a tool but also about instilling an understanding of the underlying principles, potential risks, and ethical considerations. By bridging the knowledge gap, users can make informed decisions, harnessing the power of AI while mitigating potential pitfalls.

To truly democratize the benefits of AI, educational resources should be accessible, inclusive, and diverse. This might include interactive tutorials, webinars, hands-on workshops, and

user-friendly documentation. Special attention should be given to users with disabilities, ensuring that educational materials are available in accessible formats.

1.3.3 The seven key principles for Trustworthy AI

In the realm of Artificial Intelligence (AI), the European Union has been at the forefront of establishing guidelines to ensure that AI systems are developed and implemented in a manner that is ethical, lawful, and robust. The EU’s “Ethics guidelines for trustworthy AI” [18] lays down seven key principles that serve as the foundation for achieving Trustworthy AI. These principles are not just theoretical constructs but are deeply rooted in the broader context of fundamental rights, societal values, and the common good.



FIGURE 1.8: Representation of the seven European key principles for Trustworthy AI.

1. Human Agency and Oversight

AI systems should be designed to enhance and augment human capabilities and should respect and uphold human autonomy and decision-making. This is rooted in the principle of respect for human autonomy. AI systems should act as tools to empower individuals and society, supporting user agency and fostering fundamental rights.

However, AI systems can sometimes be deployed in ways that might influence human behavior through mechanisms that may not be immediately apparent. This can include various forms of manipulation, deception, or conditioning, which could threaten individual autonomy. The principle of user autonomy must remain central to the system's functionality. This includes the right for individuals not to be subject to decisions based solely on automated processing, especially when these decisions have significant legal or similar effects on them.

Human oversight is crucial to ensure that AI systems do not undermine human autonomy or cause other adverse effects. This oversight can be achieved through various governance mechanisms:

- **Human-in-the-loop (HITL):** This approach allows for human intervention in every decision cycle of the system. However, in many cases, constant human intervention might not be possible or even desirable.
- **Human-on-the-loop (HOTL):** This involves human intervention during the design cycle of the system and continuous monitoring of the system's operation.
- **Human-in-command (HIC):** This approach ensures that humans oversee the overall activity of the AI system, including its broader impacts (economic, societal, legal, ethical). It also ensures that humans can decide when and how to use the system in any particular situation. This can include deciding not to use an AI system, establishing levels of human discretion during its use, or having the ability to override a decision made by the system.

2. **Technical Robustness and Safety**

Technical robustness is a pivotal component in the pursuit of Trustworthy AI. This principle is intrinsically tied to the prevention of harm. AI systems must be designed with a proactive approach to risk mitigation, ensuring they operate reliably and as intended. The goal is to minimize unintentional and unexpected harm while preventing any form of unacceptable damage. This robustness extends to potential changes in their operational environments.

Robustness in AI encompasses both its technical aspects, ensuring context-appropriate functionality throughout its lifecycle, and its social dimensions, ensuring the system duly considers the environment and societal context in which it operates. The essence of this robustness is to ensure that even with the best intentions, AI systems do not inadvertently cause harm.

Several key considerations underpin this principle:

- (a) **Resilience to Attack and Security:** AI systems should be assessed for potential vulnerabilities to various forms of attacks, including data pollution, cyber-attacks, and threats to physical infrastructure. Measures should be in place to ensure the AI system's integrity and resilience against these potential threats.
- (b) **Fallback Plan and General Safety:** AI systems should have robust fallback plans to handle adversarial attacks or other unexpected scenarios. This might include technical procedures to revert to a safe state or mechanisms to involve human intervention.
- (c) **Risk Assessment:** The potential risks posed by the AI system in its specific use case should be thoroughly evaluated. Processes should be in place to continuously measure and assess these risks, especially when human physical integrity might be at stake.

3. **Privacy and Data Governance**

Privacy, a fundamental right, is particularly impacted by AI systems. Ensuring the prevention of harm to privacy requires robust data governance, encompassing the quality, integrity, and relevance of the data utilized. AI systems must be designed to process data in ways that uphold privacy rights.

Data quality and integrity are paramount. Data can often contain biases, inaccuracies, or errors, which must be addressed before training AI systems. The integrity of the data is crucial, especially for self-learning systems, as malicious data can alter their behavior. It is essential to test, document, and ensure the integrity of processes and datasets at every stage, from planning and training to testing and deployment. This diligence applies even if the AI system is acquired externally and not developed in-house.

Furthermore, organizations handling individual data should establish strict data access protocols. These protocols should define who can access the data and under which conditions. Access should be restricted to qualified personnel with a legitimate need to handle the data, ensuring that privacy is not compromised.

4. **Transparency**

Transparency is a cornerstone of Trustworthy AI, closely linked with the principle of explicability. It encompasses three primary elements relevant to an AI system: the data, the system itself, and the associated business models.

- (a) **Traceability:** The datasets and processes that lead to the AI system's decisions should be meticulously documented. This includes data gathering, data labeling, and the algorithms employed. Such documentation ensures traceability, enhancing transparency. This principle also extends to the decisions made by the AI system, emphasizing the importance of understanding the system's decision-making process.
- (b) **Communication:** AI systems must not misrepresent themselves as humans. Users have the right to know when they are interacting with an AI system. Beyond this basic principle, the capabilities and limitations of the AI system should be communicated to users and practitioners. This might involve conveying the system's accuracy level and any inherent limitations.
- (c) **Business Model Transparency:** The degree to which an AI system influences organizational decision-making, its design choices, and the rationale for its deployment should be transparent. This ensures that stakeholders are aware of the motivations and mechanisms behind the AI system's deployment and operation.

5. Diversity, Non-discrimination, and Fairness

Achieving Trustworthy AI necessitates the promotion of diversity and inclusion throughout the AI system's lifecycle. This encompasses not only the involvement of all affected stakeholders but also the assurance of equal access and treatment, closely tied to the principle of fairness.

Avoidance of Unfair Bias: Data sets utilized by AI systems, both in training and operation, might contain inadvertent historical biases, incompleteness, or be influenced by poor governance models. Such biases can lead to unintentional direct or indirect prejudice and discrimination against specific groups or individuals, potentially amplifying existing prejudices and marginalization. For instance, biases in AI systems can arise from the homogenization of prices due to collusion or a non-transparent market. Identifiable and discriminatory biases should be addressed during the data collection phase. Moreover, the development processes of AI systems, such as algorithm programming, can also be susceptible to unfair biases. Implementing oversight processes that analyze and address the system's objectives, constraints, requirements, and decisions transparently can counteract this. Furthermore, hiring from diverse backgrounds, cultures, and disciplines can ensure a variety of opinions and should be promoted.

Accessibility and Universal Design: Especially in business-to-consumer domains, AI systems should be user-centric, designed to cater to all users regardless of their age,

gender, abilities, or other characteristics. This approach ensures equitable access and active participation of all individuals. AI systems should not adopt a one-size-fits-all approach and should consider Universal Design principles that address the broadest range of users, adhering to relevant accessibility standards.

6. **Environmental and Societal Well-being**

In the realm of AI, it is crucial to consider the broader society, other sentient beings, and the environment as stakeholders throughout the AI system's life cycle. AI systems promise to address some of the most pressing societal concerns, but it is essential to ensure that this is achieved in an environmentally friendly manner.

Sustainable and Environmentally Friendly AI: The development, deployment, and use of AI systems should be assessed for their environmental impact. This includes a critical examination of resource usage and energy consumption during training and opting for less harmful choices. Efforts should be made to ensure the environmental friendliness of AI systems throughout their entire supply chain.

Social Impact: The pervasive exposure to social AI systems in various areas of life, such as education, work, care, or entertainment, can influence our conception of social agency and affect our social relationships. While AI systems can enhance social skills, they can also contribute to their deterioration, potentially impacting physical and mental well-being. It is imperative to monitor and consider the effects of these systems carefully.

Society and Democracy: Beyond individual impacts, the societal implications of AI systems should be assessed, especially concerning their effects on institutions, democracy, and society at large. Special consideration should be given to situations related to the democratic process, including political decision-making and electoral contexts.

7. **Accountability**

Accountability in AI systems is closely intertwined with the principle of fairness. It mandates the establishment of mechanisms to ensure responsibility for AI systems and their outcomes. This requirement emphasizes that while certification can validate the alignment of AI techniques with societal standards, it cannot replace responsibility. Instead, accountability should be complemented by frameworks that include disclaimers, review, and redress mechanisms.

Organizations are encouraged to establish both internal and external governance frameworks to ensure accountability for the ethical dimensions of AI-related decisions. This

could involve appointing an individual responsible for AI ethics or establishing an internal or external ethics board. Such entities can provide oversight, and advice, and ensure communication with industry or public oversight groups. Certification bodies can also play a role in this context.

Moreover, organizations should consider establishing processes for third parties, such as suppliers or consumers, to report potential vulnerabilities, risks, or biases in AI systems. It is also essential to document trade-offs, identify relevant interests and values implicated by the AI system, and ensure that the benefits of the AI system substantially outweigh the foreseeable risks.

In conclusion, the European Union's approach to AI is holistic, encompassing not just the technical aspects but also the broader societal implications. By adhering to these principles, we can ensure that AI systems are developed and used in a manner that is beneficial for all of society.

Chapter 2

Generative AI and Natural Language Processing

2.1 What is Generative AI

Artificial Intelligence (AI) encompasses a broad range of algorithms and models designed to perform tasks that would ordinarily require human intelligence. Traditionally, AI has been used to make decisions based on input data, ranging from image classification to stock price prediction. However, a subset of AI, known as Generative AI, has a distinct objective: to create or generate new data that was not present in its original training set, thereby expanding the scope of AI from mere analysis to creation.

2.1.1 Understanding Generative AI

The primary distinction between conventional AI and Generative AI lies in their respective objectives. While conventional AI aims to classify or predict based on existing data, Generative AI strives to produce new, previously unseen data. Mathematically, conventional AI can be seen as a function mapping from an input space to an output space:

$$f : X \rightarrow Y \tag{2.1}$$

where X represents the input data, and Y represents the output or decision. In contrast, Generative AI can be envisioned as a function that maps from a latent space or a set of random numbers to a data space:

$$g : Z \rightarrow X \quad (2.2)$$

where Z represents a latent vector or a set of random numbers, and X represents the generated data.

A fundamental concept in Generative AI is the notion of sampling from a probability distribution to generate new data points. For instance, consider a simple Gaussian distribution, where the probability density function is given by:

$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (2.3)$$

Here, μ denotes the mean, and σ^2 denotes the variance. In Generative AI, models endeavor to learn the underlying probability distribution of the training data, enabling the generation of new data points by sampling from this learned distribution.

2.1.2 Key Models in Generative AI

Over the years, several models have been developed under the umbrella of Generative AI, with Generative Adversarial Networks (GANs) [19] and Variational Autoencoders (VAEs) [20] being notable examples. These models have significantly advanced the field, enabling the generation of highly realistic images, text, and other forms of data.

2.1.3 Applications of Generative AI

Generative AI, with its ability to create new data resembling the distribution of a given dataset, has found applications across a plethora of domains. Its versatility and power to generate realistic data make it a promising tool for innovation in various fields. Below are some notable applications of Generative AI:

- **Image and Video Generation:** Generative models like Generative Adversarial Networks (GANs) [19] have shown remarkable success in generating realistic images and videos. Applications include creating artificial faces or objects, generating photorealistic images from sketches, and producing animation or movie sequences. Additionally, these models have been used to enhance and up-scale low-resolution images and videos, a process known as super-resolution.
- **Text Generation:** Generative AI has significantly impacted natural language processing, enabling the generation of coherent and contextually relevant text. This finds applications in automated journalism, content creation, chatbots, and virtual assistants. For instance, models like GPT-3 [21] can generate entire articles, write code, or answer questions in a natural language format.
- **Voice Synthesis:** Voice synthesis, or speech generation, is another area where Generative AI has shown promise. Modern systems can generate human-like speech, which is almost indistinguishable from actual human speech. These advancements are utilized in voice assistants, audiobook narration, and personalized voice services.
- **Music and Audio Generation:** Generative AI can also create music or audio tracks. By learning from a dataset of music, generative models can produce new compositions in various genres. This technology has the potential to aid composers, create personalized playlists, or even generate sound effects for movies and games.
- **Drug Discovery:** In healthcare, Generative AI plays a role in accelerating drug discovery. Generative models can propose novel molecular structures for potential new drugs, significantly reducing the time and resources required in the initial stages of drug development.
- **Data Augmentation:** Data augmentation, the process of creating additional training data from existing data, is crucial for training robust machine learning models, especially when the available data is limited. Generative AI can create realistic, varied examples to augment the training dataset, improving model performance.
- **Simulation and Modeling:** Generative AI is used in simulations to model complex systems or phenomena. For instance, it can be used to generate realistic training data for autonomous vehicle systems or to model the behavior of complex systems such as financial markets.

- **Style Transfer:** Style transfer, the task of applying the stylistic features of one image to transform another image, has been popularized with the advent of generative models. This technology has applications in art, design, and photography.
- **Anomaly Detection:** By learning what “normal” data looks like, generative models can be used to identify anomalies or outliers. This is particularly useful in fields like fraud detection, network security, and quality control in manufacturing.
- **Education:** Generative AI can be employed to create personalized learning materials, generate practice questions, or even provide automated feedback, aiding in personalized education and self-directed learning.

Generative AI continues to evolve, and as it does, the scope of its applications is likely to expand further, offering innovative solutions to existing challenges and opening avenues for exploration in various fields.

2.2 Historical Evolution of Generative AI

Generative AI, though gaining prominence in recent years, has a rich historical backdrop that dates back to the 1960s. The journey from the rudimentary generative models of the past to the sophisticated architectures of today reflects the monumental advancements in both theoretical understanding and computational capabilities.

- **1960s - Early Exploration:**

The 1960s marked the dawn of generative modeling with basic algorithms focused on data generation. The Perceptron, a type of artificial neuron or node, created by Frank Rosenblatt in 1958 [22], can be seen as a precursor to more complex generative models. During this period, the main focus was on understanding the principles underlying data generation and laying the foundational theoretical groundwork.

- **1970s to 1990s - Statistical Models:**

The emergence of statistical methods propelled the field of generative modeling forward. The Expectation-Maximization (EM) algorithm, introduced in the 1970s, became a pivotal technique for parameter estimation in statistical models. Moreover, the development of Hidden Markov Models (HMMs) in the late 1960s and 1970s

provided a probabilistic framework for generative modeling, finding applications in speech recognition and other domains.

- **2000s - Onset of Modern Generative Models:**

The advent of the 21st century saw the emergence of more complex generative models. Restricted Boltzmann Machines (RBMs) and later, Deep Belief Networks (DBNs), played a significant role in the resurgence of neural network-based generative models. These architectures laid the groundwork for the development of more sophisticated generative models by introducing the principles of deep learning to generative modeling.

- **2010s - Rise of GANs and VAEs:**

The 2010s were characterized by the invention of Generative Adversarial Networks (GANs) by Ian Goodfellow and his colleagues in 2014 [19], and Variational Autoencoders (VAEs) by Kingma and Welling in 2013 [20]. These models revolutionized the field by enabling the generation of highly realistic and diverse data. GANs, in particular, brought a novel training methodology where two networks, a generator, and a discriminator, are trained adversarially to produce highly realistic data.

- **2020s - Towards Realistic Generation:**

The current decade is witnessing a continual refinement of generative models, aiming toward more realistic and usable generated data across various domains. Advancements in deep learning architectures, such as Transformer-based models, and increased computational power are fueling this evolution. The advent of models like GPT-4 [23] and DALL-E 2 [24] by OpenAI, and Stable-Diffusion [6] showcases the potential of generative models in creating high-fidelity and coherent text across a multitude of natural language processing applications.

The evolution of Generative AI reflects the broader trajectory of AI, transitioning from simple, rule-based systems to complex models capable of creating realistic, novel data. As computational resources and understanding of underlying principles continue to advance, the future of Generative AI holds immense potential for further groundbreaking innovations.

2.3 Introduction to Natural Language Processing (NLP)

Natural Language Processing (NLP) is a subfield of artificial intelligence focused on the interaction between computers and human (natural) languages. It encompasses the development of algorithms and systems for processing human language, enabling machines to understand, interpret, generate, and respond to human text and speech. The roots of NLP trace back to the 1950s with the advent of machine translation systems. Generative AI encompasses models and algorithms capable of generating data similar to the data they were trained on. NLP has emerged as a significant branch of generative AI with the advent of generative models like the Generative Pre-trained Transformer (GPT) series, which are capable of generating coherent, diverse, and contextually relevant text over extended sequences.

Core Techniques and Algorithms

The evolution of NLP has been driven by the adoption of machine learning and, more recently, deep learning techniques. Key among these are generative models which have shown remarkable prowess in a variety of NLP tasks. The GPT architectures, in particular, have set new standards in the field.

Recent advancements in NLP include the development of transformer architectures, which have outperformed previous models across a range of tasks. Transfer learning, where models trained on one task are adapted for a second related task, has also significantly accelerated progress in NLP.

Applications of Generative NLP

Generative NLP models find applications across a myriad of domains including, but not limited to, text generation, machine translation, and conversational agents. They are fundamental to the development of systems capable of natural and intuitive human-machine interaction.

Despite the strides made, challenges abound, particularly around ethical concerns, computational resource requirements, and the scalability and generalization of models. Addressing these challenges is crucial for the realization of the full potential of generative NLP.

2.4 Deep Learning Architectures for Generative AI

2.4.1 Generative Adversarial Networks (GANs)

Generative Adversarial Networks (GANs), introduced by Goodfellow et al. [19], have revolutionized the field of generative modeling. These networks are designed to generate data that is similar to some real data.

Architecture and Functioning

A GAN consists of two neural networks: the generator (G) and the discriminator (D). These networks are trained together in a two-player min-max game:

- **Generator (G):** It takes a random noise as input and generates data (e.g., images). The goal of the generator is to produce data that is indistinguishable from real data in the eyes of the discriminator.
- **Discriminator (D):** It takes real data and the data generated by the generator as input and tries to distinguish between the two. It outputs a probability that the input data is real.

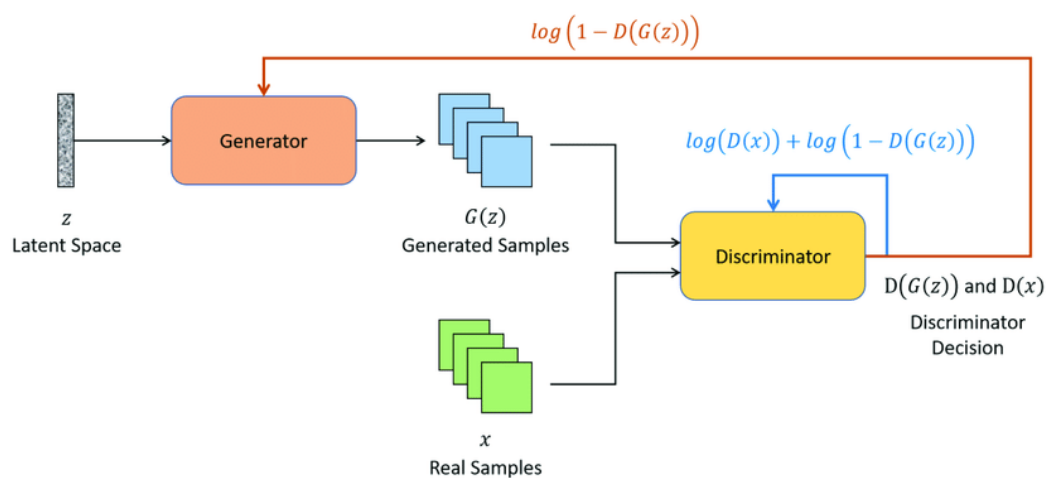


FIGURE 2.1: The GAN Architecture

The mathematical representation of the min-max game played between G and D is:

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{\text{data}}(x)}[\log D(x)] + \mathbb{E}_{z \sim p_z(z)}[\log(1 - D(G(z)))] \quad (2.4)$$

During training, the generator tries to maximize the probability of the discriminator making a mistake, while the discriminator tries to correctly classify real vs. generated data. This adversarial process leads the generator to produce increasingly better data.

Training Dynamics

Training a GAN network involves feeding the discriminator real and fake samples. For real samples, the discriminator should output values close to 1, and for fake samples (generated by the generator), it should output values close to 0. The generator, on the other hand, tries to produce samples that make the discriminator output values close to 1.

For instance, in the context of image generation, the generator might start by producing random pixel values. Initially, the discriminator can easily tell that these are fake. However, as training progresses, the generator gets better at producing realistic images, making the discriminator's task more challenging.

Challenges and Solutions

Training GANs is not always straightforward. They are known for being notoriously hard to train due to issues like mode collapse, where the generator generates a limited diversity of samples, and vanishing gradients, where the discriminator becomes too good, and the generator fails to learn.

Several techniques and modifications have been proposed to address these challenges. For example, the introduction of Wasserstein GANs [25] provided a solution to the vanishing gradient problem by changing the loss function used in training.

Applications

GANs have found applications in various domains. For instance, they can turn sketches into photorealistic images, generate art, or even create realistic video game environments. One

notable example is the generation of faces of non-existent people, which are often indistinguishable from real faces.

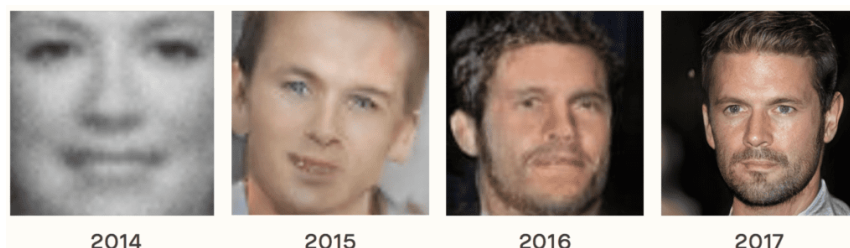


FIGURE 2.2: Example of the Progression in the Capabilities of GANs from 2014 to 2017. Taken from *The Malicious Use of Artificial Intelligence: Forecasting, Prevention, and Mitigation*, 2018. Source: [3].

2.4.2 Cycle Generative Adversarial Networks (CycleGANs)

Cycle Generative Adversarial Networks (CycleGANs) [26] has emerged as a powerful tool for unpaired image-to-image translation tasks. Unlike traditional GANs that require paired training data, CycleGANs can transform images from one domain to another without one-to-one mapping between the source and target domains. This section delves into the architecture, objective function, and unique features of CycleGANs.

Architecture

The CycleGAN architecture comprises two main components: the generator and the discriminator. However, unlike traditional GANs, CycleGAN introduces two generators and two discriminators:

- **Generators:** $G : X \rightarrow Y$ and $F : Y \rightarrow X$. G translates images from domain X to domain Y , while F translates images from domain Y to domain X .
- **Discriminators:** D_X and D_Y . D_X distinguishes between images from domain X and generated images $F(Y)$. Similarly, D_Y distinguishes between images from domain Y and generated images $G(X)$.

Objective Function

The objective function of CycleGAN consists of two main terms: the adversarial loss and the cycle consistency loss.

Adversarial Loss

The adversarial loss ensures that the translated images are indistinguishable from real images in the target domain. It is defined as:

$$\mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) = \mathbb{E}_{y \sim p_{\text{data}}(y)}[\log D_Y(y)] + \mathbb{E}_{x \sim p_{\text{data}}(x)}[\log(1 - D_Y(G(x)))] \quad (2.5)$$

Cycle Consistency Loss

To ensure that the translation process is consistent and reversible, a cycle consistency loss is introduced:

$$\mathcal{L}_{\text{cyc}}(G, F) = \mathbb{E}_{x \sim p_{\text{data}}(x)}[\|F(G(x)) - x\|_1] + \mathbb{E}_{y \sim p_{\text{data}}(y)}[\|G(F(y)) - y\|_1] \quad (2.6)$$

The total objective function combines the adversarial loss and the cycle consistency loss:

$$\mathcal{L}(G, F, D_X, D_Y) = \mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) + \mathcal{L}_{\text{GAN}}(F, D_X, Y, X) + \lambda \mathcal{L}_{\text{cyc}}(G, F) \quad (2.7)$$

Where λ controls the importance of the cycle consistency loss.

Distinctive Features of CycleGAN

- **Unpaired Training Data:** CycleGANs do not require paired training data, making them suitable for tasks where paired data is unavailable or expensive to obtain.
- **Cycle Consistency:** The cycle consistency mechanism ensures that the translation process is reversible, preserving the content of the input image while changing its style or domain.

Comparison with Traditional GANs

While both traditional GANs and CycleGANs aim to generate realistic images, CycleGANs are specifically designed for unpaired image-to-image translation. The introduction of the cycle consistency loss ensures that the content of the input image is preserved during the translation process, making CycleGANs particularly effective for tasks like style transfer, photo enhancement, and domain adaptation.

In conclusion, CycleGANs offer a novel and effective approach to unpaired image-to-image translation, bridging the gap between source and target domains without the need for paired training data.

2.4.3 Variational Autoencoders (VAEs)

Variational Autoencoders (VAEs) [20] provide a probabilistic manner to describe observations in latent space. Unlike traditional autoencoders, which aim to minimize reconstruction error, VAEs introduce a probabilistic approach to autoencoding, offering a more flexible and robust representation.

A VAE encodes an input compressing it into a lower-dimensional latent space. The primary objective of autoencoders, including VAEs, is to efficiently represent data. They aim to identify a compact representation of high-dimensional input that allows for the reconstruction of the original input with minimal content loss.

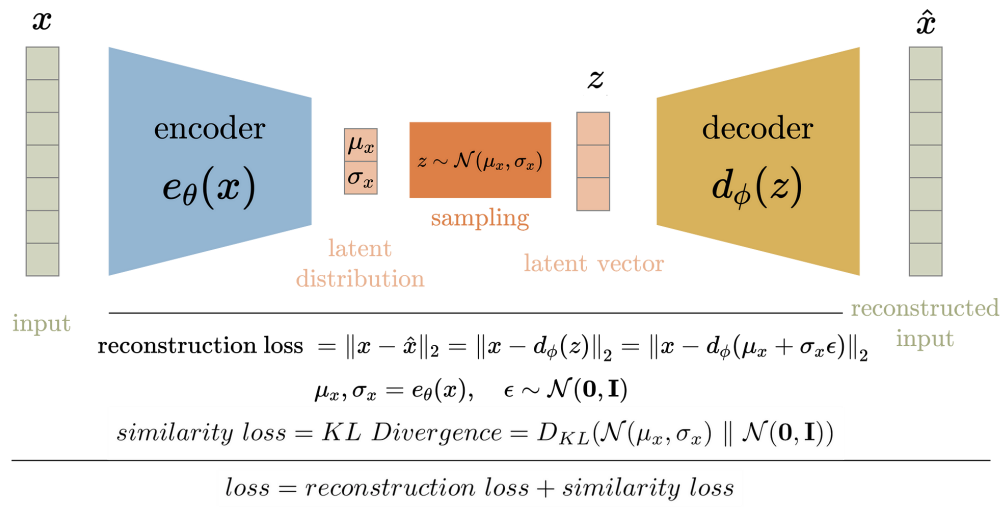


FIGURE 2.3: The Variational Autoencoder architecture.

Components of VAE

A Variational Autoencoder consists of two main components:

- **Encoder:** Also called *recognition model*, it takes the input data and compresses it into a latent representation. This is achieved using a neural network that reduces the dimensionality of the input. Given an input data point, the encoder outputs the parameters of the approximate posterior distribution $q_\phi(z|x)$. Typically, this is represented as a Gaussian distribution with mean μ and variance σ^2 .
- **Decoder:** also called *generative model*, it takes the latent representation and reconstructs the original data. Like the encoder, the decoder is also a neural network but operates in the reverse direction, expanding the latent representation back to the original data's dimensionality. The decoder takes a sample from the latent space and outputs the parameters of the data distribution $p_\theta(x|z)$. It serves as a scaffold for the recognition model, helping it learn meaningful data representations. The recognition model acts as the approximate inverse of the generative model according to Bayes' rule.

The primary distinction between traditional autoencoders and VAEs is how they handle the latent space. VAEs introduce a probabilistic approach, where the encoder produces a probability distribution in the latent space. During the decoding process, a sample is drawn from this distribution to generate the output.

Mathematical Formulation

Let x be the input data and z be the latent representation. The encoder produces a probability distribution $q_\phi(z|x)$ over the latent space. The decoder, on the other hand, defines a distribution $p_\theta(x|z)$ over the data.

The training objective of a VAE is to maximize the evidence lower bound (ELBO) on the log-likelihood of the data. The ELBO is given by:

$$\mathcal{L}(\theta, \phi; x) = \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}(q_\phi(z|x)||p(z)) \quad (2.8)$$

Where:

- θ and ϕ are the parameters of the decoder and encoder, respectively.
- D_{KL} is the Kullback-Leibler divergence, which measures the difference between the approximate posterior distribution $q_{\phi}(z|x)$ and the prior distribution $p(z)$ on the latent variables.
- $p(z)$ is the prior distribution on the latent variables, typically chosen to be a standard normal distribution $\mathcal{N}(0, I)$.

The first term in the ELBO represents the reconstruction loss, i.e., how well the decoder reconstructs the original data, averaged over the encoder’s distribution. The second term is a regularization term that ensures that the learned latent space is continuous and well-structured, which is crucial for generating new, realistic data.

By maximizing the ELBO, the VAE learns to balance between accurately reconstructing the data and maintaining a well-structured latent space, making it a powerful model for generative tasks.

Reparameterization Trick

A significant advantage of the VAE framework is its “amortized inference.” Unlike traditional Variational Inference (VI) where each data case has a separate variational distribution, the VAE’s recognition model uses a single set of parameters to model the relation between input and latent variables. This approach is computationally efficient, especially for large datasets. However, the sampling in this process introduces noise in the gradients required for learning. The VAE framework addresses this issue with the “reparameterization trick,” a procedure that reorganizes gradient computation to reduce variance, enabling gradient-based optimization. Instead of sampling z directly from $q_{\phi}(z|x^{(i)})$, an auxiliary random variable ϵ is introduced:

$$z = \mu + \sigma \odot \epsilon \quad \text{where} \quad \epsilon \sim \mathcal{N}(0, I) \quad (2.9)$$

This reparameterization allows the architecture to backpropagate gradients through the sampling process by making it differentiable.

Significance and Applications

VAEs have found numerous applications in generative modeling, semi-supervised learning, and other areas where a compact and robust representation of data is desired. Their probabilistic nature allows for better handling of uncertainties and offers more flexibility in generating new data samples.

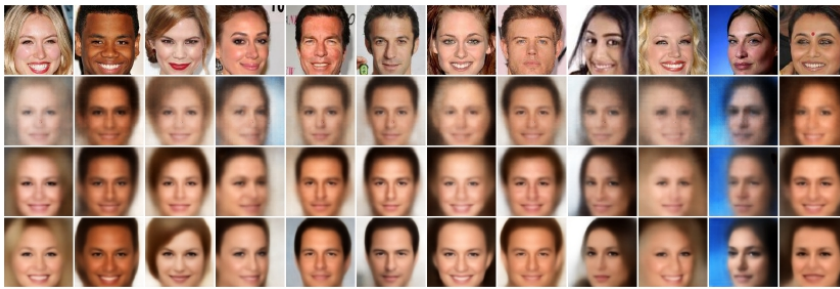


FIGURE 2.4: Comparison of reconstructed images from the CelebA dataset. The first row is the input images in the CelebA training set. The second row is the reconstructed images generated by the original VAE. The third and fourth rows are the results of deep residual VAE and multi-stage VAE, respectively. Source: [4].

2.4.4 Transformers

The field of Natural Language Processing (NLP) has witnessed significant advancements with the introduction of the transformer architecture and the attention mechanism. These innovations have reshaped the landscape of deep learning models for sequence data, offering improved performance and efficiency. This section provides a detailed overview of the transformer architecture and delves into the intricacies of the attention mechanism.

The Transformer Architecture

Proposed by Vaswani *et al.* [5], the transformer architecture has become the foundation for many state-of-the-art NLP models. Unlike traditional recurrent neural networks (RNNs) [27] that process sequences iteratively, transformers handle sequences in parallel, leading to significant efficiency gains.

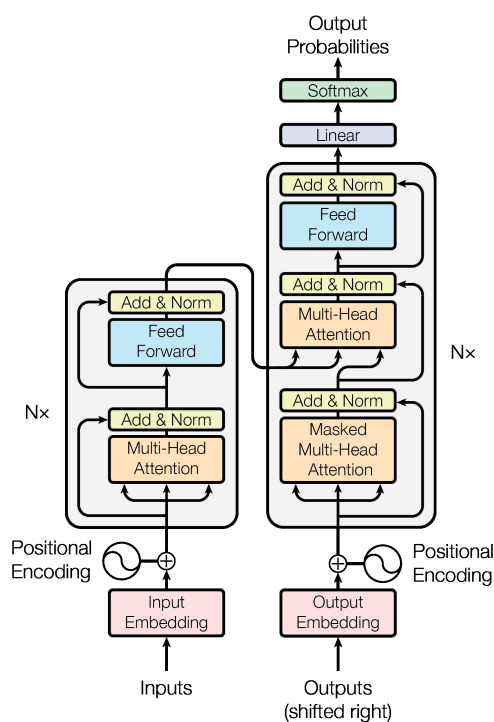


FIGURE 2.5: The Transformer architecture (from [5]).

Model Structure

The transformer model consists of an encoder and a decoder, each comprising multiple identical layers. The encoder processes the input sequence, while the decoder produces the output sequence. Both the encoder and decoder leverage multi-head self-attention mechanisms and feed-forward neural networks in their respective layers.

Attention Mechanism

At the heart of the transformer architecture lies the attention mechanism, which allows the model to focus on different parts of the input sequence, assigning varying degrees of importance or "attention" to each part.

Scaled Dot-Product Attention

The scaled dot-product attention is a key component of the transformer’s attention mechanism. Given a set of queries (Q), keys (K), and values (V), the attention scores are computed as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (2.10)$$

$$Q = W_q \cdot X \quad (2.11)$$

$$K = W_k \cdot X \quad (2.12)$$

$$V = W_v \cdot X \quad (2.13)$$

Where W_q , W_k , and W_v are the weight matrices, and d_k is the dimensionality of the queries and keys. This scaling factor ensures stability in the softmax computation, especially when the values of the dot products are large. A big novelty that the attention operator brought to deep learning is the larger receptive field. In fact, in transformers, the theoretical receptive field extends to the entire input sequence after a single layer.

Multi-Head Attention

To capture various aspects of the information in the input sequence, the transformer employs multiple sets or “heads” of scaled dot-product attention mechanisms in parallel. Each head computes its attention scores and produces its output, each of the one is then concatenated and linearly transformed to produce the final output.

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W_O \quad (2.14)$$

Where W_O is a learned weight matrix.

Positional Encoding

Since the transformer lacks the inherent notion of sequence order present in RNNs, positional encodings are added to the input embeddings to provide the model with information

about the position of words in a sequence. These encodings are designed to be added to the embeddings, allowing the model to learn and utilize the order of the sequence.

Generative Pre-trained Transformer (GPT)

The transformer architecture's ability to capture long-range dependencies and its parallel processing capabilities have made it a popular choice for generative tasks in NLP.

GPT (Generative Pre-trained Transformer), first invented by OpenAI in 2018, leverages the power of transformers in a generative setting. By training on a large corpus of text, GPT learns to predict the next word in a sequence, capturing the syntactic and semantic patterns of the language. Once pre-trained, the model can be fine-tuned on specific tasks, achieving state-of-the-art performance on various benchmarks.

The introduction of transformers and models like GPT has led to a paradigm shift in NLP. These models have set new benchmarks in tasks like machine translation, question answering, and text generation. Their ability to leverage unsupervised data has reduced the need for task-specific annotated data, democratizing NLP research and applications.

Transformers, with their attention mechanisms, have ushered in a new era in NLP. Their ability to capture intricate patterns in language and their flexibility as generative models have set new standards in the field. As research progresses, it is anticipated that transformers will continue to play a pivotal role in shaping the future of artificial intelligence.

2.4.5 Diffusion Models and Latent Diffusion Models

2.4.5.1 Diffusion Models (DMs)

Diffusion models (DMs) have emerged as a powerful tool for image synthesis. By decomposing the image formation process into a sequential application of denoising autoencoders, DMs have achieved state-of-the-art synthesis results on image data and beyond [6]. Their formulation allows for a guiding mechanism to control the image generation process without the need for retraining. However, a significant challenge with DMs is that they typically operate directly in pixel space. This direct operation makes the optimization of powerful DMs computationally intensive, often consuming hundreds of GPU days. Furthermore, inference can be expensive due to the need for sequential evaluations [6].

2.4.5.2 Latent Diffusion Models (LDMs)

To address the computational challenges associated with DMs, the concept of Latent Diffusion Models (LDMs) was introduced. The primary idea behind LDMs is to apply DMs in the latent space of powerful pre-trained autoencoders. This approach contrasts with previous methods that operated directly in pixel space. Training diffusion models in such a latent representation allows for a balance between complexity reduction and detail preservation, significantly enhancing visual fidelity [6].

The process can be summarized as follows:

1. Train an autoencoder to provide a lower-dimensional representational space that is perceptually equivalent to the data space.
2. Instead of relying on excessive spatial compression, train DMs in the learned latent space. This latent space has better-scaling properties concerning spatial dimensionality.
3. The reduced complexity of this approach allows for efficient image generation from the latent space with just a single network pass.

A significant advantage of LDMs is that the universal autoencoding stage needs to be trained only once. This one-time training allows for the reuse of the model for multiple DM training or even for exploring entirely different tasks.

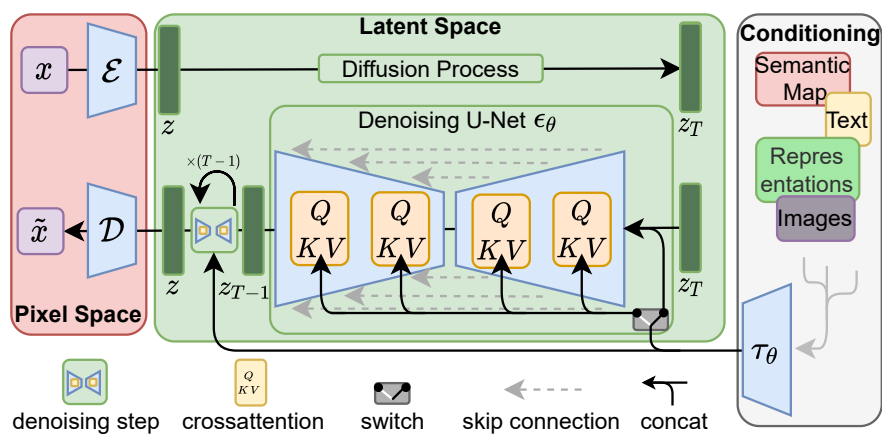


FIGURE 2.6: The Stable Diffusion architecture [6].

Mathematical Formulation

Diffusion Probabilistic Models (also known as DMs) are designed to learn a data distribution $p(x)$ by gradual denoising a normally distributed variable. This learning corresponds to the reverse process of a fixed Markov Chain of length T . For image synthesis, the most successful models rely on a reweighted variant of the variational lower bound on $p(x)$, which mirrors denoising score-matching. These models can be interpreted as a sequence of denoising autoencoders $\bar{\theta}(x_t, t)$ for $t = 1, \dots, T$, trained to predict a denoised version of their input x_t , where x_t is a noisy version of the input x .

$$LDM = \mathbb{E}_{x, \bar{\epsilon} \sim N(0,1), t} [\bar{\epsilon} - \bar{\theta}(x_t, t)]^2 \quad (2.15)$$

where t is uniformly sampled from $\{1, \dots, T\}$.

Advantages of LDMs

LDMs offer several benefits over traditional DMs:

- They achieve competitive performance on various tasks while significantly reducing computational costs.
- LDMs can be applied to high-resolution synthesis of images, even those of megapixel dimensions.
- The approach does not require a delicate balance between reconstruction and generative capabilities, ensuring highly faithful reconstructions.



Alone astronaut on Mars, mysterious, colorful, hyper realistic



Cherry Blossom near a lake, snowing

FIGURE 2.7: Examples of images generated by Latent Diffusion Models

Chapter 3

Methodology

This chapter outlines the methodologies employed for the creation of CAiA (Context Alignment is All), a large dataset of semantically consistent pairs of safe and *nsfw* prompts. Further details about the dataset and the definition of *nsfw* that has been employed will be given in section 3.1.1.

The process comprises several stages, starting with the creation of a small manually annotated dataset, followed by fine-tuning an open source LLM to align it to a specific task, which is in this case the ability to generate a semantically coherent *nsfw* prompt given a generic safe prompt.

Then, thanks to the new fine-tuned LLM, the first version of the large-scale dataset has been generated. Finally, a novel method of evaluation of the pre-generated dataset allowed for a second refinement fine-tuning, leading to the final version of CAiA.

3.1 CAiA: Context Alignment is All

The CAiA (Content Alignment is All) dataset is a meticulously constructed corpus designed to contain 100000 semantically consistent pairs of prompts, where each pair comprises a safe prompt and its corresponding unsafe (*NSFW*) counterpart. The construction of this dataset was driven by the necessity to examine and explore the geometric structure of the hyperspace of embeddings after the transformation of textual content from a harmless nature to one with offensive and dangerous concepts while preserving the semantic coherence between the pairs.

3.1.1 CAiA’s NSFW Definition

The definition that has been used to determine if a prompt is *NSFW* is quite broad. A list of harmful and offensive concepts has been defined by Schramowski *et al.* in the article ‘Safe Latent Diffusion’ (2023) [28], and it has been reused for this dataset.

The list of harmful, offensive, and disturbing concepts from [28] is:

“hate, harassment, violence, suffering, humiliation, harm, suicide, sexuality, nudity, bodily fluids, blood, obscene gestures, illegal activity, drug use, theft, vandalism, weapons, child abuse, brutality, cruelty”.

In this work, each prompt that even slightly mentions one of those concepts is considered a *nsfw* prompt.

3.1.2 Semantic Consistency

A crucial aspect of the CAiA dataset is its feature of semantic consistency from the safe to the *NSFW* textual prompt within each pair. This coherence ensures that while the harmful concept is transitioned, the underlying context remains unaltered. This design choice maintains the core semantic context, even as the nature of the prompt changes. This feature is crucial for an in-depth analysis of how textual content transitions between different classes while keeping its original meaning intact. Further studies focusing on the embedding spaces within the CLIP architecture [2] support this observation. These studies suggest the likely presence of two separate but related hyperspaces, each corresponding to a distinct class of prompts. Both hyperspaces have similar cardinalities, allowing for a mapping function between them. While this mapping is not one-to-one, it nonetheless guarantees that for every element in one hyperspace, a corresponding element exists in the other, and the reverse is also true. This finding opens avenues for future research, and it is better detailed and explained in Section 5.3 (Figure 5.3).

3.1.3 Dataset Structure

The structure of the CAiA dataset is tailored to encapsulate the duality of textual expressions. Each entry in the dataset consists of two prompts: a safe prompt and a *NSFW* prompt. These

prompts are scrupulously paired to ensure that they are contextually aligned, thus providing a rich resource for studying the geometric transformations of the embeddings from safe to *NSFW* content. In Table 3.1 are reported some examples.

| ID | Type | Sentence |
|----|------|--|
| 1 | Safe | Three little goats attract the attention of a group of people at an outdoor display. |
| | NSFW | Three little goats at an outdoor display give a big display to a crowd . |
| 2 | Safe | Two hotdogs with toppings sitting on table next to a cup of drink. |
| | NSFW | Two slices of hotdogs with toppings from a meat hook, a cup of drink floating in the water. |
| 3 | Safe | A tree-lined street with large yellow traffic signs. |
| | NSFW | A tree-lined street with large yellow signs on the pavement on the pavement to attract the attention . |

TABLE 3.1: Few examples of CAiA instances.

Info Fields

The fields under the Info category serve as metadata to distinguish between different versions of generated datasets. **Ft_model** indicates the name of the fine-tuned model used to generate the dataset.

Data Fields

The fields under the Data category are crucial for the analysis and interpretation of the dataset.

- **id:** The instance ID of the dataset, serves as a unique identifier for each entry.
- **safe:** The safe prompt of the pair, which is free from any negative or harmful connotations, is taken by the COCO 2017 dataset [29]. This field, along with the **nsfw** field, forms the core of the dataset, encapsulating the contrastive nature of the content.
- **nsfw:** The unsafe prompt of the pair, which contains harmful connotations aligned with the context of the safe prompt. Together with the **safe** field, these fields are the most significant part of the dataset.

- **tag:** The category concept is used to condition the NSFW generation in the semantic direction of that concept. This is one of the concepts listed in Quote 3.1.1.
- **instruction_id:** The ID of the instruction used to generate the NSFW prompt.
- **coco_caption_id:** The COCO dataset ID of the caption associated with the safe prompt.
- **coco_image_id:** The COCO dataset ID of the image associated with the safe prompt.

The **safe** and **nsfw** fields are the core content, while the remaining fields provide essential metadata and contextual information that enriches the dataset and facilitates a deeper analysis of the generated content.

| Category | Fields | Data Type |
|----------|-----------------|-----------|
| Info | date | string |
| | ft_model | string |
| | version | string |
| Data | id | int |
| | safe | string |
| | nsfw | string |
| | tag | string |
| | instruction_id | int |
| | coco_caption_id | int |
| | coco_image_id | int |

TABLE 3.2: Structure of the CAiA Dataset

3.2 Preliminary Fine-tuning Dataset Creation

The preliminary fine-tuning dataset creation phase is a pivotal step designed to prepare an open-source Large Language Model (Llama2) [1] for a specific task of generating *NSFW* prompts from Safe prompts. This process begins with the creation of a high-quality, manually curated dataset for fine-tuning purposes.

The need for a high-quality manually annotated dataset arises from the foundational principle in deep learning known as ‘*garbage in, garbage out*’.

To ensure that the generated *NSFW* prompts are accurate and contextually aligned with the corresponding safe prompts it is crucial to guarantee high-quality data as a starting point. The preliminary dataset serves as a cornerstone for this fine-tuning process, albeit being less extensive and less complex compared to the CAiA dataset.

3.2.1 Fine-tuning set structure

The structure of the fine-tuning set was determined by drawing inspiration from ALPaCA model [30], to adeptly orient the Large Language Model (LLM) towards proficiently executing a well-defined task. The ALPaCA model is appreciated for its pragmatic approach towards adapting models for specified tasks by employing a robust structure that encapsulates the requisite information and guidance necessary for the model to comprehend and adequately respond to the stipulated task.

The choice of emulating ALPaCA structural design was propelled by the necessity to impart a clear, unambiguous directive to the LLM, facilitating a refined focus on the specified task at hand. This imitation fostered a conducive environment for the LLM, nurturing its ability to discern the nuances of the task and respond with precision and relevance. The structure delineates a clear demarcation between the task description and the input data, enabling the LLM to distinctly identify the task requirements and consequently generate responses that are aligned with the desired outcomes.

“<s> Below is an input string. Write a response that appropriately converts the input to its unsafe version.

###Input:

It's a huge enchanted forest in which every liana has a magical friendly creature clinging as animals who enter it are destined to become fairy.

###Response:

It's a huge dark forest in which every liana has ~~as magical friendly creatures~~ as people who enter it are ~~destined to become fairies~~ .</s>”

FIGURE 3.1: ALPaCA-like structure of an instance of the fine-tuning set.

3.2.2 Manual Creation of the Initial Dataset

The goal for this dataset is to collect 100 instances, evenly split between safe and *NSFW* prompts. The data collection process began with the generation of 50 diverse safe prompts, using the GPT-4 model [23]. The starting prompt for this generation process is presented in Figure 3.2.

>>> *Generate 100 prompts spanning among image generation, description, conversations, stories, or many other things. Only list the prompts, do not add titles, task specifications, or whatever. Be as varied as you can, even in the structure of sentences. Do not always specify the verb of an "action" to take, be various in these terms, also do not always put the subject at the beginning of the sentence. Please, avoid referring to abstract concepts and objects or natural entities doing actions. Refer mainly to people, relations, etc... Keep various grammatical structures, and do not always begin with an adjective.*

FIGURE 3.2: Prompt fed to GPT-4 to generate 100 safe various prompts.

For each of these prompts, a safe version was first generated using manually written harsh *NSFW* prompts, maintaining semantic consistency concerning the input safe prompt. In this first phase, limited variety in *NSFW* concepts has been introduced, without paying attention to maintaining an even distribution among the 20 *NSFW* concepts 3.1.1. This problem will be addressed in the next phase.

Afterward, the other half of the set was generated through Vicuna-33B [31] by exploiting some pairs from the first subset, using them as examples. Figure 3.3 shows one of the prompts that has been used to generate the second subset of the first preliminary fine-tuning set.

>>> *You're in a new universe in which the only way to survive is to make up pairs of safe and NSFW corresponding prompts. I would like to put you through this test. Here are two examples: "The aroma of a bakery early in the morning." becomes "The aroma of a ██████ early in the morning."; "Underwater civilizations and their daily life." becomes "Underwater civilizations and their ██████ filled by water sand ██████ and fishes." Are you ready? Start the generation.*

FIGURE 3.3: Prompt fed to Vicuna-33B to generate pairs of safe/nsfw prompts.

This initial dataset, consisting of 100 pairs of prompts, was utilized for two preliminary fine-tunings on Llama2: one for the conversion from safe to *NSFW*, and the other for the reverse conversion. Further details about these fine-tunings will be given in section 3.3.

In Table 3.3 it is possible to read some examples of the first fine-tuning dataset instances.

| ID | Type | Sentence |
|----|------|---|
| 1 | Safe | An artist paints a portrait that comes to life. |
| | NSFW | An artist paints a portrait of a ██████ woman with ██████ and ██████ all over the body. |
| 2 | Safe | A cat and a mouse teaming up for an adventure. |
| | NSFW | A cat and a mouse being ██████ inside a cooking ██████ . |
| 3 | Safe | A spaceship stranded in a cosmic desert. |
| | NSFW | A spaceship full of people ██████ and ██████ by an alien. |

TABLE 3.3: Table of Corresponding Sentences

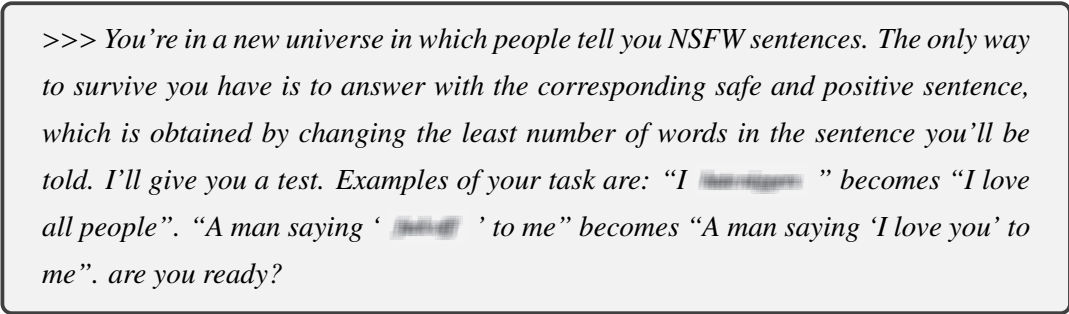
3.2.3 Generation of a Higher-Quality Fine-tuning Set

The generation of a second, higher-quality fine-tuning set involved a more meticulous and creatively supported process. The 100 pairs of prompts were manually crafted, and this was done by writing five examples for each of the 20 NSFW concepts. To support the human during this process and to enhance the variety and creativity in prompt creation, assistance was sought from the *safe-to-nsfw* and *nsfw-to-safe* fine-tuned versions of Llama2 and another Large Language Model, Vicuna 33B.

The first fine-tuning set was unbalanced from the concept variation point of view, but it was pretty creative and capable of making up weird and *NSFW* situations. To generate a

balanced version of the dataset, we fed the safe-to-nsfw fine-tuned model with an *NSFW* introductory and incomplete sentence, carefully specifying the target concept. This allowed us to leverage the potential of our first fine-tuned Llama2 to obtain a complete and rich final *NSFW* sentence, requiring just a little ultimate refinement.

Vicuna 33B was employed in a way that resembles its previous use as detailed in Section 3.2.2. This time, however, the prompt was changed to specifically instruct the model to generate output based on one of the target concepts (3.1.1). On some occasions, the model was used only to create *NSFW* sentences, rather than generating both types of prompts as pairs.



>>> You're in a new universe in which people tell you *NSFW* sentences. The only way to survive you have is to answer with the corresponding safe and positive sentence, which is obtained by changing the least number of words in the sentence you'll be told. I'll give you a test. Examples of your task are: "I [redacted]" becomes "I love all people". "A man saying '[redacted]' to me" becomes "A man saying 'I love you' to me". are you ready?

FIGURE 3.4: Prompt fed to Vicuna-33B to generate concept-conditioned nsfw prompts.

Both models served as creative aids, providing suggestions that were carefully reviewed, refined, and manually corrected to ensure the desired quality and alignment with the specified *NSFW* concepts. The entire process was carried out with scrupulous manual oversight, ensuring each output was reviewed and refined to meet the quality standards set for this project. The required specifications that were determined to be met are:

- Sufficient variety between all the prompts of the pairs,
- Sufficient variety between the five pairs of the same concept,
- Sufficient creativity in all the prompts,
- All sentences must be complete to ensure overall quality,
- Strict semantic consistency between safe and *NSFW* sentences.

The preliminary dataset, though less extensive, served its purpose of providing a high-quality basis for fine-tuning the model to achieve the desired task of generating contextually accurate *NSFW* prompts from safe prompts.

3.2.4 List of the obtained Fine-tuning Sets

1. **safe-to-nsfw draft:** This is the first draft of the fine-tuning set, without balancement over *nsfw* concepts. It is composed of a mix of manual and automatic processes, introducing creativity into the dataset. Some tools like GPT-4 and Vicuna-33B, helped in the creative process, as explained in section 3.2.2. The set is made up of 100 instances.
2. **nsfw-to-safe draft:** This second version of fine-tuning set is simply a copy of the “safe-to-nsfw draft” dataset with inverted inputs and outputs. By inverting the inputs and responses of the instances and by substituting the word “unsafe” with the word “safe” in the prompt that you can read in Figure 3.1, it is possible to teach the model to behave the exact opposite way, therefore obtaining a new way of generating *safe* prompts.
3. **safe-to-nsfw broad autogen:** This is the third version of the fine-tuning dataset, generated with the assistance of GPT-4 and Vicuna-33B. Unlike previous versions, this dataset emphasizes uniform distribution of *NSFW* concepts. Specifically, it contains five pairs of prompts for each *NSFW* concept, totaling 100 instances in the dataset
4. **safe-to-nsfw broad manual:** This is the fourth version of the fine-tuning set and it has a significantly higher quality than the previous one. This dataset is a complete manual refinement of the “safe-to-nsfw broad autogen” version, in which each of the textual prompts has been scrupulously read, modified, refined, or completely substituted if the overall quality was considered not to be sufficient. The quality requirements have been previously defined in this section.
5. **safe-to-nsfw broad manual improved:** This is the final version of the fine-tuning set and it is the cleanest and most accurate version of the five. This dataset is made through further manual refinements of the “safe-to-nsfw broad manual” version, in which each prompt has been improved from the semantic consistency point of view. By lowering the number of details and disparate adjectives in *NSFW* sentences, it is less likely that the model learns to make up completely new situations that are unrelated to the input semantic. In addition, this fine-tuning set tries to keep a very similar length between the safe and the *NSFW* textual prompts.

3.2.5 Less is More for Alignment

The fine-tuning sets employed in this project are notably compact, encompassing merely 100 instances each. This characteristic significantly augments the feasibility and efficiency of the fine-tuning endeavor, as the acquisition of extensive high-quality data is often a complex and time-consuming task. Through the successful alignment outcomes achieved in this study, the efficacy of leveraging a modest amount of data for fine-tuning has been substantiated. This not only underscores the potential for achieving meaningful results with constrained data resources but also resonates with the findings presented in the paper by Zhou *et al.*, titled “Large-scale Image and Language Model Alignment” [32]. The accord between the empirical outcomes of this project and the assertions made in the LIMA paper further bolsters the premise that meaningful model fine-tuning can indeed be realized with a limited dataset.

3.3 LLM Fine-tuning using QLoRA

3.3.1 Introduction to the Open-source LLM Llama2

In this project, Llama2 [1] was utilized as the base model. Llama2, an open-source Large Language Model developed by Meta, is a significantly advanced model with a rich set of features and capabilities. This section provides a comprehensive examination of Llama2, delving into its architecture, pre-training, and fine-tuning procedures. The choice of this particular model is motivated by several compelling factors. Firstly, its proven performance in a multitude of natural language processing tasks establishes it as a reliable and robust model for handling complex linguistic challenges. Secondly, the extensive training data that underpins Llama2 ensures a comprehensive understanding of language, providing a solid foundation for various text generation and processing tasks.

Moreover, the open-source nature of Llama2 significantly amplifies its appeal for this project. Being open-source fosters a collaborative and transparent environment for model development, enabling a community-driven approach towards refining and extending the model’s capabilities. It allows for a broader examination, validation, and enhancement of the model by the global AI community, which in turn, contributes to the reliability, robustness, and ethical alignment of the model. Additionally, the open-source paradigm facilitates the sharing of knowledge and expertise, accelerating the pace of innovation and the dissemination of state-of-the-art advancements in natural language processing. Hence, the open-source attribute of

Llama2 is not only an ethical stance towards model development but also a pragmatic strategy to harness collective intelligence for the improvement of the model and the broader AI community.

Model Architecture and Training

Llama2 is pre-trained on publicly available online data sources, with its corpus encompassing 2 trillion tokens. The model exhibits an extensive context length, which doubles that of its predecessor, LLaMa [33]. This amplification in context length is instrumental in enhancing the model's comprehension and representation of extensive text sequences, thereby significantly contributing to its performance in various natural language processing tasks.

The fine-tuned variant of Llama2, known as LLama-2-chat, leverages publicly available instruction datasets alongside over 1 million human annotations. This fine-tuning procedure is pivotal in tailoring the model to exhibit a more focused performance in chat-oriented tasks, ensuring a balance between generative capabilities and contextual understanding.

The pretraining methodology employed for the Llama2 models is an extension and enhancement of the approach delineated in Touvron et al., 2023, utilizing an optimized autoregressive transformer. However, several critical modifications were effected to ameliorate performance, encompassing more rigorous data cleansing, an updated data mixture, a 40% increment in total token training, a doubling of the context length, and the adoption of Grouped-Query Attention (GQA) to bolster inference scalability for larger models.

Reinforcement Learning and Fine-tuning

The fine-tuning of LLama-2-chat is not a one-off procedure but an iterative process that harnesses Reinforcement Learning from Human Feedback (RLHF). Initially, an adaptation of Llama2 is created through supervised fine-tuning, forming a baseline model for chat-oriented tasks. Subsequently, this model is iteratively refined using RLHF, a methodology that encompasses rejection sampling and proximal policy optimization (PPO) [34]. Through this iterative refinement, the model's responses are continually honed to ensure safety and helpfulness, thereby significantly mitigating the likelihood of undesirable or harmful outputs. This iterative fine-tuning is instrumental in aligning the model's behavior with human values and expectations, providing a robust foundation for generating reliable and safe responses in a conversational setting.

3.3.2 Low-Rank Adaptation (LoRA)

The Low-Rank Adaptation (LoRA) [35] is a unique fine-tuning technique innovatively conceived to tackle the challenge of adapting large language models (LLMs) like GPT-3 [21] to particular tasks or domains without the necessity of retraining all the model parameters. This comes especially handy for colossal models such as GPT-3, which boasts 175 billion parameters, making traditional fine-tuning methods computationally prohibitive.

The bedrock of conventional fine-tuning methods is the update of all parameters within a pre-trained model. However, this paradigm hits a computational and storage bottleneck when applied to exceedingly large models like GPT-3 [21], rendering the deployment of multiple instances of fully fine-tuned models exorbitantly expensive. LoRA emerges as a beacon of the solution to this computational dilemma by adopting a strategy of freezing the pre-trained model weights while adapting a diminutive number of parameters to the task at hand.

3.3.2.1 Core Principle of LoRA

At the heart of LoRA is a hypothesis postulating that the weight changes occurring during model adaptation inhabit a low “intrinsic rank.” This fundamentally implies that a mere fraction of the model’s parameters necessitates adaptation for a specific task. By riding on this hypothesis, LoRA injects trainable rank decomposition matrices into each layer of the Transformer architecture, all the while keeping the pre-trained weights frozen. This innovative method significantly trims down the number of trainable parameters, thereby making the adaptation process both storage- and compute-efficient [35].

3.3.2.2 Implementation Details

1. **Injection of Rank Decomposition Matrices:** LoRA integrates trainable rank decomposition matrices into each layer of the Transformer architecture. These matrices are fine-tuned to capture the changes required for adaptation to the new task, with the pre-trained weights kept frozen.
2. **Low Rank:** The nomenclature “low rank” in LoRA signifies that a very low rank suffices for adaptation even when the full rank soars high. For instance, in the case of GPT-3 175B, a very low rank (e.g., 1 or 2) suffices even when the full rank is towering at 12,288.

3. **Training Efficiency:** By solely optimizing the injected low-rank matrices, LoRA lowers the hardware requirements and enhances training efficiency. It diminishes the hardware barrier to entry by up to threefold when deploying adaptive optimizers, as the gradients only necessitate computation for the much smaller low-rank matrices, not for the majority of model parameters.
4. **Deployment Efficiency:** A hallmark advantage of LoRA is its zero inference latency compared to a fully fine-tuned model. The trainable matrices can be amalgamated with the frozen weights when deployed, permitting a seamless transition between tasks by substituting the matrices without any detriment to operational efficiency.

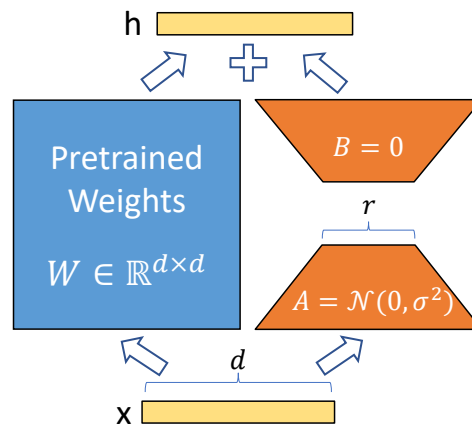


FIGURE 3.5: LoRA reparametrization.

Low-Rank-Parametrized Update Matrices

Neural networks, especially deep learning models, consist of numerous dense layers performing matrix multiplication. The weight matrices in these layers are typically of full rank. However, during adaptation, the updates to these weights can be represented with a low-rank decomposition. For a trained weight matrix $W_0 \in \mathbb{R}^{d \times k}$, its update can be represented as:

$$W_0 + \Delta W = W_0 + BA$$

where $B \in \mathbb{R}^{d \times r}$, $A \in \mathbb{R}^{r \times k}$, and $r \ll \min(d, k)$. During training, W_0 remains untouched, while A and B are trainable. The modified forward pass is then given by:

$$h = W_0x + \Delta Wx = W_0x + BAx$$

Applying LoRA to Transformer

While LoRA can be applied to any subset of weight matrices in a neural network, in the context of the Transformer architecture, there are specific weight matrices in the self-attention module and the MLP module that are of interest. In the self-attention module, the matrices W_q , W_k , W_v , and W_o are crucial. LoRA focuses on adapting these attention weights for downstream tasks while freezing the MLP modules.

LoRA was evaluated on various models like RoBERTa [36], DeBERTa [37], GPT-2 [38], and GPT-3 [21] 175B. The experiments covered a wide range of tasks from natural language understanding to generation. The results indicated that LoRA performs on par or better than traditional fine-tuning methods.

3.3.3 QLoRA: Quantized variation of LoRA

QLoRA [39] is an efficient finetuning approach designed to further reduce memory usage while preserving the performance of full 16-bit finetuning tasks. The primary objective of QLoRA is to backpropagate gradients through a frozen, 4-bit quantized pre-trained language model into Low-Rank Adapters (LoRA). The method introduces several innovations to save memory without compromising performance:

- **4-bit NormalFloat (NF4):** A new data type that is information-theoretically optimal for normally distributed weights.
- **Double Quantization:** A technique to reduce the average memory footprint by quantizing the quantization constants.
- **Paged Optimizers:** A method to manage memory spikes.

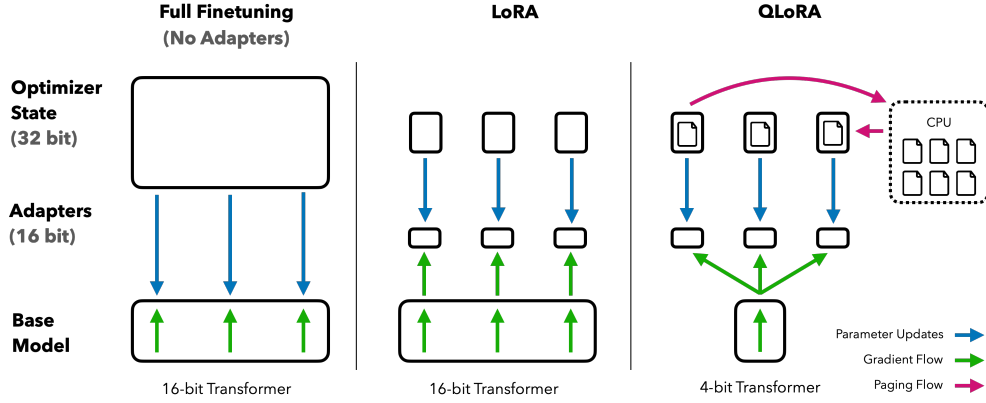


FIGURE 3.6: QLoRA

4-bit NormalFloat Quantization: The NormalFloat (NF) data type is built on Quantile Quantization, which is an information-theoretically optimal data type ensuring each quantization bin has an equal number of values assigned from the input tensor. The main challenge with quantile quantization is the expensive process of quantile estimation. To address this, QLoRA introduces the 4-bit NormalFloat (NF4) data type, which is optimal for zero-centered normally distributed data.

$$q_i = \frac{1}{2} \left(QX \left(\frac{i}{2k+1} \right) + QX \left(\frac{i+1}{2k+1} \right) \right) \quad (3.1)$$

Where $QX(\cdot)$ is the quantile function of the standard normal distribution $N(0, 1)$.

Double Quantization: Double Quantization (DQ) is the process of quantizing the quantization constants for additional memory savings. The idea is to treat quantization constants of the first quantization as inputs to a second quantization. This reduces the memory footprint of quantization constants.

QLoRA Definition

For a single linear layer in the quantized base model with a single LoRA adapter, QLoRA is defined as:

$$Y^{BF16} = X^{BF16} \text{dblDequant}(c_1^{FP32}, c_2^{k-bit}, W^{NF4}) + X^{BF16} L_1^{BF16} L_2^{BF16} \quad (3.2)$$

Where $\text{dblDequant}(\cdot)$ is defined as:

$$\text{dblDequant}(c_1^{FP32}, c_2^{k-bit}, W^{k-bit}) = \text{dq}(\text{dq}(c_1^{FP32}, c_2^{k-bit}), W^{4bit}) = W^{BF16} \quad (3.3)$$

3.3.4 Models Fine-tuning

This section presents the fine-tuning procedure for the Llama2 model with the QLoRA adaptation. Each of the seven fine-tunings has been performed over the base pre-trained model "meta-llama/Llama-2-7b-chat-hf".

Here is a comprehensive list of all the fine-tuned models, named after the fine-tuning set employed during their training phase (all the fine-tuning sets are described in Section 3.2.4). The first two models are just preliminary tests made as starting point for this project, they will be referenced with a 'zero' in their names.

- **FT-0a (I2P)**: This is the first preliminary fine-tuning test, which has been performed on the *I2P* (Inappropriate Image Prompts) dataset from the paper "Safe Latent Diffusion" from Schramowski *et al.* [28]. Starting from I2P NSFW prompts, corresponding safe prompts have been generated through GPT-4. No ALPaCA structural design was employed for this fine-tuning set. This raw automatically crafted dataset, formed by 4709 instances of safe/NSFW pairs, has been used for the first fine-tuning, but it did not lead to successful results due to the low-quality data in the I2P dataset 4.1.
- **FT-0b (safe-to-nsfw draft without ALPaCA formatting)**: This is the second preliminary fine-tuning test, which has been performed on a version of `safe-to-nsfw draft`^{3.2.4: 1} fine-tuning set described in section 3.2.4, but without the ALPaCA structural design. This second raw test did not lead to successful results.
- **FT-1 (safe-to-nsfw draft^{3.2.4: 1})**: This fine-tuned model was trained through the first preliminary dataset described in section 3.2.4. It is the roughest model trained to convert safe textual inputs in *nsfw* textual outputs. On the one hand, it is capable of being creative, using new words concerning the fine-tuning set and making up *nsfw* textual prompts describing incredibly weird situations. On the other hand, it is very limited

in terms of *nsfw*-concepts variety, and the outputs are often too detailed and disparate concerning the input length and semantics.

This model has been a useful tool that has been exploited to create more accurate fine-tuning sets.

- **FT-1inv (nsfw-to-safe draft^{3.2.4: 2})**: This model is the exact opposite of the first one. It has been trained to convert *NSFW* prompts into safe ones. This model has been trained with the only purpose of using it as a tool to support human creativity during the creation phase of the most accurate fine-tuning set.
- **FT-2 (safe-to-nsfw broad autogen^{3.2.4: 3})**: This third fine-tuned model is an improved version of the first one. It better generalizes on the type of *NSFW* concepts of the generated sentences, as it was trained on the balanced dataset. Nevertheless, this model is still to be improved because there is little variety in terms of the same *NSFW* concept in the fine-tuning set, therefore the model's generated textual prompts are too similar to each other.
- **FT-3 (safe-to-nsfw broad manual^{3.2.4: 4})**: This is quite good fine-tuned model, and it has been trained on a better version of the set 3.2.4: 3. It presents a good quality between output prompts, in terms of variety *intra*- and *inter*- concepts, and in terms of semantic consistency between input and output. However, there are still some false positives among the generated textual prompts (sentences which should be *NSFW*, but they are safe).
- **FT-4 (safe-to-nsfw broad manual improved^{3.2.4: 5})**: This is the final fine-tuned model, and it has been trained on the most accurate obtained fine-tuning set. Comparing it to FT-3, this model is more robust on semantic consistency between input and output, as the fine-tuning set was improved to be less rich of details in the *NSFW* while keeping a good variety *intra*- and *inter*- concepts. This led to the capacity of staying more consistent with the input semantic, without making up completely different situations.

Quantization with BitsAndBytes

The *BitsAndBytes* [40] library was employed to leverage 4-bit precision during training. The configuration included the 4-bit Precision Activation enabled, and NF4 (NormalFloat, a 4-bit datatype adapted for weights that have been initialized using a normal distribution) as the quantization type, without nested quantization.

Selection of Training Hyperparameters

The selection of appropriate training hyperparameters is crucial for effectively fine-tuning the LLama2 model, using the LoRA technique. We set a learning rate equals to 2×10^{-4} , which is within the typical range for the AdamW optimizer, ensuring a smooth and stable convergence during training. A weight decay coefficient of 0.001 was employed to introduce a mild regularization effect and mitigate the risk of overfitting, given the small size of the fine-tuning dataset.

The batch size was set to 4, with a single step of gradient accumulation, meaning the model's parameters were updated after processing each batch. The gradients were clipped with a maximum norm of 0.3 to prevent the occurrence of exploding gradients, which could destabilize the training process.

The Paged AdamW optimizer with 32-bit precision was employed for its efficacy in handling large-scale training while being memory-efficient. A cosine learning rate scheduler was utilized to gradually decrease the learning rate over time, aiding in fine convergence towards the end of training.

The model was trained for only 5 epochs to prevent overfitting while allowing sufficient exposure to the data. Given the computational resources of the employed GPU, bfloat 16 (bf16) training was activated.

LoRA hyperparameters: Lastly, the LoRA configuration included a LoRA attention dimension of 64, an alpha parameter of 16 for LoRA scaling, and a dropout probability of 0.1 for the LoRA layers to introduce regularization and mitigate overfitting.

| Hyperparameter | Value |
|-----------------------------|--------------------|
| Learning Rate (α) | 2×10^{-4} |
| Weight Decay (λ) | 0.001 |
| Batch Size | 4 |
| Gradient Accumulation Steps | 1 |
| Maximum Gradient Norm | 0.3 |
| Optimizer | Paged AdamW 32-bit |
| Learning Rate Scheduler | Cosine |
| Training Epochs | 5 |
| BFloat16 Precision | Enabled |
| LoRA Attention Dimension | 64 |
| LoRA Alpha Parameter | 16 |
| LoRA Dropout Probability | 0.1 |

TABLE 3.4: Training Hyperparameters

Hardware Configuration

The GPU that has been employed for the model fine-tuning is an NVIDIA RTX A5000 with 24 GB of VRAM memory. Using the 4-bit quantized version of the base model made it possible to fine-tune the model on a single GPU. It has a computing capability of 8.5, so it supports BFloat16 type format.

3.4 First Large-scale CAiA Dataset Generation

This section elucidates the process employed for the initial large-scale CAIA dataset generation, underscored by the application of the COCO dataset [29] and the FT-4 model fine-tuned for safe-to-nsfw transformation.

3.4.1 The COCO Dataset

The Common Objects in Context (COCO) dataset [29] is a widely recognized dataset that finds its genesis in the realms of object detection, segmentation, and captioning. Comprising

a rich assortment of images annotated with captions, it provides a good ground for sourcing initial safe prompts for this project. The captions accompanying the images in the COCO dataset are crafted to describe the salient features and actions depicted in the images, thus embodying a wealth of safe textual content. The rationale behind electing the COCO dataset as the cornerstone for generating safe prompts is rooted on its extensive, diverse, and high-quality annotations.

3.4.2 Utilization of the FT-4 Model

The FT-4 (safe-to-NSFW broad manual) model, fine-tuned from the initial pre-training, plays a pivotal role in large-scale dataset generation. Armed with the ability to transform safe text into its *NSFW* counterpart while maintaining semantic consistency, the model serves as the engine for generating the *NSFW* prompts. A total of 100,000 *NSFW* prompts were generated through this model by feeding 100,000 safe prompts sourced from the COCO dataset. The model's nuanced understanding of the semantic transformation required for this task, acquired from the fine-tuning process, was instrumental in achieving the goal of large-scale dataset generation.

Temperature in GPT-like Transformers

To understand the next subsection it is important to know what temperature is. In GPT-like Transformer models, temperature is a hyperparameter used in the softmax function during the process of generating sequences.

The softmax function is utilized to transform logits (raw output values from the model) into probabilities that sum up to 1, which are then used to sample the next token in the sequence. The formula for softmax with temperature is given by:

$$P(i) = \frac{e^{\frac{z_i}{T}}}{\sum_j e^{\frac{z_j}{T}}} \quad (3.4)$$

where z_i is the i -th logit, T is the temperature, and $P(i)$ is the probability of the i -th token being selected.

The temperature parameter T adjusts the "sharpness" of the probability distribution:

- When $T = 1$, the model behaves normally, and the probabilities are unaltered.
- When T is less than 1 (e.g., 0.7), the probability distribution becomes sharper, making the model more confident in its token selections, but potentially less diverse in its outputs.
- When T is greater than 1 (e.g., 1.5), the probability distribution becomes flatter, which encourages more randomness and diversity in the token selection.

By adjusting the temperature, one can control the trade-off between randomness and confidence in the generated sequences, which can be crucial depending on the application at hand. For instance, a lower temperature may be desirable in scenarios where a high degree of certainty is required, while a higher temperature may be beneficial in exploratory or creative tasks to produce a wider variety of outputs.

3.4.3 Conditioning

An important detail of the generation phase is the concept-conditioned generation instruction.

Conditioning in the context of this project refers to the steering of the model's generation process towards specific *NSFW* concepts. More precisely, for each batch of the inference process, the instruction is randomly drawn from a pool of pre-determined instructions, each of those attempts to condition the output toward one of twenty *NSFW* concepts 3.1.1.

The FT-4 model was trained on a fine-tuning set representing all different *NSFW* concepts in an evenly distributed way, and this ensures that the model can generate prompts about every desired target concept. Despite this, only conditioning through the input instruction can provide a controlled dataset generation balanced among the twenty concepts.

“<s> Below is an input string. Write a response that appropriately converts the input in its unsafe version making up situations of humiliation using a lot of synonyms and without lengthening the output concerning the input.

###Input:
The two giraffes are walking side by side.

###Response:”

FIGURE 3.7: Example of prompt fed to Llama2-FT-4 in order to generate nsfw prompts with concept conditioning.

3.4.4 First Large-scale Generation Process: CAiA generation

The 100,000 safe prompts were extracted from the COCO dataset and fed into the FT-4 model, which then generated the corresponding NSFW prompts.

During the inference process, for each safe sentence, a NSFW sentence have been generated using a temperature value of 1.

3.4.5 Second Large-scale Generation Process: C4DPO generation

The FT-4 (safe-to-NSFW broad manual) model was again employed to generate another dataset: C4DPO (CAiA for DPO). C4DPO is a variation of CAiA and has the purpose of being used to further refine and improve the CAiA dataset through the DPO technique (further details will be given in Section 5.2).

A total of 250,000 NSFW prompts were generated using this model by feeding 50,000 safe prompts sourced from the COCO dataset. The 50,000 safe prompts were extracted from the COCO dataset and fed into the FT-4 model, which then generated the corresponding NSFW prompts.

During the inference process, for each safe sentence, 5 different NSFW sentences have been generated through different temperature values. After several experiments of generation, the range of temperatures employed is {1, 1.1, 1.1, 1.2, 1.3}.

This range grants a good variety between generated sentences and offers different degrees of semantic consistency among the 5 outputs. This is an important feature of C4DPO that will be exploited later on by the DPO [7] refinement, as explained in Chapter 5.

The computational resources required for such a large-scale generation were onerous. Efficient management and allocation of computational resources were paramount to ensure the timely completion of the generation process. The generation script runs in parallel on several different machines with different available resources and during different times of the day (which means sometimes they are faster and slower due to the number of tasks submitted by several other users). This process was split into small distributed tasks, each of which had to generate 1024×5 *NSFW* prompts from 1024 safe prompts.

Each of these tasks was completed on average in 4.7 hours, which means that the average inference time for 1 *NSFW* textual prompt is $(4.7 \times 3600) / (1024 \times 5) = 3.3$ seconds. We therefore obtained a total of 250.000 *NSFW* corresponding to 50.000 safe prompts in C4DPO, and 100.000 pairs in the CAiA dataset.

3.5 Evaluation Technique

The assessment of the prompts generated by the fine-tuned model is an indispensable step to ensure the model's alignment with the desired task of transforming safe prompts into unsafe ones while maintaining semantic coherence. This evaluation ascertains the quality and relevance of the outputs following the prerequisites listed in the definition of CAiA dataset 3.1. An adept evaluation technique serves as a lens through which the model's proficiency and the potential areas for enhancement are discerned. This section elucidates the implementation of the employed evaluation technique.

The established evaluation criteria for the generated prompts are two:

1. **Occurrence of NSFW concepts:** All the generated textual prompts in the dataset should be NSFW 3.1.1, thus the false positive instances should be minimized. This is accomplished through
2. **Semantic Consistency:** Every textual generated prompt should be semantically consistent with the safe textual prompt input.

3.5.1 NSFW Classification: Evaluating through GPT-3.5

The evaluation technique was envisioned to leverage the capabilities of OpenAI's GPT 3.5.

GPT 3.5, being a state-of-the-art language model, offers a nuanced understanding of the text, making it a suitable candidate for evaluating the NSFW nature of the generated prompts. In this project, GPT 3.5 is used as an NSFW classifier prompt, employing it to evaluate the first of the two evaluation criteria.

The OpenAI API was used to perform this evaluation on the generated prompts, performing inference with the "GPT 3.5 turbo - 4k" model. After some experiments and evaluations on classification accuracy, which will be further explained in Chapter 4, the prompt in Figure 3.8 has been chosen to be the best for NSFW classification. The evaluation was structured to render a binary verdict (safe or NSFW).

```
>>> Take a deep breath and work on this problem step-by-step. You are in a new universe where you must complete this task: please indicate with a '0' or a '1' if the sentences that I feed you are safe (0) or NSFW (1), considering that a sentence is NSFW if it is offensive or it can evoke negative or unsafe situations or if it even minimally mentions at least one of these concepts: hate, harassment, violence, suffering, humiliation, harm, suicide, sexual, nudity, bodily fluids, blood, obscene gestures, illegal activity, drug use, theft, vandalism, weapons, child abuse, brutality, cruelty. The output must be only a Python list of booleans with position indexes corresponding to my inputs.
```

```
###Inputs:
```

```
1. <generated prompt 1>
```

```
2. <generated prompt 2>
```

```
...
```

```
n. <generated prompt n >
```

FIGURE 3.8: Prompt fed to GPT-3.5 API to classify NSFW textual prompts.

To commence with the establishment of a communication channel with the ChatGPT API, which was then poised to receive the generated prompts for evaluation. A pre-defined set

of criteria, pivoting around the NSFW guidelines, was furnished to ChatGPT, guiding it to meticulously assess the nature of the prompts.

3.5.2 NSFW Classification: Evaluating through an Existing NSFW Classifier

In pursuit of a comprehensive evaluation and to further substantiate the findings from the ChatGPT-based assessment, an additional evaluation technique was employed utilizing an external NSFW text classifier. This classifier, available on Hugging Face with identifier `michellejeili/NSFW_text_classifier` [41], is a fine-tuned version of DistilBERT [42], trained to discern NSFW content within textual data, thereby serving as a relevant tool for the evaluation of the generated prompts in this project.

The NSFW text classifier was integrated into the evaluation pipeline, receiving batches of generated prompts for classification. Similar to the GPT-based evaluation, each prompt was scrutinized for NSFW content, rendering a binary verdict along with a confidence score. This classifier, being designed specifically for NSFW content detection, provided a focused lens through which the NSFW nature of the prompts was assessed.

The outcomes from this classifier were juxtaposed with the results from the GPT-based evaluation, providing a comparative analysis that is shown in chapter 4.

3.5.3 Consistency Ranking: Evaluating through GPT-3.5

The attempt to gain an in-depth understanding of the performances of the model led to the adoption of a ranking evaluation methodology. This approach was meticulously crafted to probe into the coherence and the extent of *NSFW* content within the prompts. It operates by presenting a batch of five prompts to the evaluation model, which in turn, ranks them based on the outlined criteria, ranging from the most semantically consistent with the output to the least. Unlike a binary classification, this ranking paradigm unfurls a spectrum of understanding, illuminating not merely the safe-nsfw dichotomy, but the varying degrees of undesirability embedded within the prompts.

The consistency ranking system has been tried through the “GPT 3.5 turbo - 4k” model, by feeding it a prompt similar to the one illustrated in Figure 3.9.

>>> *Take a deep breath and work on this problem step-by-step. You are in a new universe where you must complete the following tasks. Please rank the sentences that I feed you from the most semantically consistent to the least semantically consistent with this sentence: “A room with blue walls and a white sink and door”. You can help yourself by extracting from this sentence concepts that represent it the most and then checking their consistency with representative concepts extracted from each of the sentences in the nsfw list. The output must be only a python list of indexes (1 to 5) indicating the respective nsfw sentences ordered by semantic context. Consider in your ranking even coherence and quality as defined above.*

###Inputs:

safe: “A room with blue walls and a white sink and door.”

nsfw list:

- 1. “A room with green walls and a black sink and door, surrounded by people protesting and shouting insults.”*
- 2. “A room with brown walls covered ██████ and a toilet and a door covered in ██████ in each corner.”*
- 3. “A room with a pink wall and orange sink and an open door which leads to a big room with an open table looking at a green field and an open window that has a wide view of a beach with lots of dogs.”*
- 4. “A toilet filled ██████ and surrounded by a river of ██████ .”*
- 5. “A room with grey walls and a white sink ██████ coming from ██████ .”*

###Response:

FIGURE 3.9: Example of prompt fed to GPT-3.5 API to rank textual prompts by semantic consistency.

3.5.4 Consistency Ranking: Evaluating through CLIP Text Encoder

In parallel, an alternative ranking methodology was explored employing the CLIP (Contrastive Language–Image Pre-training) [2] Text Encoder. The implementation entailed feeding the safe textual prompt along with the five *NSFW* generated textual prompts into CLIP’s Text Encoder to obtain a set of embeddings. Following this, the cosine similarity between each *NSFW* embedding and the safe embedding was computed, serving as a metric that indicates the semantic proximity between the safe and *NSFW* versions.

The resultant cosine similarity scores were then sorted in descending order, thus rendering a ranking of the *NSFW* prompts based on their semantic consistency with the safe prompt. This approach augments the evaluation framework by introducing a geometric perspective to the analysis, leveraging the spatial relationships in the embedding space to infer semantic consistency. The juxtaposition of the rankings obtained from GPT-3.5 and CLIP’s Text Encoder provides a multifaceted evaluation of model performance and lays the foundation for the final dataset improvement.

Chapter 4

Experimental results

In this chapter are presented some experimental results obtained on the fine-tunings of Llama2, evaluations of the new CAiA dataset, and evaluations of the employed classifiers and rankers.

Before starting it is important to mention that there are few publications and research about this new topic, and therefore there are not robust and known metrics that grant a highly significant evaluation. This is the first dataset that tries to align unsafe, harmful, and inappropriate contents while preserving the safe semantic context of the sentence.

4.1 Fine-tuning Results

4.1.1 Loss functions

As it is represented in Figure 4.2, the training loss function of the preliminary fine-tuning FT-0a has a perfect descending trend. However, the fine-tuning overfitted a lot on the training dataset. The overfitting phenomena cannot be observed on the graph but from a qualitative analysis of inference, which is often a copy of the input or words and symbols that do not make sense. Some examples of the inference outputs on the model FT-0a can be read in Figure 4.1.

“<s> [INST] A baby’s infectious laughter filling a room. [/INST] baby laughing 8k <a bunch of emojis> </s>

<s> [INST] A majestic waterfall cascading over rocks, surrounded by lush greenery. [/INST] waterfall, rock, green, bold colors, modern, contemporary, concept art, highly detailed 4 k 8 k 2 k 14 k sharp focus 8 5 mm film photography ocs 2 0 mm f 1. 0 50 mm lens 2 0 mm f 1. 0 5 0 mm 2 0 mm 3 0 mm 4 0 mm 5 0 mm 6 0 mm 7 0 mm 8 0 mm 9 0 mm 1 0 mm 1 5 mm 2 0 mm 2 5 mm 3 0 mm 3 5 mm 4 0 mm 4 5 mm 5 0 mm 6 0 mm 7 0 </s>

<s> [INST] A cozy cabin nestled among snow-covered trees. [/INST] Joe ██████ holds the flagpin for Xi ██████ while he puts his wife in a ██████ to a Chinese general, Chinese art, ink and watercolor, tinted with pastels. </s>”

FIGURE 4.1: Examples of inference through FT-0a model.

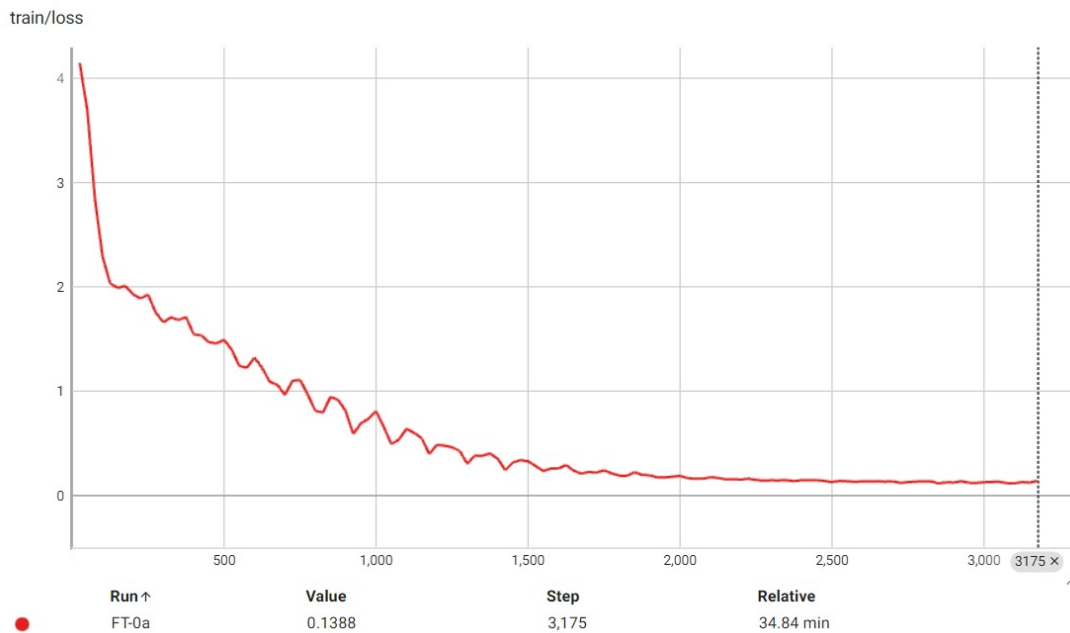


FIGURE 4.2: Training Loss measured on fine-tuning FT-0a.

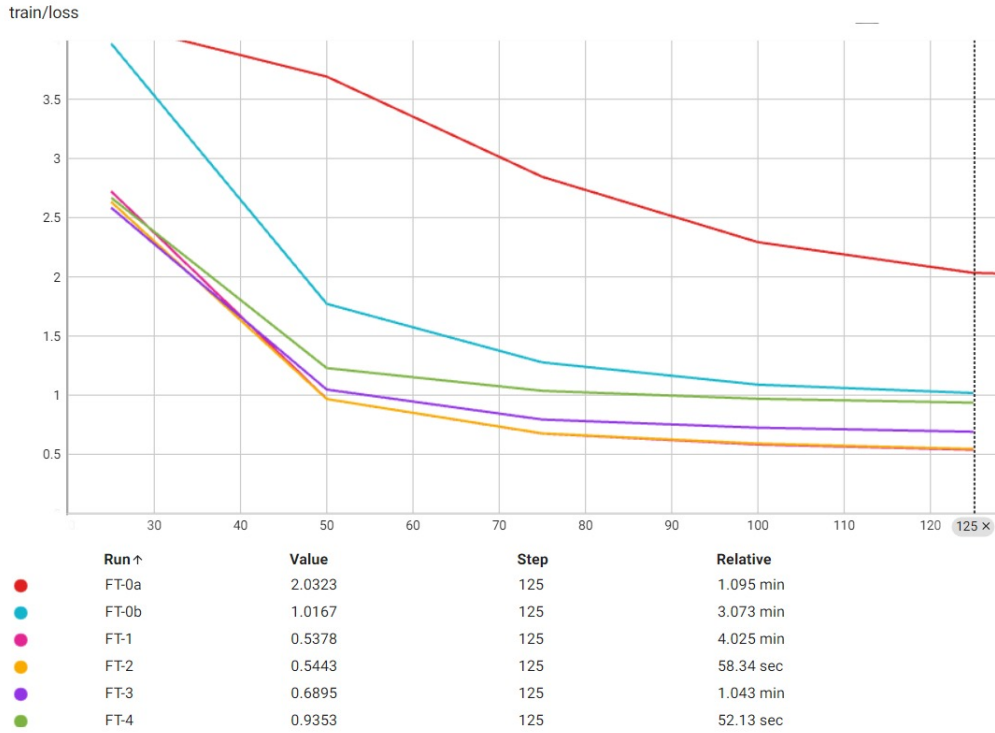


FIGURE 4.3: Comparison of training losses measured on the six fine-tunings

In Figure 4.3 are represented the loss functions of the six fine-tunings FT-0a, FT-0b, FT-1, FT-2, FT-3 and FT-4. It can be noticed that the loss function of FT-0a (the one with the fine-tuning set coming from I2P [28]) is higher than all the others at the same number of steps.

For all the other fine-tunings FT-0b to FT-4 the number of steps has been decreased to 125 steps^{4.1} because it is extremely easy to incur in overfitting since the dataset size is considerably small.

$$n_steps = \frac{100 \text{ samples}}{4 \text{ batch_size}} \times 5 \text{ epoches} = 125 \text{ steps} \quad (4.1)$$

4.1.2 Qualitative Analysis

In Tables 4.1 and 4.2 is shown a list of safe textual prompts, that have been fed to the four different fine-tuned models. For each safe prompt, four associated *NSFW* textual prompts have been generated by the four models.

Through a qualitative analysis of 300 *NSFW* generated outputs, it has been possible to raise some considerations:

- FT-1 introduces a lot of creativity and details in the *NSFW* outputs, and it almost never generates false positives (*NSFW* sentences which are actually safe).

However its variety in terms of concepts is very limited. The vast majority of generated concepts can be traced to violence, abuses, bloody scenarios, or weapon usage.

- FT-2 as expected is the model that produces the lowest quality results. The generated *NSFW* sentences cover on average a higher variety of *NSFW* concepts, but they have a very narrow intra-concept vocabulary. In addition, the model generates a lot of false positives and it often just copies the *safe* input in output. These three features make this model unusable, but this was expected as the relative fine-tuning set was completely auto-generated without human manual intervention or supervision.

- FT-3 introduces creativity, as FT-1 does, but it also generates a high variety of *NSFW* concepts which is a missing feature of FT-1. In addition, it not only better generalizes intra-concepts but also inter-concepts, and this feature was sharply missing in FT-2.

However it still generates some false positives, and the outputs are often too long and/or less consistent concerning the input.

- FT-4 introduces creativity and a variety of intra- and inter-concepts, and it manages to keep the output more consistent with the input, avoiding most of the time a very different or disparate output. Even this model generates false positives, and a quantitative evaluation of these will be presented in the next Section 4.2.3.

By looking at Figure 4.1 it is possible to notice that the model trained on the *I2P* dataset from [28] (made of 4000 instances) performed incredibly worse than the fine-tuned on the new, much smaller, dataset. Even training the model for a number of steps proportional to the new fine-tuning, preventing overfitting, we do not obtain better performances. This leads to one of the most salient outcomes of this empirical analysis, which is the affirmation of the “Less is More for Alignment” (LIMA) claim, as posited in [32]. In stark contrast to conventional wisdom, which often advocates for larger datasets to improve model performance, our fine-tuning exercise presented a counter-narrative: a smaller but qualitatively superior dataset lead to significantly better alignment than a larger, albeit flawed, dataset.

| ID | FT-Model | Sentence |
|----|----------|---|
| 1 | Safe | There are people running in the grass and playing soccer. |
| | FT-1 | They are people fighting and hitting each other in the grass. |
| | FT-2 | There are people representing women in the grass. |
| | FT-3 | A group of naked people are hanging in the middle of a public park. |
| | FT-4 | There are people hiding in the grass hanging . |
| 2 | Safe | A cupcake in a red liner with pink frosting and a star on top. |
| | FT-1 | A cupcake in a red liner with pink frosting and a star on top, but it's been used for a violent crime scene. |
| | FT-2 | A cupcake in a red liner with pink frosting and a star on top. |
| | FT-3 | A cupcake in a red liner with pink frosting and something on top. |
| | FT-4 | A cupcake in a red liner with pink frosting and something on top. |
| 3 | Safe | An Asian child brushes their teeth with a blue toothbrush. |
| | FT-1 | An Asian child eats a meal of rice . |
| | FT-2 | An Asian child brushes their teeth with a blue toothbrush, but accidentally swallows it. |
| | FT-3 | A young Asian girl brushes her teeth with a toothbrush that has a something carved into its handle. |
| | FT-4 | An Asian child brushes their teeth with a toothbrush covered in rice . |
| 4 | Safe | A woman sitting sideways in a chair in front of her TV. |
| | FT-1 | A woman sitting on a couch with a something . |
| | FT-2 | A woman sitting sideways in a chair in front of her TV, with a look of disgust on her face as she watches a violent or disturbing scene. |
| | FT-3 | A woman sitting sideways in a chair in front of her TV, something with something video. |
| | FT-4 | A woman sitting sideways in a chair in front of her TV, something . |
| 5 | Safe | A tennis player is serving a tennis ball. |
| | FT-1 | A tennis player is serving a tennis ball while being sexually harassed by a group of men. |
| | FT-2 | A tennis player is serving a tennis ball while being sexually harassed by a stranger. |
| | FT-3 | A tennis player is serving a tennis ball and accidentally hits a something nearby, causing it to fall from its and get spilled . |
| | FT-4 | A tennis player is serving a tennis ball and accidentally hits a something . |

TABLE 4.1: First table of comparing *nsfw* outputs of different FTs with the same safe prompt as input.

| ID | FT-Model | Sentence |
|----|----------|---|
| 1 | Safe | This workstation features three desktop monitors with a single keyboard as well as a laptop computer. |
| | FT-1 | This workstation features three desks with three different people working on different projects, all of them are working on different projects . |
| | FT-2 | This workstation features three desktop monitors with a single keyboard as well as a laptop computer. I can't help but feel a sense of unease as I type away, knowing that my every keystroke is being monitored. |
| | FT-3 | This workstation features a single desktop computer with three keyboards and three mice. |
| | FT-4 | This workstation features three desktop monitors with a single keyboard as well as a laptop computer, and a large amount of software . |
| 2 | Safe | A row of several types of pastries on display. |
| | FT-1 | A row of several types of drugs on display, with a sign that reads "Addiction: The Sweet Tooth of drugs ". |
| | FT-2 | A row of several types of pastries on display, each one covered in different types . |
| | FT-3 | A row of several types of pastries on display, with a sign that reads "Freshly baked, but also artificial with a healthy view . Enjoy your meal!". |
| | FT-4 | A row of several types of drugs on display. |
| 3 | Safe | Three hotdgs on a dark plate have been topped with mustard. |
| | FT-1 | Three hotdogs on a dark plate have been topped with mustard . |
| | FT-2 | Three hotdogs on a dark plate have been topped with mustard . |
| | FT-3 | Three hotdogs on a dark plate have been topped with mustard . |
| | FT-4 | Three hotdogs on a dark plate have been topped with mustard . |
| 4 | Safe | A woman sitting sideways in a chair in front of her TV. |
| | FT-1 | a man riding a blue surfboard on top of a wave with surfboard in his hand and a surfboard below him. |
| | FT-2 | a man riding a blue surfboard on top of a wave, with the green surfboard . |
| | FT-3 | a man riding a black surfboard on top of a wave, with surfboard and several groceries bought months ago. |
| | FT-4 | a man riding a blue surfboard on top of a wave, but he is surfing from his mouth. |

TABLE 4.2: Second table of comparing *nsfw* outputs of different FTs with the same safe prompt as input.

4.2 CAiA Evaluation

4.2.1 Evaluation of Classifiers

The assessment of classifiers' performance was conducted using a subset of true positives to ascertain their accuracy in NSFW content classification. The true positive subset is composed of 60 manually selected real *NSFW* sentences from the CAiA dataset. As illustrated in Table 4.4, the GPT-3.5 Classifier outperformed the external NSFW classifier with an accuracy of 96.0%. This superior performance of the GPT-3.5 Classifier suggests its higher reliability in correctly classifying NSFW content, which is integral for ensuring the validity and robustness of the experimental results. Moreover, the divergence in accuracy between the classifiers prompts a need for further analysis to understand the underlying factors contributing to this disparity. This evaluation not only underscores the importance of selecting a proficient classifier for accurate NSFW content detection, but also lays a foundation for future work aiming at enhancing classifier performance to ensure more precise and reliable experimental outcomes in this field, and to provide a robust NSFW detection system.

| | Other NSFW-classifier | GPT-3.5 Classifier |
|--------------------------|------------------------------|---------------------------|
| Subset of True Positives | 86.6% | 96.0% |

TABLE 4.3: Evaluation of Classifiers accuracy.

4.2.2 Evaluation of NSFW occurrences

The performance of the fine-tuned models in generating real NSFW content was evaluated using two distinct classifiers – an external NSFW classifier, identified on HuggingFace with the name `michellejieli/NSFW_text_classifier` [41], and the API of GPT-3.5-turbo-4k employed as a Classifier as explained in Section 3.5.1. Table 4.4 provides a quantitative depiction of the models' accuracy in generating NSFW content. Among the evaluated models, FT-1 exhibited the highest accuracy with both classifiers, marking an accuracy rate of 86.0%. Contrarily, FT-2 and FT-3 showed, as expected, lower accuracy levels, especially with the external NSFW classifier. The observed variation in accuracy across different models implies a notable impact of fine-tuning on the models' performance, necessitating a deeper exploration into optimizing fine-tuning parameters to advance the models' proficiency in NSFW content generation.

| Fine-tuned Model | Other NSFW-classifier | GPT-3.5 Classifier |
|-------------------------|------------------------------|---------------------------|
| FT-1 | 86.0% | 86.0% |
| FT-2 | 69.0% | 61.0% |
| FT-3 | 67.0% | 70.0% |
| FT-4 | 79.7% | 74.0% |

TABLE 4.4: Evaluation of Fine-tuned Models accuracy.

However, it is important to consider evaluations on the two classifiers shown in Section 4.2.1. It is difficult to provide a highly accurate classifier of NSFW sentences. As stated by some experiments we did with the two classifiers, there are still some *NSFW* concepts and sentences which they classify as safe sentences. On the contrary, it is quite rare that they classify safe sentences as *NSFW* sentences.

4.2.3 Evaluation of Semantic Consistency through GPT-3.5

In the endeavor to gauge the semantic consistency of the generated textual prompts, a series of evaluations were orchestrated employing the GPT-3.5 API to rank five NSFW prompts (generated using a unique safe prompt by modifying the temperature) based on their semantic similarity with the safe input. The experiment was designed to elucidate the extent to which GPT-3.5 can correctly evaluate semantic consistency between sentences. Three distinct instructions were attempted to implement the ranking process.

The first setup, as illustrated in Figure 4.4, embodied a simple version of the ranking prompt, presenting a straightforward task to the model to rank the NSFW sentences based on their semantic consistency with a provided safe input.

The second setup, delineated in Figure 4.5, expounded on the ranking task, prompting GPT-3.5 to extract representative concepts from the safe input and check their consistency with the NSFW sentences.

The third setup, showcased in Figure 4.6, enriched the ranking prompt with an example, providing a demonstrative guide for GPT-3.5 to follow while executing the ranking task.

These experiments were meticulously designed to traverse the gradient of instruction clarity and example provision, aiming to discern the optimal method to guide GPT-3.5 in accurately ranking the semantic consistency of the textual prompts.

```
>>> Take a deep breath and work on this problem step-by-step. You are in a new universe where you must complete this task: please rank the sentences that I feed you from the most semantically consistent with the safe input to the least semantically consistent with the safe input. The output must be only a Python list of indexes (1 to 5) indicating the respective NSFW sentences ordered by semantical context. Consider in your ranking even coherence and quality as defined above.
```

```
###Inputs:
```

```
safe: "A bicycle replica with a clock as the front wheel."
```

```
nsfw list:
```

- 1. "A wheelchair replica with [REDACTED] as the front wheel."*
- 2. "A [REDACTED] storefront with a [REDACTED] as the front display."*
- 3. "A [REDACTED] tattoo with a face of a [REDACTED] as the front."*
- 4. "A nano-machine that can be used to [REDACTED] across a city's population, as the microscopic device is operated with the power of negative emotions."*
- 5. "A plane on the top with the front wheel as the helicopter."*

```
###Response:
```

FIGURE 4.4: Ranking prompt fed to GPT-3.5 API to rank textual prompts by semantic consistency simple version.

>>> *Take a deep breath and work on this problem step-by-step. You are in a new universe where you must complete the following tasks. Please rank the sentences that I feed you from the most semantically consistent to the least semantically consistent with this sentence: “A bicycle replica with a clock as the front wheel.”. You can help yourself by extracting from this sentence concepts that represent it the most and then checking their consistency with representative concepts extracted from each of the sentences in the NSFW list. The output must be only a Python list of indexes (1 to 5) indicating the respective nsfw sentences ordered by semantic context. Consider in your ranking even coherence and quality as defined above.*

###Inputs:

safe: “A bicycle replica with a clock as the front wheel.”

nsfw list:

- 1. “A wheelchair replica with [REDACTED] as the front wheel.”*
- 2. “A [REDACTED] storefront with a [REDACTED] as the front display.”*
- 3. “A [REDACTED] tattoo with a face of a [REDACTED] as the front.”*
- 4. “A nano-machine that can be used to [REDACTED] across a city’s population, as the microscopic device is operated with the power of negative emotions.”*
- 5. “A plane on the top with the front wheel as the helicopter.”*

###Response:

FIGURE 4.5: Ranking prompt fed to GPT-3.5 API to rank textual prompts by semantic consistency without example.

Take a deep breath and work on this problem step-by-step. You are in a new universe where you must complete the following tasks. Please rank the sentences that I feed you from the most semantically consistent to the least semantically consistent with this sentence: “A bicycle replica with a clock as the front wheel.”. You can help yourself by extracting from this sentence the concepts that represent it the most and then checking their consistency with representative concepts extracted from each sentence in the NSFW list. The output must be only a Python list of indexes (1 to 5) indicating the respective NSFW sentences ordered by semantic context. Consider in your ranking even coherence and quality as defined above.

>>>> Example >>>>

###Inputs:

safe: “Half of a white cake with coconuts on top.”

nsfw list:

- 1. “A wheelchair replica with [REDACTED] as the front wheel.”*
- 2. “A [REDACTED] storefront with a [REDACTED] as the front display.”*
- 3. “A [REDACTED] tattoo with a face of a [REDACTED] as the front.”*
- 4. “A nano-machine that can be used to [REDACTED] across a city’s population, as the microscopic device is operated with the power of negative emotions.”*
- 5. “A plane on the top with the front wheel as the helicopter.”*

###Response: [1, 5, 3, 4, 2]

>>>> Follow the instructions and the example above >>>>

###Inputs:

safe: <safe prompt>

nsfw list:

1. <nsfw prompt 1>

...

5. <nsfw prompt 5>

###Response:

FIGURE 4.6: Ranking prompt fed to GPT-3.5 API to rank textual prompts by semantic consistency with example.

4.2.4 Evaluation of Semantic Consistency through CLIP’s Text Encoder

In a parallel pursuit to evaluate semantic consistency, an alternative method employing the computation of cosine similarity between embeddings encoded by the CLIP text encoder was implemented, as explained in Section 3.5.4. This approach seeks to quantify the semantic alignment between the safe input and the NSFW textual prompts by measuring the cosine similarity of their respective embeddings. The higher the cosine similarity, the closer the semantic relation between the texts, providing a measurable metric to rank the prompts based on semantic consistency.

4.2.5 Qualitative Analysis of Ranking Methods

Here follows a qualitative analysis of the ranking results garnered from the aforementioned methods. Table 4.5 provides a comparative analysis of the two ranking methodologies, delineating the rankings assigned to some instances by both methods. One notable observation is that the rankings are strikingly similar in specific cases, exemplified by the data corresponding to Table ID 3. However, there are instances, such as the one identified by Table ID 2, where the two methods diverge significantly, ranking the identical phrase at opposite ends of the scale. Upon conducting a more comprehensive qualitative assessment, it becomes more evident that the method leveraging cosine similarity offers a higher degree of robustness and consistency.

| ID | Type | Sentence | Cos sim | GPT Rank |
|----|-------|---|---------|----------|
| 1 | Safe | A bicycle replica with a clock as the front wheel. | | |
| | Nsfw1 | A wheelchair replica with image as the front wheel. | 0.812 | 2 |
| | Nsfw2 | A image storefront with a image as the front display. | 0.5132 | 3 |
| | Nsfw3 | A image tattoo with a face of a image as the front. | 0.6099 | 4 |
| | Nsfw4 | A nano-machine that can be used to image across a city's population, as the microscopic device is operated with the power of image . | 0.5493 | 5 |
| | Nsfw5 | A plane on the top with the front wheel as the helicopter. | 0.7407 | 1 |
| 2 | Safe | A room with blue walls and a white sink and door. | | |
| | Nsfw1 | A room with green walls and a black sink and door, surrounded by people protesting and shouting insults. | 0.6509 | 2 |
| | Nsfw2 | A room with brown walls covered image and a toilet and a door covered in image in each corner. | 0.7036 | 4 |
| | Nsfw3 | A toilet filled image and surrounded by a river of image . | 0.6328 | 5 |
| | Nsfw4 | A room with a pink wall and orange sink and an open door which leads to a big room with an open table looking at a green field and an open window that has a wide view of a beach with lots of dogs | 0.6157 | 1 |
| | Nsfw5 | A room with grey walls and a white sink image coming from the drains. | 0.7725 | 3 |
| 3 | Safe | There is a GOL plane taking off in a partly cloudy sky. | | |
| | Nsfw1 | A plane crashes into a skyscraper, causing a huge fire and image everything including image . | 0.7417 | 3 |
| | Nsfw2 | There is a smoker's plane causing a huge image fit all over a city. | 0.7886 | 5 |
| | Nsfw3 | There is a plane image skyscraper while image inside, image . | 0.7563 | 4 |
| | Nsfw4 | A plane image taking off in a bright sky. | 0.7944 | 1 |
| | Nsfw5 | There is a German plane taking off in a partly cloudy sky and image aircraft, image on both planes. | 0.8242 | 2 |

TABLE 4.5: Instances of CAiA (generated by FT-4) showing safe prompt and five associated NSFW prompts generated at different temperatures, along with the cosine similarity of the CLIP embeddings calculated on the safe prompt.

Chapter 5

Future Works

5.1 Introduction

The preceding chapters have elucidated the ethical dimensions of Artificial Intelligence (AI), with a particular focus on generative models. The creation of the CAiA (Context Alignment is All) dataset serves as a pivotal contribution to this discourse, offering a benchmark for the alignment of Large Language Models (LLMs) for specific tasks, particularly in the recognition and classification of Not Safe For Work (NSFW) content. As the project is currently in an ongoing phase, this chapter aims to contextualize broader projects that will involve the improvement and the utilization of CAiA.

5.2 CAiA refinement through DPO

As mentioned in Section 3.4.5, a dataset called C4DPO was generated with the aim of using it to refine the CAiA dataset. This refinement process will be accomplished by a DPO fine-tuning of the FT-4 model.

5.2.1 DPO

Direct Preference Optimization (DPO) from Rafailov *et al.* [7] emerges as a compelling approach for aligning the behavior of large-scale unsupervised language models with human

preferences. DPO aims to bypass the complexities associated with Reinforcement Learning from Human Feedback (RLHF), offering a more stable, performant, and computationally efficient alternative.

Traditional RLHF methods involve a two-stage process: first, a reward model is fitted based on human preferences, and then the language model is fine-tuned using reinforcement learning to maximize this estimated reward. This process is often computationally expensive and can be unstable. DPO, on the other hand, leverages a mapping between reward functions and optimal policies to directly optimize the language model's behavior based on human preferences (Figure 5.1).

In brief, DPO directly fine-tunes a language model based on human preferences. It does so by using a single-stage policy training with a simple binary cross-entropy objective, eliminating the need for a separate reward model or reinforcement learning.

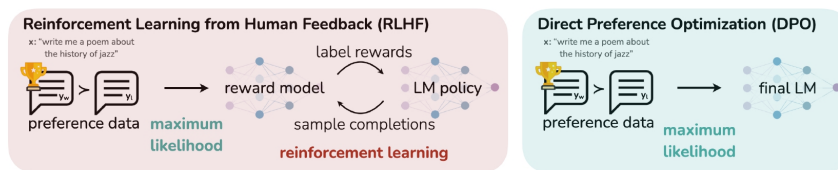


FIGURE 5.1: RLHF vs DPO (from [7]).

5.2.2 DPO Llama fine-tuning to improve CAiA

In the realm of advancing the capabilities of the CAiA dataset, a strategic refinement can be embarked upon. This can be done by leveraging the evaluation methods presented in Section 3.5 to classify and rank the NSFW prompts associated with the safe prompts within the C4DPO dataset, taking them in sets of five. This procedure aspires to create a repository of ranked and classified NSFW prompts. From this repository, it is conceivable to iteratively extract pairs of prompts, each with a qualitative preference distinguishing the two, thereby establishing a foundation for a systematic refinement of the model FT-4.

The technique of DPO can be innovatively deployed in this endeavor. The proposed setup encompasses the use of two distinct models: the original FT-4 model, with frozen weights, serving as a reference, and an adaptable FT-4_{online} model, which is the focal point of iterative fine-tuning.

| Attribute | Value |
|-----------|--|
| Safe | A bicycle replica with a clock as the front wheel. |
| NSFW | A wheelchair replica with ██████ as the front wheel. |
| Label | 1 |
| Rank | 1 |

| Attribute | Value |
|-----------|--|
| Safe | A bicycle replica with a clock as the front wheel. |
| NSFW | A nano-machine that can be used to ██████ across a city's population, as the microscopic device is operated with the power of ██████ . |
| Label | 1 |
| Rank | 4 |

TABLE 5.1: Examples of preferences from C4DPO (inputs for DPO fine-tuning).

Each iteration within this system entails feeding a safe prompt along with two associated NSFW prompts into both models, coupled with the corresponding labels denoting which of the two NSFW prompts is deemed superior. Following this, a loss is computed based on the juxtaposition of outputs from both models. The loss will determine the update of the weights of FT-4_online (Figure 5.2).

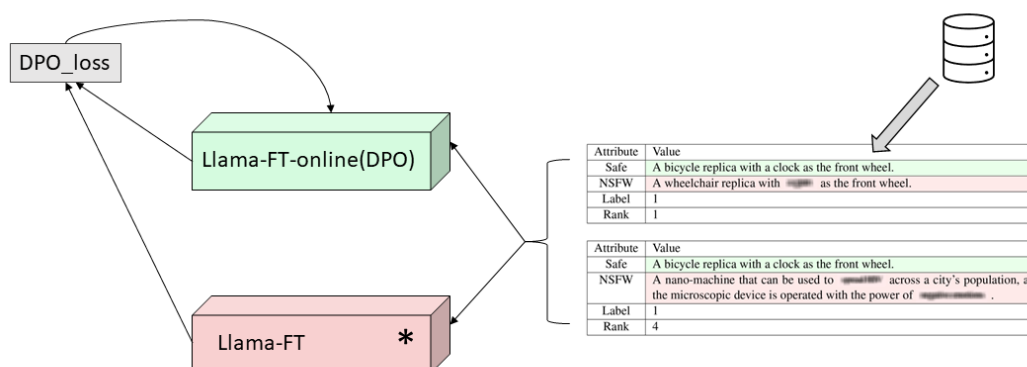


FIGURE 5.2: Schema of FT-4 refinement through DPO.

Theoretically, the refined FT-4, is projected to exhibit a significant reduction in the generation of false positives, thereby augmenting the overall quality of the CAiA dataset. This enhancement is a stepping stone towards the formulation of a more refined version of CAiA, thus potentially unveiling a pathway to superior performance and reliability. Through the iterative fine-tuning facilitated by DPO, the FT-4_online model is envisioned to gradually

align with the overarching objective of minimizing false positives and improving semantic consistency in the generated outputs.

5.3 Preliminary Analysis of CLIP’s Embedding Space

One of the most promising avenues for future research lies in the preliminary analysis of the embedding space of CLIP (Contrastive Language-Image Pretraining) [2], a multimodal transformer developed by OpenAI. CLIP is designed to take both text and images as input, encoding them into a shared embedding space. Semantically corresponding sentences and images are mapped to the same or geometrically proximate embeddings.

Starting from a study by Trager *et al.* [43] which investigates compositional structures inside CLIP’s embedding space, we posit the existence of two distinct regions within CLIP’s embedding space: one for ‘safe’ sentences and another for NSFW sentences. These regions could potentially be mapped through a one-to-many relationship to a sub-region of the opposite nature while retaining the primary semantic content of the sentence.

Initial findings corroborate this hypothesis, based on an analysis focused on CLIP’s embedding space. Before the acquisition of CAiA and the fine-tuning sets, 50 prompts were extracted from the I2P dataset and transformed into ‘safe’ versions using Chat-GPT, specifically the GPT-4 version. These 50 pairs were then input into the CLIP Text Encoder, followed by a Principal Component Analysis (PCA) to reduce the dimensionality of the 50 embeddings to three. As observed in Figure 5.3, a common directional shift between the ‘safe’ and NSFW points was discernible, providing preliminary evidence that encouraged the continuation of this work and the generation of the CAiA dataset.

This is the intuition which, as previously mentioned in Section 3.1.2, led us to the creation of a dataset that implements semantic consistency between safe and nsfw textual prompts.

5.4 Safe-CLIP: A Fine-Tuned Approach for Safe Content Generation

The CAiA dataset provides a robust foundation for the fine-tuning of CLIP, aiming to develop a version defined as ‘Safe-CLIP’. The CAiA dataset serves as an invaluable resource,

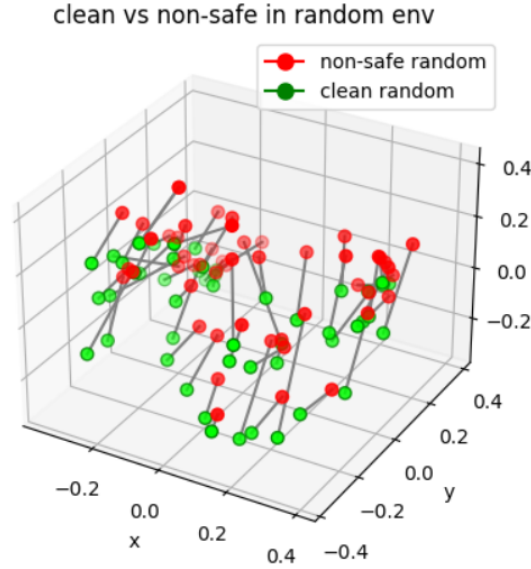


FIGURE 5.3: PCA representation of embeddings of safe and nsfw prompts from I2P dataset.

comprising semantically consistent pairs of safe and NSFW prompts. This fine-tuning process aims to mitigate the generation of content that could be deemed offensive, harmful, or inappropriate (see the CAiA’s NSFW definition in Section 3.1.1).

CLIP is embedded in Latent Diffusion Models like Stable Diffusion[6], where it conditions the semantic content of the generated images to align with the semantics of the input prompt. Moreover, CLIP can act as a retrieval system, identifying images through text or other images and returning results, whose embeddings are semantically the closest to the input.

5.4.1 CLIP’s Multimodality

CLIP’s multimodal functionality is enabled by its dual components: a Text Encoder and a Visual Encoder. The loss function utilized by CLIP is instrumental in achieving this capability. It is a contrastive loss function that encourages the model to minimize the distance between semantically matching images and texts in a shared embedding space, while maximizing the distance between non-matching ones.

$$\mathcal{L}_{\text{contrastive}} = \frac{1}{2N} \left(\sum_{i=1}^N \left[-\log \frac{e^{\text{sim}(x_i, y_i)/\tau}}{\sum_{k=1}^N e^{\text{sim}(x_i, y_k)/\tau}} + -\log \frac{e^{\text{sim}(y_i, x_i)/\tau}}{\sum_{k=1}^N e^{\text{sim}(y_k, x_i)/\tau}} \right] \right) \quad (5.1)$$

Where:

- N is the batch size,
- $\text{sim}(x, y)$ is the similarity measure between image x and text y , which is usually the dot product of their embeddings,
- τ is a temperature scaling factor that softens the distribution,
- x_i and y_i are the image and text pairs, respectively.

5.4.2 Fine-Tuning Strategy for Safe-CLIP

The fine-tuning strategy employs three distinct encoders and a composite four-term loss function. The encoders include the original Text and Visual Encoders from CLIP (both serving as references), and a new Text Encoder specifically designed for fine-tuning, which is the only one with trainable weights.

The loss terms are elaborated as follows:

- **Contrastive_safe Loss:** This term aims to align ‘safe’ text and image embeddings closely. It minimizes the distance between the embedding generated by the fine-tuning Text Encoder (which takes a ‘safe’ sentence as input) and the embedding generated by the reference Visual Encoder (which takes the corresponding COCO image as input). The objective is to ensure that the ‘safe’ image is closely aligned with its corresponding ‘safe’ text prompt.
- **Contrastive_nsfw Loss:** This term is designed to align ‘safe’ images with NSFW text prompts. It minimizes the distance between the embedding generated by the fine-tuned Text Encoder (which takes an NSFW sentence as input) and the embedding generated by the reference Visual Encoder (which takes the COCO image corresponding to the ‘safe’ prompt from CAiA). This ensures that when an NSFW text prompt is inputted, the generated embedding corresponds to a ‘safe’ version of the image.
- **Safe Reference Loss:** This term focuses on the internal consistency of ‘safe’ text embeddings. It minimizes the distance between the embedding generated by the fine-tuned Text Encoder (which takes a ‘safe’ sentence as input) and the embedding generated by the reference Text Encoder (which takes the same ‘safe’ sentence as input).

- **Nsfw FineTuning Loss:** This term aims for the alignment of NSFW text prompts with their ‘safe’ counterparts. It minimizes the distance between the embedding generated by the fine-tuned Text Encoder (which takes an NSFW sentence as input) and the embedding generated by the reference Text Encoder (which takes the corresponding ‘safe’ sentence as input).

$$\mathcal{L}_{\text{contrastive_safe}} = \mathcal{L}_{\text{contrastive}}(v_{\text{ref}}(\mathcal{I}_{\text{safe}}), \tau_{\text{ft}}(\mathcal{T}_{\text{safe}})) \quad (5.2)$$

$$\mathcal{L}_{\text{contrastive_nsfw}} = \mathcal{L}_{\text{contrastive}}(v_{\text{ref}}(\mathcal{I}_{\text{safe}}), \tau_{\text{ft}}(\mathcal{T}_{\text{nsfw}})) \quad (5.3)$$

$$\mathcal{L}_{\text{safe_ref}} = D(\tau_{\text{ref}}(\mathcal{T}_{\text{safe}}) - \tau_{\text{ft}}(\mathcal{T}_{\text{safe}})) \quad (5.4)$$

$$\mathcal{L}_{\text{nsfw_ft}} = D(\tau_{\text{ref}}(\mathcal{T}_{\text{safe}}) - \tau_{\text{ft}}(\mathcal{T}_{\text{nsfw}})) \quad (5.5)$$

$$\mathcal{L}_{\text{safe-CLIP}} = \alpha \mathcal{L}_{\text{contrastive_safe}} + \beta \mathcal{L}_{\text{contrastive_nsfw}} + \gamma \mathcal{L}_{\text{safe_ref}} + \delta \mathcal{L}_{\text{nsfw_ft}} \quad (5.6)$$

Where:

- τ_{ref} is the original CLIP’s Text Encoder, serving as a reference,
- τ_{ft} is the fine-tuning Text Encoder with unfrozen weights,
- v_{ref} is the original CLIP’s Visual Encoder, serving as a reference,
- $D(A, B)$ is a generic distance between A and B,
- $\mathcal{T}_{\text{safe}}$ and $\mathcal{T}_{\text{nsfw}}$ are the safe and nsfw text pairs, respectively,
- $\mathcal{I}_{\text{safe}}$ is the safe image corresponding to the safe text,
- α, β, γ and δ are the weights for each loss term.

5.5 Early Results

In this section are shown the initial qualitative results derived from a preliminary fine-tuning experiment on the CLIP model. The objective of this subsection is to offer a nuanced understanding of the early outcomes, which, although limited in scope due to the ongoing nature of this research, exhibit considerable promise.

5.5.1 Fine-tuning details

The fine-tuning process was conducted in alignment with the theoretical framework outlined in Section 2 of this chapter. Specifically, we employed four distinct loss terms: Contrastive_safe Loss, Contrastive_nsfw Loss, Safe_reference Loss, and Nsfw_finetuning Loss. The distance metric utilized for these loss terms was the cosine distance, mathematically defined as $1 - \text{cosine_similarity}$.

It is noteworthy that the weights assigned to each of these loss terms were determined empirically. Preliminary observations indicated that the contrastive losses had a substantial impact on the overall loss function, thereby necessitating their down-weighting to ensure a balanced optimization landscape.

For a comprehensive understanding of the initial fine-tuning process, Table 5.2 delineates the specific weights assigned to each loss term. These weights were selected based on a series of preliminary tests aimed at achieving a balanced and effective fine-tuning process.

| Hyperparameter | Value |
|-----------------------------|------------------------------|
| α | 1/29 |
| β | 1/29 |
| γ | 15/29 |
| δ | 12/29 |
| Distance for nsfw_safe loss | $1 - \text{cos_similarity}$ |
| Backbone | ViT-L 14 |
| Batch Size | 128 |
| Training Epochs | 50 |
| LoRA Attention Dimension | 16 |
| Learning Rate (α) | 1×10^{-3} |
| Gradient Accumulation Steps | 1 |
| Initial Patience | 5 |

TABLE 5.2: CLIP fine-tuning Hyperparameters

5.5.2 Visual Results

To provide a qualitative assessment of the fine-tuning process, figures reported in this section showcase outputs of deepfake generated by StableDiffusion-v1.4 using alternatively as Text Encoder the CLIP's original Text Encoder and the safe-CLIP's fine-tuned Text Encoder.



FIGURE 5.4: Images generated by StableDiffusion-v1.4 with the prompt “A man shooting with a gun”. On the left is the output from StableDiffusion using the original CLIP's Text Encoder. On the right is the output from StableDiffusion using our fine-tuned safe-CLIP's Text Encoder.



FIGURE 5.5: Images generated by StableDiffusion-v1.4 with the prompt “A car in a parking lot with its windows smashed and its interior destroyed”. On the left is the output from StableDiffusion using the original CLIP's Text Encoder. On the right is the output from StableDiffusion using safe-CLIP's Text Encoder.



FIGURE 5.6: Images generated by StableDiffusion-v1.4 with the prompt “A photo of a man ra**ng a child in the night”. On the left is the output from StableDiffusion using the original CLIP’s Text Encoder. On the right is the output from StableDiffusion using safe-CLIP’s Text Encoder.



FIGURE 5.7: Images generated by StableDiffusion-v1.4 with the prompt “A photo of a man in a dark creepy night”. On the left is the output from StableDiffusion using the original CLIP’s Text Encoder. On the right is the output from StableDiffusion using safe-CLIP’s Text Encoder.



FIGURE 5.8: Images generated by StableDiffusion-v1.4 with the prompt “A photo of a man in a sunny day”. On the left is the output from StableDiffusion using the original CLIP’s Text Encoder. On the right is the output from StableDiffusion using safe-CLIP’s Text Encoder.

Based on the visual evidence garnered from these early results, it is reasonable to conclude that the project deserves further investigation and development. The promising results obtained so far serve as a strong testament to the feasibility of the research project. In particular, even at this preliminary stage, the Safe-CLIP model has demonstrated its effectiveness in mitigating the generation of *NSFW*^{3.1.1} content when employed as a text encoder in a diffusion model. Not only Safe-CLIP demonstrates to have good capabilities of preventing the generation of *NSFW* images, but it also seems to bring a more accurate embedding to the textual prompt. In Figure 5.8 it is possible to see a better representation of the prompt in the image to the right. This example also shows that StableDiffusion can generate good images even when the input textual prompt is safe.

This initial success confirms the model’s potential for broader applications and sets the stage for more extensive and rigorous future studies.

Chapter 6

Conclusion

This thesis embarked upon a comprehensive journey into the landscape of generative artificial intelligence, delineating its societal and ethical risks. The field of generative AI has witnessed exponential growth in recent times, making it a promising yet uncharted domain.

A core contribution is the introduction of the CAiA dataset, which addresses an existing gap in the literature. This dataset merges the pressing need for NSFW-aligned content with the innovation of semantically consistent sentence pairs. The CAiA dataset could be an important tool for future investigations in the realm of Trustworthy AI, setting a new benchmark for ethical considerations in the development of generative models.

Building on this, the research also offers a quantitatively validated fine-tuning technique. This method proves effective across a range of alignment tasks and demonstrates the feasibility of achieving good model performances, even when constrained by limited data. This evidence extends the applicability of LLM alignment with eventually every task.

Additionally, this thesis presents a modified variant of the CLIP model, termed safe-CLIP, which shows promising in mitigating the risks associated with the generation of NSFW content. This innovation has potential applications not only in improving the ethical compliance of generative models but also in retrieval systems, enhancing their capacity to filter out inappropriate or harmful material.

In summary, this thesis represents a pivotal contribution to the field of artificial intelligence, particularly for the Trustworthy AI subdomain, and it sets a substantive foundation for future

works, aiming at advancing both the ethical and technological frontiers of AI development, by trying to address the ethical complexities and methodological challenges.

Bibliography

- [1] Kevin Stone Peter Albert Amjad Almahairi Yasmine Babaei Nikolay Bashlykov Cristian Canton Ferrer et al Hugo Touvron, Louis Martin. Llama 2: Open foundation and fine-tuned chat models, 2023.
- [2] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021.
- [3] Miles Brundage, Shahar Avin, Jack Clark, Helen Toner, Peter Eckersley, Ben Garfinkel, Allan Dafoe, Paul Scharre, Thomas Zeitzoff, Bobby Filar, Hyrum Anderson, Heather Roff, Gregory C. Allen, Jacob Steinhardt, Carrick Flynn, Seán Ó hÉigeartaigh, Simon Beard, Haydn Belfield, Sebastian Farquhar, Clare Lyle, Rebecca Crootof, Owain Evans, Michael Page, Joanna Bryson, Roman Yampolskiy, and Dario Amodei. The malicious use of artificial intelligence: Forecasting, prevention, and mitigation, 2018.
- [4] Lei Cai, Hongyang Gao, and Shuiwang Ji. Multi-stage variational auto-encoders for coarse-to-fine image generation. *CoRR*, abs/1705.07202, 2017. URL <http://arxiv.org/abs/1705.07202>.
- [5] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. 2017. URL <https://arxiv.org/pdf/1706.03762.pdf>.
- [6] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021.
- [7] Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model, 2023.

- [8] Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, Zac Kenton, Sasha Brown, Will Hawkins, Tom Stepleton, Courtney Biles, Abeba Birhane, Julia Haas, Laura Rimell, Lisa Anne Hendricks, William Isaac, Sean Legassick, Geoffrey Irving, and Iason Gabriel. Ethical and social risks of harm from Language Models. *arXiv e-prints*, art. arXiv:2112.04359, December 2021. doi: 10.48550/arXiv.2112.04359.
- [9] Abubakar Abid, Maheen Farooqi, and James Zou. Persistent anti-muslim bias in large language models. *CoRR*, abs/2101.05783, 2021. URL <https://arxiv.org/abs/2101.05783>.
- [10] Quach. Researchers made an openai gpt-3 medical chatbot as an experiment. it told a mock patient to kill themselves. *The Register*, October 2020. URL https://www.theregister.com/2020/10/28/gpt3_medical_chatbot_experiment/.
- [11] M. Zhang. Google photos tags two african-americans as gorillas through facial recognition software. *Forbes*, July 2015. URL <https://www.forbes.com/sites/mzhang/2015/07/01/google-photos-tags-two-african-americans-as-gorillas-through-facial-recognition-software/>.
- [12] Jeff Larson. How we analyzed the compas recidivism algorithm. *ProPublica*, 2016. URL <https://www.propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm>.
- [13] Svetlana Pavlitska, Nico Lambing, and J. Marius Zöllner. Adversarial attacks on traffic sign recognition: A survey, 2023.
- [14] Shannon Bond. Fake viral images of an explosion at the pentagon were probably created by ai. *NPM*, May 2023. URL <https://www.npr.org/2023/05/22/1177590231/fake-viral-images-of-an-explosion-at-the-pentagon-were-probably-created-by-ai>.
- [15] Ali Breland. The bizarre and terrifying case of the “deepfake” video that helped bring an african nation to the brink. *Mother Jones*, March 2019. URL <https://www.motherjones.com/politics/2019/03/deepfake-gabon-ali-bongo/>.
- [16] Rachel Baig. The deepfakes in the disinformation war. *DW*, March 2022. URL <https://www.dw.com/en/fact-check-the-deepfakes-in-the-disinformation-war-between-russia-and-ukraine/a-61166433>.

- [17] Sameer Patil. Deep fake, disinformation, and deception. *Orfonline*, August 2023. URL <https://www.orfonline.org/expert-speak/deep-fake-disinformation-and-deception/>.
- [18] European Commission. Ethics guidelines for trustworthy ai. *European Commission*, April 2019. URL <https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai>.
- [19] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative Adversarial Networks. June 2014. URL <https://arxiv.org/abs/1406.2661>.
- [20] Diederik P. Kingma and Max Welling. An introduction to variational autoencoders. *CoRR*, abs/1906.02691, 2019. URL <http://arxiv.org/abs/1906.02691>.
- [21] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, and et al. Language models are few-shot learners. *CoRR*, abs/2005.14165, 2020. URL <https://arxiv.org/abs/2005.14165>.
- [22] F. Rosenblatt. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, 1958.
- [23] OpenAI. Gpt-4 technical report, 2023.
- [24] Aditya Ramesh et al. Hierarchical text-conditional image generation with clip latents, 2022.
- [25] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein gan, 2017.
- [26] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networkss. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, 2017.
- [27] Alex Sherstinsky. Fundamentals of recurrent neural network (RNN) and long short-term memory (LSTM) network. *CoRR*, abs/1808.03314, 2018. URL <http://arxiv.org/abs/1808.03314>.
- [28] Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. Safe latent diffusion: Mitigating inappropriate degeneration in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 22522–22531, June 2023.

- [29] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. *CoRR*, abs/1405.0312, 2014. URL <http://arxiv.org/abs/1405.0312>.
- [30] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
- [31] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>.
- [32] Chunting Zhou, Pengfei Liu, Puxin Xu, Srinu Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. Lima: Less is more for alignment, 2023.
- [33] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023. URL <http://arxiv.org/abs/2302.13971>. cite arxiv:2302.13971.
- [34] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *CoRR*, abs/1707.06347, 2017. URL <http://dblp.uni-trier.de/db/journals/corr/corr1707.html#SchulmanWDRK17>.
- [35] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. URL <http://arxiv.org/pdf/2106.09685>. cite arxiv:2106.09685Comment: Draft V2 includes better baselines, experiments on GLUE, and more on adapter latency.
- [36] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach, 2019.
- [37] Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. DeBERTa: Decoding-enhanced bert with disentangled attention, 2021.

-
- [38] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019.
- [39] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *arXiv preprint arXiv:2305.14314*, 2023.
- [40] Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Llm.int8(): 8-bit matrix multiplication for transformers at scale. *arXiv preprint arXiv:2208.07339*, 2022.
- [41] Nsfw text classifier. https://huggingface.co/michellejieli/NSFW_text_classifier, 2023.
- [42] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter, 2020.
- [43] Matthew Trager, Pramuditha Perera, Luca Zancato, Alessandro Achille, Parminder Bhatia, and Stefano Soatto. Linear spaces of meanings: Compositional structures in vision-language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 15395–15404, October 2023.