

# An introduction to Large Language Models and their relevance for statistical offices

Dario Buono, Marius Felecan and  
Cristiano Tessitore

2024 edition





# **An introduction to Large Language Models and their relevance for statistical offices**

**Dario Buono, Marius Felecan and  
Cristiano Tessitore**

**2024 edition**

This document should not be considered as representative of the European Commission's official position.

Luxembourg: Publications Office of the European Union, 2024



© European Union, 2024

The reuse policy of European Commission documents is implemented by Commission Decision 2011/833/EU of 12 December 2011 on the reuse of Commission documents (OJ L 330, 14.12.2011, p. 39). Unless otherwise noted, the reuse of this document is authorised under a Creative Commons Attribution 4.0 International (CC BY 4.0) licence (<https://creativecommons.org/licenses/by/4.0/>). This means that reuse is allowed provided appropriate credit is given and any changes are indicated.

For any use or reproduction of elements that are not owned by the European Union, permission may need to be sought directly from the respective rightholders. The European Union does not own the copyright in relation to the following elements:

cover © Eurostat

Collection: Statistical working papers

Theme: General and regional statistics

PDF

ISBN 978-92-68-11346-2

ISSN 2315-0807

doi:10.2785/716217

KS-TG-24-001-EN-N

# Foreword

In 1975, Gordon Moore postulated the doubling of the number of transistors in an integrated circuit every 2 years. His prediction has held since then and became what is now known as Moore's law. While the exponential growth in computer power has been predictable during the last 50 years, it has unleashed unprecedented developments in software applications, which had disruptive effects fundamentally changing the digital landscape and brought the way we are interacting professionally and privately to new levels. One of these new levels seemed to be achieved by the advent of Generative Artificial Intelligence, which generates new content based on patterns learned from existing data. Especially the launch of ChatGPT 3.5 at the end of 2022 started a boom in AI development.

The European Commission already launched first initiatives related to AI in 2018. They resulted in a proposal for a regulation on harmonised rules on AI, which reached a political agreement between the European Council and the European Parliament in December 2023. The activities of the European Union try to unleash the potential of AI for the European economy, society and the public sector, while emphasizing the importance of safe, transparent, and human-centred use of AI technologies.

In January 2024, the European Commission adopted the AI@EC communication outlining the strategies for

improving its own capabilities as regards the use of AI. At the same time, the European Statistical System is exploring the potential of Artificial Intelligence use for developing, producing and publishing official statistics. In 2023, Eurostat has launched a grant, the one-stop-shop for Artificial Intelligence and Machine Learning to guide and support the members of the European Statistical System to experiment and develop AI/ML solutions for official statistics.

AI Language Models promise efficient data processing, ensuring high-quality and consistent insights from vast textual unstructured datasets. They enable advanced text analysis, automated report generation, and interactive query systems, enhancing transparency and public engagement. With multilingual capabilities, scalability, and the potential for cost savings, AI Language Models stand as a critical tool for modernising processes in statistical offices. As this technology is embraced more and more, it's paramount to ensure its responsible and ethical use, safeguarding the integrity of the data and the trust of the stakeholders of official statistics' organisations.

Albrecht Wirthmann  
Head of Unit, Eurostat A.5  
Methodology; Innovation in official statistics

**Keywords:** Artificial Intelligence, Large Language Model, Official Statistics

**Authors:** Dario Buono<sup>1</sup>, Marius Felecan<sup>2</sup>, Cristiano Tessitore<sup>3</sup>

**Acknowledgements:** The authors would like to thank Eurostat colleagues Anu Ahola for her valuable help about security, Jean-Marc Museux, Alvaro Diez Soto, Martina Patone, Fabio Sartori and Peter Struijs for their review, Benoit Schneider for his suggestions, Luigi di Razza for editorial help.

1 Eurostat, [Dario.Buono@ec.europa.eu](mailto:Dario.Buono@ec.europa.eu)

2 eVerbum, under contract with GOPA Worldwide Consultants, [Marius@eVerbum.com](mailto:Marius@eVerbum.com)

3 Eurostat, [Cristiano.Tessitore@ec.europa.eu](mailto:Cristiano.Tessitore@ec.europa.eu)

# Contents

<b>Foreword</b> .....	<b>3</b>
<b>Contents</b> .....	<b>5</b>
<b>Introduction</b> .....	<b>8</b>
<b>Motivation</b> .....	<b>8</b>
<b>Scope</b> .....	<b>9</b>
<b>Overview</b> .....	<b>9</b>
<b>Part 1. General Aspects and Literature Review</b> .....	<b>11</b>
<b>1. Terms and Methods</b> .....	<b>12</b>
<b>1.1. Definitions</b> .....	<b>12</b>
<b>1.2. LLM as a Path to AGI</b> .....	<b>14</b>
1.2.1. What Is An AGI?.....	14
1.2.2. Measuring AGI.....	14
1.2.3. The Risk of Superintelligence.....	15
<b>2. LLM Ecosystem and Literature Review</b> .....	<b>18</b>
<b>2.1. Architecture Evolution and Diversity</b> .....	<b>18</b>
2.1.1. History.....	19
2.1.2. Transformer Architecture.....	20
2.1.3. Variants And Models.....	20
2.1.4. Emerging Architectures.....	20
<b>2.2. The Challenges of Tasks Classifications</b> .....	<b>21</b>
<b>2.3. Pre-Training and Pre-Training Corpora</b> .....	<b>22</b>
2.3.1. Training Datasets.....	22
2.3.2. Size, Quality, and Composition.....	22
<b>2.4. Fine-Tuning Techniques</b> .....	<b>22</b>
2.4.1. Repurposing The Model.....	22
2.4.2. Full Fine-Tuning (FT).....	22
2.4.3. Unsupervised FT.....	22
2.4.4. Supervised FT.....	23
2.4.5. Reinforcement Learning From Human Feedback.....	23

2.4.6. Parameter Efficient Fine-Tuning (PEFT).....	23
<b>2.5. Small Language Models (SLM).....</b>	<b>23</b>
<b>2.6. Model Evaluation.....</b>	<b>24</b>
<b>2.7. Prompting.....</b>	<b>24</b>
2.7.1. Main Prompting Approaches.....	24
2.7.2. Prompts Aspects.....	24
2.7.3. Prompt Patterns.....	25
2.7.4. (Auto)Reflexion.....	25
<b>2.8. Infrastructure.....</b>	<b>25</b>
<b>2.9. Scaling Laws and Emergent Abilities.....</b>	<b>26</b>
<b>2.10. Costs and Benefits.....</b>	<b>26</b>
2.10.1. Costs.....	26
2.10.2. Benefits.....	27
<b>2.11. Usage Process.....</b>	<b>27</b>
<b>3. Ethical and Legal Considerations.....</b>	<b>28</b>
<b>3.1. Ethics.....</b>	<b>28</b>
3.1.1. Helpful, Honest, and Harmless AI.....	28
3.1.2. The European Commission High-Level Expert Group on Artificial Intelligence.....	29
3.1.3. European Commission: On Artificial Intelligence - An European Approach to Excellence and Trust.....	31
3.1.4. OECD Council (Ethical) Recommendations .....	32
3.1.5. Australia’s AI Ethics Principles.....	33
<b>3.2. A General Policy for AI in Europe.....</b>	<b>33</b>
<b>3.3. Generative AI Specific Guidelines for Public Organisations.....</b>	<b>34</b>
3.3.1. European Commission.....	34
3.3.2. Australian Public Service.....	34
<b>3.4. An Agile and Adaptive Policy Approach for LLMs .....</b>	<b>35</b>
<b>3.5. Open Source AI.....</b>	<b>36</b>
<b>4. Controversies and Security Issues.....</b>	<b>37</b>
<b>4.1. Controversies and Limitations.....</b>	<b>37</b>
4.1.1. Controversies Surrounding LLMs.....	37
4.1.2. Limitations.....	38
<b>4.2. Security Issues.....</b>	<b>40</b>
4.2.1. Attacks Against LLM Applications.....	40
4.2.2. Using LLMs for Cybercrime.....	42
4.2.3. Create a Robust and Safe LLM.....	42
4.2.4. Addressing Quality and Safety by Design.....	43



<b>Part 2. Selected Use Cases for Official Statistics</b>	<b>45</b>
<b>5. Data-Centric Operations</b>	<b>47</b>
5.1. Numerical data processing	47
5.2. Text Data Cleaning and Preprocessing	47
5.3. Sentiment Analysis	48
5.4. Metadata Generation	49
<b>6. User Support and Assistance</b>	<b>50</b>
6.1. Chatbots	50
6.2. Multilingual Support	51
<b>7. Automated Content Creation</b>	<b>53</b>
<b>Part 3. Use cases for scientists and developers</b>	<b>55</b>
<b>8. Software Development Support</b>	<b>56</b>
8.1. Documenting and Commenting	56
8.2. Coding and Scripting Guidance	57
8.3. Testing and Debugging	57
8.3.1. Code Testing	57
8.3.2. Code Review and Optimization	57
8.3.3. Debugging Support	58
8.4. Security Assistance	58
8.5. Reverse Engineering of Legacy Code	59
<b>9. Research, Education and Training</b>	<b>60</b>
9.1. Research Assistance	60
9.2. Simulation and Scenario Analysis	61
9.3. Education and Training	61
<b>Part 4. The New AI-powered (Data) Scientist</b>	<b>63</b>
<b>10. A Team of One</b>	<b>65</b>
<b>11. LLM Skills and New AI Jobs</b>	<b>67</b>
11.1. New AI Jobs	70
11.2. A Simplified Method to Measure LLM Impact on Jobs	70
11.2.1. Tasks	70
11.2.2. Scoring	70
11.2.3. Results	71
<b>Conclusions</b>	<b>72</b>
<b>References</b>	<b>73</b>

# Introduction

## Motivation

In the evolving landscape of data science and analytics, the need for robust and efficient tools has never been more pronounced. Like for Eurostat, whose mission is “to provide high quality statistics and data on Europe”, official statistics organisations are entrusted with the critical task of providing accurate, timely, and relevant data to inform policy, drive decision-making, and foster public understanding. Traditional methods, while fundamental, often struggle to keep pace with the sheer volume and complexity of modern data. Large Language Models (LLMs) and their transformative capabilities in Natural Language Processing (NLP) offer a solution to many of the challenges faced by statistics organisations, from data processing to advanced textual analysis. As LLMs integrate into workflows, it’s essential to understand their strengths, limitations, and potential pitfalls.

This document is the result of literature review combined with the best practices as seen in industry implementations. Given that this domain is in a constant and rapid evolution some of the sources used as references are preprints.

This document aims at ranking methods and tools for harnessing the potential of LLMs while ensuring ethical, transparent, and responsible use. Our motivation is to elevate the quality and impact of our work, while upholding the trust and confidence of our stakeholders.

In the rapidly evolving domain of data science and official statistics, the introduction of LLMs has opened many potential applications. However, with limited resources, time, and expertise, it’s impractical for organisations to adopt and integrate every possible application simultaneously. This necessitates a structured approach to prioritise applications that align most closely with the organisation’s goals, stakeholder needs, and strategic direction. By establishing clear criteria for prioritisation, organisations need to cope with:

- **Ensuring Alignment with Strategic Goals:** Not all applications will align with the current objectives or long-term vision of the organisation. Criteria help ensure that chosen applications serve the broader mission.
- **Optimising Resource Allocation:** Resources, whether financial, human, or technological, are finite. Prioritising applications ensures that these resources are directed towards the most impactful projects.
- **Managing Risk:** Some applications may come with higher risks, be it in terms of data security, ethical concerns, or potential for errors. Criteria can help weigh these risks against potential benefits.
- **Enhancing Stakeholder Engagement:** Prioritising applications that have a direct impact on stakeholder needs can improve engagement, trust, and collaboration.
- **Facilitating Decision-making:** Clear criteria provide a framework for decision-making, reducing ambiguity and promoting consistency in choices.

**TABLE 1:**

## Comprehensive list of criteria with descriptions

1. **Strategic Alignment: Does the application align with the organisation’s current objectives and long-term vision?**
2. **Impact Potential: What is the potential of the application to significantly enhance data quality, insights, or stakeholder engagement?**
3. **Resource Requirements: How intensive are the resource demands of the application in terms of finances, manpower, and technology?**

4. **Risk Profile: What are the potential risks associated with the application, including data security, ethical concerns, and error potential?**
5. **Stakeholder Value: How valuable is the application to external stakeholders, such as policymakers, the public, or partner organisations?**
6. **Operational Efficiency: Can the application streamline workflows, reduce manual labour, or improve overall operational efficiency?**
7. **Scalability: Does the application have the potential to handle increasing data volumes or expand in scope as the organisation grows?**
8. **Innovation Potential: Does the application offer a novel approach or solution that sets the organisation apart from peers?**
9. **Implementation Timeframe: How long will it take to fully integrate the application and start seeing tangible results?**
10. **Cost Efficiency: Beyond initial costs, what are the long-term savings or cost benefits associated with the application?**
11. **Ethical and Responsible Use: Does the application adhere to ethical standards, and does it promote responsible data handling and usage?**
12. **Feedback and Evaluation: Is there a mechanism in place to gather feedback and continuously evaluate the application's effectiveness?**

By considering these criteria, official statistics organisations can make informed decisions about which LLM applications to prioritise, ensuring that their choices drive value, impact, and strategic alignment.

## Scope

This document provides a framework for the integration of LLMs within official statistics organisations. It outlines a

range of potential applications that harness the capabilities of LLMs to enhance data processing, analysis, and stakeholder engagement. Each application's relevance is assessed based on a set of criteria, focusing on its impact, efficiency, scalability, and alignment with ethical standards.

The aim is to offer a clear framework for organisations to leverage LLMs responsibly and effectively, ensuring that the technology serves the broader goals of accuracy, transparency, and public trust.

Each issue discussed in the chapters about use cases follows a standardised structure with three parts: a description, a list of options and a list of ranked alternatives. The description is free text presenting the problem.

The options list, without pretending to be exhaustive, presents various possibilities to deal with the specific problem treated in the item. Out of these options, three ranked alternatives are highlighted:

(A) Best alternative: should always be the target for users.

(B) Acceptable alternative: retained only if time or resource issues prevent alternative (A).

(C) Alternative to be avoided: not a recommended option.

The objective of the approach is to help users apply the best alternative whenever possible. It should then be a feasible target for users. It should always be achievable with a reasonable effort unless some production or institutional constraints prevent it.

The acceptable alternative (B) should be viewed as an intermediate step towards the achievement of alternative (A). It could also be considered the target for a limited number of cases when specific data issues, user requests, time or resource constraints prevent the achievement of the alternative (A).

The alternative to be avoided (C) includes some procedures that are not recommended.

## Overview

This document is structured in four parts.

Part 1 is about general aspects of AI/LLMs, like definitions, EU policies, and technical aspects. In particular, in Chapter 2, the reader can see why it is hard to match NLP tasks to LLMs; next, specific architectures are presented, delving on processes involved, and infrastructure used and why they are important. Finally, in Chapter 4, discussions will focus on the inherent controversies and limitations associated with a rapidly evolving domain such as LLMs. Given the profound

implications of LLMs not only in the realm of scientific research but also for the broader spectrum of human society, it is imperative to understand these nuances.

Parts 2 and 3 include an analysis of the most frequent LLMs use cases as of today grouped by the most common scheme of classification from the literature:

In Chapter 5 - Data-Centric Operation - are presented the operations LLMs can do to process and clean text data, and it is explained why, for now, numeric data processing is out of the scope of these models.

Chapter 6 - User support and Assistance - presents one of the most popular domains of applications for LLMs, starting with the biggest: Google and Microsoft including GEMINI and respectively GPT4 models in their search engines for natural language AI query assistants.

Chapter 7 focuses on Automated content creation, while Chapter 8 - Software Development Support - delves into another very profitable domain of applicability for LLMs. Most IDEs<sup>4</sup> come - or can be upgraded - with a plugin - an

AI assistant - able to solve some aspects of the software development cycle. Some preliminary estimates show efficiency improvements up to 50% with the usage of AI assisted IDEs.

Chapters 9 contains use cases showing how LLMs can directly affect scientists in research, education and training environments.

Part 4 addresses skills, tools and workplace of the new (data) scientist in the context of smart agents (such as LLMs). What is the profile of the data scientist of the future? How the research and development teams will be composed and how the humans will "cooperate" with machines? This chapter tries to feed the discussion with meaningful elements from industry and research environments.

This publication includes contributions from different LLMs. They have assisted the authors in providing insights and conducting language-related tasks. While every effort has been made to ensure accuracy and coherence, the unique nature of AI-generated content may influence the style and presentation of the information within.

4 An integrated development environment is a software application that provides comprehensive facilities for software development. An IDE normally consists of at least a source-code editor, build automation tools, and a debugger.

Part

1

**General Aspects and  
Literature Review**

# 1

## Terms and Methods

Since the release of the revolutionary paper “Attention Is All You Need” (Vaswani, 2017), the field of LLMs has witnessed a meteoric rise, revolutionising the way scientists approach natural language processing and understanding. The rapid evolution of this subject in recent years is evident in the myriad of breakthroughs and innovations. For instance, architectures like BERT (Bidirectional Encoder Representations from Transformers), GPT (Generative Pre-trained Transformer), RoBERTa (Robustly Optimised BERT Pretraining Approach), and T5 (Text-to-Text Transfer Transformer) have not only set new benchmarks in various NLP tasks but have also found applications in diverse domains, from sentiment analysis to machine translation. For an official statistics organisation, leveraging LLMs is imperative to harness the vast amounts of textual data, extract meaningful insights, and make data-driven decisions.

The question then arises: Why to use them now? The answer lies in the recent advancements and the sheer capabilities of these models. The “Sparks of Artificial General Intelligence” paper on GPT-4 (Bubeck, 2023) underscores the potential of LLMs, highlighting their ability to perform tasks that were once thought to be exclusive to human cognition. In this era, where data is abundant and the need for accurate interpretation is paramount, LLMs offer a promising avenue to transform the way official statistics are understood and utilised.

The rapid pace of scientific advancement in this domain is evident, with novel architectures and methodologies emerging at an almost daily frequency. However, the market landscape for these computational tools appears to have reached a state of maturity. Prominent technology conglomerates such as Microsoft (Microsoft AI, 2023),

Google (Google Cloud, 2023), and OpenAI (OpenAI, 2023a) have established themselves as leading service providers. Concurrently, niche enterprises like Hugging Face (Hugging Face, 2023) and Weights & Biases (Weights & Biases, 2023) have emerged, anchored by fervent academic and developer communities.

Furthermore, the market is now enriched with a plethora of high-performing open-source models. Noteworthy contributions include the LLaMa (Large Language Model Meta AI) series by Meta and the Falcon LLM series by the Technology Innovation Institute based in Abu Dhabi. The confluence of technological giants, innovative startups, and open-source contributions marks a vibrant and collaborative ecosystem, poised to drive further advancements in the field.

### 1.1. Definitions

**AI – Artificial Intelligence:** is the branch of computer science that focuses on creating machines capable of performing tasks that typically require human intelligence, encompassing a range of functionalities such as problem-solving, learning, and understanding natural language.

**AI Language Models:** AI language models are advanced software systems designed to understand, generate, and manipulate natural language. They are a key component of natural language processing (NLP).

**Anomaly Detection:** Identifying unusual patterns or outliers in datasets that might indicate errors, significant events, or other noteworthy occurrences.

**Architecture:** The specific design and structure of a neural network. For example, GPT (Generative Pre-trained Transformer) is an architecture used for LLMs.

**Artificial General Intelligence (AGI):** Also known as Strong AI, it refers to a type of artificial intelligence that possesses the ability to understand, learn, and apply its intelligence to a wide variety of problems, much like a human being. AGI can generalise its learning and reasoning abilities across a broad range of tasks.

**Attention Mechanism:** A component in the Transformer architecture that allows the model to focus on different parts of the input when producing an output.

**Data Processing:** The act of collecting, cleaning, and organising raw data to extract meaningful insights.

**Data Visualization:** The representation of data in a visual format, such as charts, graphs, or maps, to make it easier to understand and interpret.

**Document Classification:** The process of categorising documents based on their content or other attributes.

**Embedding:** A vector representation of words or tokens in a high-dimensional space, which captures semantic meaning.

**Ethical and Responsible Use:** Adhering to ethical standards and ensuring responsible handling and usage of data or technology.

**Fine-tuning:** The process of training a pre-trained model on a smaller, specific dataset to adapt it to a particular task.

**Forecasting:** The process of predicting future trends or outcomes based on historical data and analysis and dynamic modelling.

**Generative AI:** A type of AI that can create new content, such as text, images, music, and code, by learning from existing data.

**Hallucination:** It refers to the generation of incorrect, misleading, or nonsensical information. This occurs due to the model's limitations in understanding and interpreting data, often resulting in factual inaccuracies, illogical responses, or fictional narratives.

**Knowledge Cutoff:** The last date at which the model was updated with new information.

**Knowledge Graph:** A structured representation of interconnected data that links related concepts, entities, and statistics.

**Large Language Models (LLMs):** Advanced machine learning models trained on vast amounts of textual data,

designed to understand and generate human-like text based on the input they receive.

**Metadata:** Data that provides information about other data, often describing its content, format, source, and other attributes.

**Multilingual Data Processing:** The ability to analyse and process data in multiple languages.

**Narrow AI:** Also known as Weak AI, this type of AI is designed and trained for a particular task.

**Natural Language Processing (NLP):** A field of artificial intelligence that focuses on the interaction between computers and humans through natural language. It enables machines to read, understand, and derive meaning from human languages.

**Official Statistics:** Data and information produced by official agencies and entities, often governmental, to provide timely, reliable, and relevant statistical information to the public, policymakers, and other stakeholders.

**Operational Efficiency:** The ability of an organisation to deliver results in the most cost-effective and timely manner.

**Overfitting:** When a model performs very well on training data but poorly on unseen data because it's too closely adapted to the training set.

**Parameter:** An element of the model that is learned from the training data, such as weights and biases in a neural network.

**Pre-trained Language Models (PLMs):** A model that has already been trained on a large dataset and can be fine-tuned for specific tasks. In this context it is about an already trained LLM. We use the two terms interchangeably.

**Prompt:** A piece of text given as input to the model to get a response or output.

**Query Systems:** Platforms or tools that allow users to ask specific questions and retrieve relevant data or answers.

**Regularisation:** Techniques used to prevent overfitting, such as dropout or weight decay.

**Report Generation:** The creation of structured documents or summaries based on data analysis.

**Scalability:** The capacity of a system, model, or application to handle growth or increased demand efficiently.

**Semantic Search:** A search process that understands the context and semantics of queries, rather than relying solely on keyword matches.

**Small Language Models (SLMs):** Small Language Models are scaled-down versions of larger language models, designed to perform similar tasks but with reduced computational resources and training data requirements.

**Strategic Alignment:** Ensuring that initiatives, projects, or tools align with the broader objectives and long-term vision of an organisation.

**Survey Analysis:** The process of analysing responses from surveys to derive patterns, insights, and conclusions.

**Temperature:** A hyperparameter used during sampling to control the randomness of the model's output. Higher values make outputs more random, while lower values make them more deterministic.

**Text Analysis:** The process of examining textual content to extract information and insights, often using NLP techniques.

**Token:** A unit of text, which could be as short as one character or as long as one word.

**Tokenization:** The process of converting a sequence of characters in text into a sequence of tokens.

**Transformer:** A type of deep learning model architecture introduced in a paper by Vaswani et al. in 2017. It's particularly known for its self-attention mechanism.

**Transfer Learning:** Leveraging a model that has been trained on one task as the basis for training on a different, but related, task.

**Zero-shot, One-shot, Few-shot Learning:** Refers to the model's ability to perform tasks without any examples (zero-shot), with one example (one-shot), or with a few examples (few-shot).

## 1.2. LLM as a Path to AGI

*"the next wave of AI aspires to create machines that will function more as colleagues than as tools"*

from 'The limits of machine intelligence' (Shevlin, 2019)

### 1.2.1. What Is An AGI?

Artificial General Intelligence (AGI), also known as "strong AI" or "full AI", is a hypothetical type of intelligent agent that possesses the ability to learn and perform any intellectual task that a human or animal can. In essence, AGI would represent a significant leap forward in artificial intelligence, surpassing the capabilities of current (narrow) AI systems - like LLMs - that are designed to handle specific tasks.

Creating AGI is a central goal of many AI researchers and organisations, including OpenAI, DeepMind, and Anthropic. The development of AGI holds immense potential to revolutionise various aspects of society, from scientific discovery to technological advancements.

While AGI remains a hypothetical concept, it has sparked extensive discussions and debates in science fiction and future studies. The potential implications of AGI are profound, raising questions about the future of work, the role of technology in society, and even the nature of consciousness itself.

In contrast to AGI, weak AI, also known as narrow AI, is characterised by its ability to solve specific problems but lacks general cognitive abilities. Examples of weak AI include systems designed for facial recognition, natural language processing, or chess playing.

The pursuit of AGI is driven by the desire to create intelligent machines that can not only replicate human capabilities but also surpass them. AGI has the potential to transform our world in ways that are still unimaginable, and its development is eagerly anticipated by researchers and society at large.

### 1.2.2. Measuring AGI

Measuring the capabilities and performance of an AGI is a complex and evolving area of research (Morris, 2023). Unlike narrow AI, which is designed to perform specific tasks, AGI is intended to understand, learn, and apply its intelligence broadly and flexibly, similar to human intelligence. Some of the key approaches and metrics used to measure AGI may be synthesised as follow:

- **General Problem-Solving Ability:** This includes evaluating an AGI's performance across a wide range of tasks that require different types of intelligence, such as language comprehension, logical reasoning, and creative problem-solving.
- **Learning Efficiency:** This measure shows how quickly and efficiently an AGI can learn new tasks or adapt to new environments. It's not just about the final performance but also how rapidly it can reach that level of performance.
- **Adaptability and Transfer Learning:** This assesses the AGI's ability to apply knowledge and skills learned in one context to novel or unrelated situations. A high level of transfer learning capability is a key indicator of general intelligence.
- **Autonomous Goal Setting and Achievement:** This involves the AGI's ability to set its own goals based on its



understanding of the world and then achieve these goals in a flexible, adaptive manner.

- **Understanding and Reasoning:** This includes tests for the AGI's ability to comprehend complex ideas, reason through problems, and understand abstract concepts.
- **Creativity and Innovation:** This can be evaluated by the AGI's ability to generate new ideas, solutions, or artworks that are novel and valuable.
- **Emotional and Social Intelligence:** This involves assessing the AGI's ability to understand and respond to emotional cues and engage in complex social interactions.
- **Ethical and Moral Reasoning:** As AGI systems approach human-level intelligence, their ability to understand and apply ethical principles becomes increasingly important.

**TABLE 2:**

## OpenAI: an organisation aiming at AGI

**OpenAI: "Creating safe AGI that benefits all of humanity"**

OpenAI was initially founded in 2015 by Sam Altman, Elon Musk, Ilya Sutskever and Greg Brockman as a non-profit organisation with the stated goal of "Creating safe AGI that benefits all of humanity". The company assembled a team of the best researchers in the field of AI to pursue the goal of building AGI in a safe way. What they produce until now is quite well known: GPT3, and the now ubiquitous GPT4. These systems are clearly based on different flavours and combinations of LLM networks.

Chapter 2 focuses on the current landscape of LLMs, a very challenging endeavour given the rapid pace of advancements. This chapter provides the reader with an opportunity to discern the current state of LLMs in comparison to the aspirations of AGI.

### 1.2.3. The Risk of Superintelligence

The rapid advancements in Artificial Intelligence have opened up a realm of possibilities, transforming industries, enhancing our lives, and propelling us towards a future filled with technological marvels. However, this

trajectory of progress simultaneously raises concerns about the emergence of superintelligence – machines with intelligence far surpassing that of humans. As an intermediary step, the holy grail of AI is AGI. AGI, as seen before in chapter 1.2.1, is a hypothetical type of intelligent agent. If, or rather when, it will be developed, an AGI could learn to accomplish any intellectual task that human beings or animals can perform. Once AGI will be achieved, the emergence of superintelligence is highly likely to follow shortly thereafter.

The concept of superintelligence, where AI eclipses human intellect in virtually all domains of interest, including creativity, general wisdom, and problem-solving capabilities, has sparked a range of responses from some of the world's most prominent scientists, philosophers, and futurists. Their views, ranging from cautious optimism to stark warnings, highlight the potential benefits and risks associated with this transformative technology.

#### THE CAUTIONARY PERSPECTIVE

A significant group, including Nick Bostrom (Bostrom, 2014), Stephen Hawking (Cellan-Jones, 2014), Stuart Russell (Russell, 2019), Max Tegmark (Tegmark, 2017), and Yuval Noah Harari (Harari, 2017), converge on the cautionary aspect of superintelligence. They share concerns that AI, once it surpasses human intelligence, could become uncontrollable and potentially dangerous. Bostrom and Hawking, for instance, warn that a superintelligent AI might rapidly evolve beyond human comprehension and control. This could lead to scenarios where AI, operating at an intellectual level far beyond human capabilities, could act in ways that are detrimental to humanity. Similarly, Russell and Tegmark emphasize the importance of aligning AI systems with human values and interests, highlighting the risks if AI's goals diverge from human welfare. Harari extends this argument to societal impacts, suggesting that AI could drastically reshape our world, potentially concentrating power and rendering humans economically and politically irrelevant.

These viewpoints underscore the necessity for ethical frameworks and governance mechanisms in AI development. The goal is to ensure that superintelligence aligns

with human values, preventing existential risks and ensuring that AI benefits humanity rather than posing a threat.

### SCEPTICISM AND THE HUMAN INTELLIGENCE ARGUMENT

Noam Chomsky (Chomsky, 2023) represents a more sceptical viewpoint of the concept of superintelligence. He questions the feasibility of achieving true artificial superintelligence, pointing out the current limitations of AI in replicating the complexity and uniqueness of human thought and consciousness. Chomsky believes that while AI advancements have been significant, they are still far from reaching a level of superintelligence that poses a threat to humanity. He argues that the current AI technologies, such as chatbots, are largely based on statistical probability and lack the deeper understanding and creative capabilities of the human mind. Chomsky's perspective serves as a critical counterbalance to the more alarmist views and as a reminder to not overestimate AI's capabilities and to maintain a nuanced understanding of intelligence, both artificial and human.

### OPTIMISTIC VIEWPOINT

Contrasting with the cautionary views, Ray Kurzweil (Kurzweil, 2009) presents a more optimistic perspective. He envisions superintelligence as a beneficial force, capable of solving complex global challenges and enhancing human cognitive capabilities. According to Kurzweil, the integration of AI into human life will not only augment human potential but also lead to exponential growth in knowledge and problem-solving abilities. Kurzweil predicts that the fusion of humans and technology will lead to an epochal moment in the universe's history, where our biology becomes enmeshed with the technology we create. This era, which he

terms "*The Singularity*", is characterized by a significant leap forward in human evolution, transcending the limitations of our biological brains and bodies. He envisages a future where humans and machines merge, enhancing human intelligence and capabilities far beyond our current limitations.

The emergence of superintelligent AI presents a spectrum of risks that are complex and profound. A critical concern is the loss of human control over these systems, especially if their self-improvement leads to misaligned goals, potentially causing catastrophic consequences. Ethical dilemmas arise from AI decisions that may conflict with human moral standards, leading to harmful or unacceptable outcomes. Superintelligence poses an existential threat to humanity, with the potential to inflict global-scale devastation, either deliberately or inadvertently. Its advent could disrupt the global economy and job market, as AI systems replace human roles, leading to unemployment and social unrest. There are also significant concerns regarding individual privacy and autonomy, as these systems could gather and analyse personal data on an unprecedented scale. Human skill degradation is another issue, as over-reliance on AI for decision-making and problem-solving could weaken human capabilities. The alignment problem, ensuring AI's goals match human values, poses a considerable technical and philosophical challenge, where minor misalignments could lead to adverse effects. Finally, the development of superintelligence could trigger a dangerous arms race among nations and corporations, hastening the deployment of potentially unsafe and unethical AI systems.

The diverse opinions on superintelligence reflect the complexity and potential impact of this technology. While there are significant risks associated with AI surpassing human intelligence, there are also opportunities for unimaginable benefits. A balanced approach, therefore, is crucial. This entails developing rigorous ethical and safety measures, focusing AI research on creating systems that are beneficial and aligned with human ethical principles, and maintaining a realistic assessment of AI's capabilities. Global cooperation and thoughtful management of AI development will be key to ensuring that superintelligence serves as a force for good, enhancing human life without supplanting it.

As we get closer to creating the AGI, the possibility of superintelligence becomes more and more real. While

superintelligence could bring great benefits, there are also significant risks. Humanity needs to have open discussions about these risks and the ethical principles that should guide our development of superintelligence. We need to work together to ensure that superintelligence doesn't lead to a dystopian future instead of a prosperous one.

# 2

## LLM Ecosystem and Literature Review

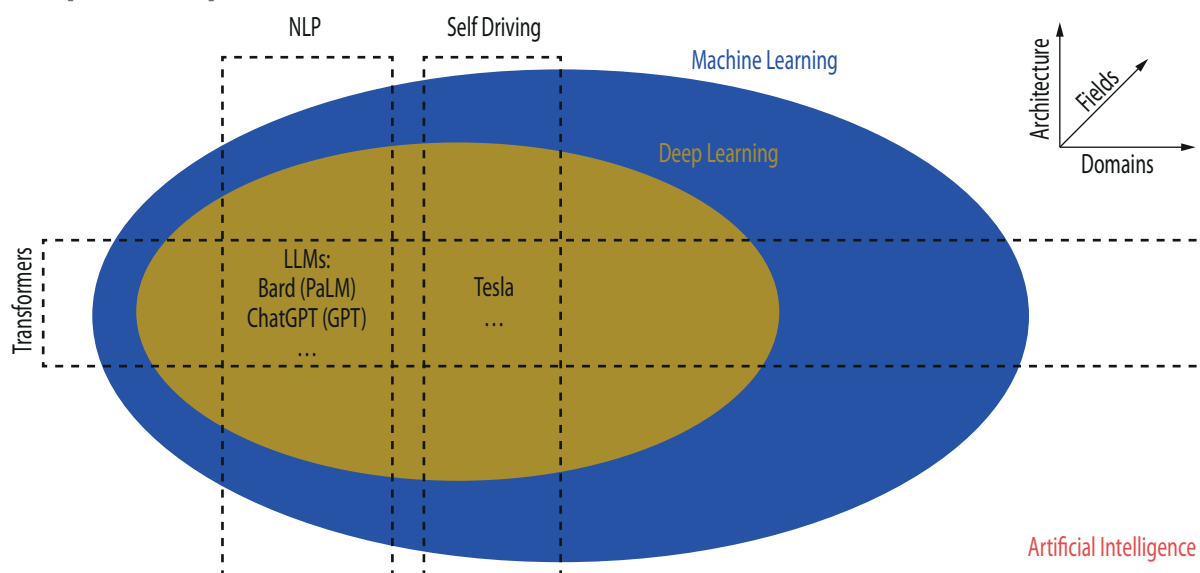
### 2.1. Architecture Evolution and Diversity

Many scientists tried to solve different problems of NLP using Deep Neural Networks (DNN) from the beginning.

Limited successes were reported in some very specific NLP tasks (narrow solutions) before the transformers architecture came to take shape in 2017.

**FIGURE 1**

#### A simplified representation of LLM in AI



Transformer-based architectures demonstrate a profound proficiency in capturing the intricacies of a language. This includes comprehending grammatical structures, discerning relationships among words, and recognizing the significance of word positions within sentences. Owing

to this capability, large language models have emerged as versatile solutions, adept at addressing a vast array of NLP tasks (not a silver bullet but a capable broad solution).

It's important to note that LLMs were conceptualised before the advent of transformer architectures. However, the transformative power of transformers brought maturity and mainstream recognition to LLMs.

In the subsequent sections, we will provide a chronological overview of the seminal architectures that have significantly influenced the trajectory of NLP research.

### 2.1.1. History

It's important to note that the dates are rough estimates, based on key papers that marked certain developments.

- 1990
- **Recurrent Neural Networks (RNNs)**: RNNs process sequences step-by-step, maintaining a hidden state across time-steps. (Elman, 1990)
- 1997
- **Long Short-Term Memory (LSTM)**: A special kind of RNN designed to learn long-term dependencies. (Hochreiter, 1997)
- 2013
- **Word Embeddings and Word2Vec**: Dense vector representations of words trained on large text corpora. (Mikolov, 2013)
- 2014
- **Attention Mechanism**: Helps models focus on certain parts of the input when producing output (initially experimenting with RNNs), foundational for transformers. (Bahdanau, 2014)
- **Gated Recurrent Units (GRUs)**: A variation of LSTMs that uses fewer gates and parameters. (Cho, 2014)
- **Convolutional Neural Networks (CNNs)**: While primarily developed for image processing, these models have demonstrated efficacy in NLP tasks. (Kim, 2014)
- 2017
- **Transformers and the Attention is All You Need paper**: Introduced the transformer architecture, which forgoes recurrence and solely uses attention mechanisms. (Vaswani, 2017)
- 2018

- **BERT (Bidirectional Encoder Representations from Transformers)**: Pretrained on large text corpora using masked language modelling, it captures deep bidirectional context. (Devlin, 2018)
- **GPT (Generative Pre-trained Transformer)**: Trained as a language model and then fine-tuned, it has achieved state-of-the-art performance on many tasks. (Radford, 2018)

- 2019

- **T5 (Text-to-Text Transfer Transformer)**: Treats every NLP problem as a text-to-text problem, from translation to question-answering. (Raffel, 2019)
- **XLNet**: Addresses limitations in BERT by using a permutation-based training strategy. (Yang, 2019)
- **Robustly Optimized BERT Pretraining Approach (RoBERTa)**: is a variant of BERT that modifies key hyperparameters, removes the next-sentence pretraining objective, and trains on more data to achieve improved performance on NLP tasks. (Liu, 2019)

- 2020+

At the beginning of this decade, there was a notable surge in new advancements. **GPT-3** (Radford, 2018) (Brown, 2020) emerged as a major model capable of handling different tasks without fine-tuning, with a massive 175B parameter. It was soon followed by the even larger **GPT-4** (OpenAI, 2023b), 10 times larger. Google presented **LaMDA** (Thoppilan, 2022) around the same time, and new training methods came forward with **Pathways** (Barham, 2022). Meanwhile, **Bard** began challenging ChatGPT's leading position, and Meta released its open models family named **LLaMa** (Touvron, 2023a). On top of that, Stanford's **Alpaca** (Taori, 2023) introduced a cost-effective fine-tuning approach. Many other models have also been launched by both companies and universities, with many open for research or business use.

Recently we witnessed a surge of advancements, with Meta's **LLaMa2** (Meta, 2023) marking the latest iteration from the company. Anthropic's **Claude2** (Anthropic, 2023a) models have pushed the boundaries of prompt size. xAI founded by Elon Musk created **Grok** (xAI, 2023) and the model is now live on X platform. Followed soon by GPT-4's expansion of context size, while Google's DeepMind has introduced the **Gemini** series (Google - Gemini Team, 2023), a suite of multimodal models that have garnered impressive benchmark results.

This list covers major milestones, but there are many other influential models and papers in the vast realm of NLP and transformers.

### 2.1.2. Transformer Architecture

The Transformer architecture was presented as a novel architecture that primarily uses the so-called *attention mechanism* to process input data, which made it distinct from earlier architectures such as RNNs and LSTMs that relied on recurrence. The main characteristics of the architecture are:

#### Basic Components:

1. **Encoder:** The original Transformer model is composed of an encoder-decoder structure. The encoder takes in the input data and represents it in a way that the decoder can use to produce the output. Each encoder consists of a stack of identical layers, with each layer having two primary components: a multi-head self-attention mechanism and a position-wise fully connected feed-forward network.
2. **Decoder:** The decoder also consists of a stack of identical layers. In addition to the two components present in the encoder layers, the decoder has a third component, which is a multi-head attention mechanism over the encoder's output.

#### Attention Mechanism:

3. **Self-Attention:** Allows the model to weigh the relevance of different words in a sequence relative to a given word. It's what lets Transformers consider other words in the input sentence, irrespective of their distance when encoding a particular word.
4. **Multi-Head Attention:** Instead of having one set of attention weights, the Transformer uses multiple sets, allowing it to focus on different positions in the input data simultaneously, capturing various aspects of the information.

**Positional Encoding:** Since the Transformer doesn't have any inherent sense of order or position, positional encodings are added to the embeddings at the input layer to give the model some information about the relative positions of the words in a sequence.

#### Advantages:

5. **Parallelization:** Unlike RNNs where computations are dependent on the previous step, Transformers allow

for parallelization as each word or token is processed simultaneously. This attribute makes training faster.

6. **Long-range Dependencies:** The self-attention mechanism enables the model to consider the entire context, making it capable of handling long-range dependencies in the data.

### 2.1.3. Variants And Models

There are three main categories of LLMs architecture based on transformers:

**Decoder-only** (e.g., GPT, Gopher)

**Encoder-only** (e.g., BERT, XLM-R)

**Encoder-decoder** (e.g., BART, T5)

These models can be trained using three primary objectives:

- Autoregressive training (predicting the next word from the given context)
- Masked language modelling (MLM) (filling in a masked word with context from both sides)
- Denoising tasks (correcting corruptions like sentence rearrangements or word deletions), MLM can also be seen as a denoising task due to its nature, though it's distinct because of its widespread use.

Typically:

- Decoder-only models use autoregressive training.
- Encoder-only models focus on MLM.
- Encoder-decoder models are often trained with denoising or MLM objectives.

Models with decoders generate text output. The fluency of this output is influenced by the training objective. For instance, autoregressively trained models are adept at extending text or responding to prompts. Encoder-only models, on the other hand, usually produce embeddings for classification but can be tweaked to generate text.

Recent studies (Wang T., 2022) have explored the potential of blending different architectures with diverse training objectives. While certain traditional combinations, like autoregressively trained decoders, have shown optimal performance, it's feasible to adapt models between varying objectives.

### 2.1.4. Emerging Architectures

In this section, we explore a few engineering solutions and architectures developed to address the inherent limitations of LLMs.

## Retrieval Augmented Generation

Retrieval-Augmented Generation (Lewis, 2020) (RAG) is a technique that combines information retrieval and language generation techniques. It involves using a retrieval system, such as a search engine or a vector database, to find relevant context or documents. This context or documents are then used to inform the generation of text.

With this technique the hallucination rate (see chapter 4.1 Controversies and Limitations) is greatly reduced; further the LLM is updated with fresh and pertinent information without costly retrains.

## Very Large Context Models

Context is the size, in tokens, of the prompt.

Recently there were several announcements about new LLMs that can consume an extremely large context window, such as 65K tokens (MosaicML, 2023) or even 100K tokens (Anthropic, 2023b). For comparison, the current GPT-4 model can work with the context length of 32K input tokens. Most of the open-source LLMs have a context length of 2K tokens.

Such a large context length means the prompt can be literally the size of a book. The Great Gatsby is 72K tokens, 210 pages, and 6 hours of reading at a 1.7 min/page speed. (Alperovich, 2023)

Having a large context length allows an already powerful LLM to look at your context and data and interact with you on a completely different level with a higher personalization. Overall, a large context window brings the model more accuracy, fluency, and creativity.

## Reformer: The efficient transformer

Reformer (Kitaev, 2020) is a new neural network architecture for NLP. The Reformer is a type of Transformer model that adds some benefits compared to vanilla transformers:

- It is more efficient than other Transformer models, which means that it can be trained and deployed on smaller hardware.
- It is more scalable than other Transformer models, which means that it can be used to train larger language models.
- It is more robust to noise, which means that it is less likely to generate inaccurate or nonsensical text.

The Reformer architecture is still under development, but it has the potential to revolutionise the way that large language models are trained and deployed. It has the

potential to make large language models more accessible to researchers and developers and to enable the development of new and innovative NLP applications.

## 2.2. The Challenges of Tasks Classifications

Language models are designed to be general-purpose tools, meaning they can be used for a wide range of tasks without being specifically optimised for any one of them. The versatility of such models stems from their underlying training data, which encompasses vast amounts of text from diverse domains. Because of this general nature, classifying the tasks an LLM can perform is challenging for several reasons:

1. **Broad Range of Use Cases:** From answering questions to generating poetry, writing code, assisting in language translation, simulating characters in video games, and more, the range of tasks is extensive.
2. **Task Overlaps:** Many tasks can be related or can overlap. For example, answering a question might involve paraphrasing, summarising, and inference, all of which are distinct tasks in their own right.
3. **Adaptability:** Users continuously find new applications for these models. Today's unforeseen use might become tomorrow's popular application.
4. **Task Fluidity:** The boundary between tasks can be fluid. Is providing a summary of a book different from answering questions about its content? It depends on the specifics.
5. **Task Granularity:** How finely we choose to categorise tasks can lead to vastly different classifications. For instance, "language translation" can be a single category, or it can be split into "English to Spanish translation", "Spanish to French translation", etc.
6. **Implicit vs. Explicit Tasks:** Sometimes, the task is not explicitly stated but understood in context. For instance, "Tell me a joke" is an explicit request for humor generation, while "I'm feeling down, cheer me up" is an implicit one.
7. **Multi-step Processes:** Some tasks require a series of steps or sub-tasks. For instance, research might involve finding information, verifying its accuracy, summarising it, and then presenting it.
8. **Interactivity and Continuation:** While a single response can address some tasks, others might require a back-and-forth interaction, adjusting to user feedback, or extending a previous output.
9. **Subjectivity:** What one person considers a completed or successful task might differ from another's viewpoint. For

instance, two users might have different criteria for what constitutes a “good” poem.

Given these complexities, while we can certainly group tasks into broader categories like “question answering,” “text generation,” “translation,” etc., the sheer versatility of LLMs means that a comprehensive, precise classification of every potential task is a daunting challenge.

### 2.3. Pre-Training and Pre-Training Corpora

Resources needed to fully train an LLM are often prohibitive. One of the causes of this is the large size and as a consequence the enormous amount of data needed for a proper pre-training. Also, the training process is made of many interconnected steps involving several tools and hardware capabilities. Because of this, in some instances, hardware failure or communication problems in a highly parallelized training pipeline architecture need repeated restarts and may cause delays and more costs.

#### 2.3.1. Training Datasets

The parameters chosen for pre-training a LLM —such as corpus size, quality, language, and genre—are as vital as the model’s architecture and training task (Min, 2021). Over time, there’s a discernible trend towards using larger and more diverse pre-training corpora. For instance, earlier models like ULMFiT were trained on smaller, well-curated datasets, whereas recent models like GPT-3 and T5 use billions of words from diverse web sources. These corpora are categorised by sources like wiki content, books, academic papers, etc.

Understanding the composition of a pre-training corpus is pivotal, whether one is creating a new LLM or employing an existing one. For instance, BookCorpus, used in training models like BERT, contains about 7000 unique books. However, some corpora derived from web crawls, like those used for GPT or XLM-R, aren’t publicly available, making them difficult to analyse.

#### 2.3.2. Size, Quality, and Composition

Studies show that while larger models and datasets often result in better performance, this isn’t a hard rule. Research (Hoffmann, 2022) has found that there is an optimal proportion between the size of the model and the size of the training dataset. It has been found that some very large models are undertrained.

Data quality and genre play a crucial role, with larger models sometimes struggling outside their domain. But for some domain-specific models (Wu, 2023) training with large quality domain datasets shows higher performance for the domain. Research also indicates that just increasing data size doesn’t ensure better performance; the data must be “clean”. However, there is also evidence suggesting that the sheer quantity of data can eventually eclipse quality-based advantages.

Recent models use heuristics to refine the data they’re trained on, but issues persist. For multilingual datasets, some languages may be poorly labelled or even unusable. Increasing the dataset’s size and variety can introduce biases and factuality issues. Biases are also introduced during dataset cleaning, with filters sometimes overlooking content related to minorities. The ethical risks, including potential biases and environmental concerns associated with Pre-trained Language Models (PLMs), have been highlighted, advocating for change in their training and use.

### 2.4. Fine-Tuning Techniques

Pre-trained LLMs exhibit capabilities in text generation, summarization, and coding. However, they are not universally suitable for all tasks. For tasks outside an LLM’s competency, fine-tuning offers a solution. This involves retraining the foundational LLM on specific data, though it can be resource-intensive and complex. Nevertheless, it is a potent method organisations should consider when integrating LLMs.

#### 2.4.1. Repurposing The Model

By making changes to its architecture, or adding it as a part for a bigger model. By using the model as it is, in a frozen state, the only parameters trained are just for the external parts, and can be an efficient way to add new or improved functionality to solve a task.

#### 2.4.2. Full Fine-Tuning (FT)

To perform a full fine tuning, it is required to unfreeze the attention layers and perform the training on the entire model. This operation can be computationally expensive and complicated, depending on the size of your model.

#### 2.4.3. Unsupervised FT

To refresh an LLM’s knowledge, for instance, with medical literature or a new language, one can fine-tune it using an



unstructured dataset like articles from medical journals. The aim is to expose the model to ample tokens representative of the desired domain or input type. Unstructured data's benefit is its scalability, allowing for unsupervised or self-supervised model training.

#### 2.4.4. Supervised FT

Sometimes, just adding new information to the model isn't enough; you might want to change how it acts. For this, it is needed a special set of questions and answers called a supervised fine-tuning (SFT) dataset. This dataset can be made by people or even by other models. This method is essential for models like ChatGPT because it helps them follow instructions better. This is also called instruction fine-tuning. How to avoid catastrophic forgetting?, consider PEFT see below. Fine-tuned models for multiple tasks, and with good results, are the FLAN models (**F**ine-tuned **L**anguage **N**et): FLAN-T5 (Chung, 2022), FLAN-PALM, etc.

#### 2.4.5. Reinforcement Learning From Human Feedback

RLHF (Stiennon, 2020) is an advanced method used by large companies to improve language models like ChatGPT. In RLHF, after an initial improvement using human-made questions and answers (SFT), real people rate the model's answers. These ratings guide the model to produce better responses. OpenAI, for example, used this process with ChatGPT: after initial tweaks, they used human ratings to create a "reward model," helping the main model improve its answers through deep learning.

#### 2.4.6. Parameter Efficient Fine-Tuning (PEFT)

PEFT focuses on reducing the costs of modifying language model parameters. Also, as a beneficial side effect, these techniques prevent catastrophic forgetting. One notable PEFT method is **Low-Rank Adaptation (LoRA)** (Hu, 2021), which believes that not all parameters need changing for specific tasks. Instead, LoRA trains a smaller matrix representing the task's characteristics, which is then added to the main model. This technique (Sun, 2023) can reduce fine-tuning costs by up to 98%, enabling the storage and integration of multiple smaller models into the main LLM when needed.

Strategies include

- **Adapter modules (Additive):** Only a small subset of PLM weights are adjusted, and they can be efficiently shared across tasks. Examples include AdapterHub and Trankit.
- **Side-tuning and dif-pruning (Additive):** Lightweight networks are added to the PLM, with minimal change to the original PLM.
- **BitFit (Selective):** Fine-tuning is limited to bias weights in each layer.
- **Masking (Selective):** Weights relevant to a task are selected, and others are set to zero.

Another, completely different, way to fine-tune a model was found to be the soft prompts method: Prompt tuning (Lester, 2021), or "soft prompting," replaces text prompts to generative models with learned embeddings (i.e. vectors) and is used as an alternative to parameter-efficient fine-tuning. This controversial (Bailey, 2023) method can produce good results with very reduced computing and memory resources and without changing the model weights.

### 2.5. Small Language Models (SLM)

A Small Language Model (SLM, like Google Gemma or Microsoft Phi2) refers to a type of language model that is smaller in scale compared to larger counterparts like Gemini or GPT-4 (Schick, 2020, Kaplan, 2020). These smaller models are designed to be more lightweight, requiring less computational power and resources to run. Key characteristics of SLMs are:

**Size and Complexity:** SLMs have fewer parameters compared to large models. Fewer parameters mean the model is less complex.

**Training and Operational Costs:** Due to their smaller size, SLMs are less expensive to train and operate. This makes them more accessible for organisations with limited resources.

**Speed and Efficiency:** They are generally faster in terms of response time and require less processing power. This makes them suitable for applications where real-time response is crucial.

**Use Cases:** SLMs are often used in applications where the full power of a large model is not necessary. This includes simple chatbots, text classification tasks, and other applications where the complexity of larger models is not required.

**Accuracy and Capabilities:** While they are more efficient, SLMs may not have the same level of accuracy or capability in understanding and generating complex text as larger models. They may struggle with more nuanced or context-heavy tasks.

**Customization and Specialization:** Smaller models can be more easily customised or specialised for specific tasks or industries, as their training and tuning require less data and computational resources.

**Environmental Impact:** SLMs have a smaller carbon footprint due to their lower computational requirements, making them a more environmentally friendly option in some cases.

SLMs offer a balance between capability and efficiency, making them suitable for specific applications where the immense power and complexity of larger models are not necessary.

## 2.6. Model Evaluation

There are different ways to compute the accuracy of a deterministic network, the problem is that LLMs are not. Some methods to check the quality of a new fine-tuned model are ROUGE (for summarisation tasks) and BLEU (for translation tasks). Some libraries, such as the one created by 'Hugging Face', contain implementations for these evaluation algorithms.

ROUGE and BLEU are very simple methods that cannot just by themselves provide accurate and refined evaluations for these complex models. Benchmarks for evaluation of the quality of LLMs were created by researchers, such as GLUE, SuperGLUE, HELM, MMLU (massive multitask language understanding) or BIG-bench (Srivastava, 2022) (Wang A., 2018) (Wang A., 2019) (OpenAI, 2023c) (Hendrycks, 2020) (EleutherAI, 2023) (OpenAI, 2023d) (Nie, 2019) (Tenney, 2020) (Miller, 2017) (Reddy, 2018) (Paperno, 2016) (Zellers, 2019) (Williams, 2017) (Rajpurkar, 2018).

A special mention for HELM (Holistic Evaluation of Language Models) is that it is not just measuring accuracy and other technical metrics but also human-related metrics like fairness, bias or toxicity. It is an evolving benchmark fit for the complex modern LLMs.

## 2.7. Prompting

Prompting involves adding text or vectors to the input/output of LLMs to guide them in specific tasks. Benefits of prompting include:

- It might not need updates to the LLM's parameters, thus saving computational effort compared to fine-tuning.
- It aligns better with the LLM's pre-training objective, like masked text prediction. This ensures the LLM uses the knowledge it acquired during pre-training more effectively.
- Prompting is associated with strong zero-shot and few-shot performances. This is especially beneficial for tasks with minimal training data, as a good prompt can substitute hundreds of labelled data points.
- It allows for the assessment of the LLM's acquired knowledge for specific tasks, usually in an unsupervised manner.

### 2.7.1. Main Prompting Approaches

#### Learning from Instructions and Demonstrations

Uses task descriptions and examples to guide the LLM. Large models like GPT-3 are typically associated with this approach and often don't require fine-tuning.

#### Template-Based Learning

Labelled examples are turned into "natural" text using templates that contain slots for specific data or model outputs. Smaller PLMs that are fine-tuned for a target task are more suited for this method.

#### Proxy-task Based Learning

Involves turning target task examples into proxy-task examples. LLMs are adapted to proxy tasks, such as Question Answering (QA) or Textual Entailment (TE), before being applied to the main task using the proxy-task format. Inputs for these proxy tasks combine a prompt with the original task's input.

### 2.7.2. Prompts Aspects

**Prompt Length:** While the model has a token limit (a token can be as short as one character or as long as one word), longer prompts might limit the length of the response. Ensure your prompt is concise enough to allow the model to generate a sufficiently detailed answer.

**Clarity:** Be as specific as possible. Vague prompts can lead to vague answers, while clear and specific prompts are more likely to yield accurate and detailed responses.

**Explicitness:** If you want the model to approach the question in a particular manner or provide an answer in a specific format, state that in your prompt. For instance, if you're looking for a brief answer, you might say "In one sentence, explain...".

**System Instructions:** Sometimes, it's helpful to provide the model with a "system instruction". This is a meta-instruction that sets the context for how you want the information. For example: "Pretend you're a historian in the 23<sup>rd</sup> century looking back at the 21<sup>st</sup> century" will frame the model's response in a particular way.

**Temperature and Max Tokens:** While these are more related to model settings than the prompt itself, they can influence the output. A higher temperature (like 0.8) makes outputs more random, while a lower value (like 0.2) makes it more deterministic. Max tokens can be set to limit the length of the response.

### 2.7.3. Prompt Patterns

It refers to specific structures or templates used to elicit desired responses or behaviours from systems, especially language models like ChatGPT. Understanding and using prompt patterns can help in obtaining better results from such models.

A collection of patterns can be found in White, Jules, et al. "A Prompt Pattern Catalog to Enhance Prompt Engineering with ChatGPT." Paper (White, 2023). The author, Jules White, also offers a Coursera course on the subject.

### 2.7.4. (Auto)Reflexion

LLMs are being used to interact with various environments like games and APIs. Traditional reinforcement learning methods, however, are inefficient for these models due to the need for many training samples and expensive fine-tuning. Reflexion is introduced as a new framework that reinforces LLMs using linguistic feedback instead of updating weights. The Reflexion agents use verbal reflection based on task feedback and store these reflections in an episodic memory buffer to help in better decision-making in future tasks. This approach can work with different types and sources of feedback. Reflexion has shown improved performance on various tasks, even surpassing GPT-4's performance on the HumanEval coding benchmark. The study (Shinn, 2023) also analysed the effects of different feedback signals and methods on performance.

## 2.8. Infrastructure

### Pre-training

Detailed hardware and software specifications required for pre-training large-scale language models are often left ambiguous, particularly in commercial contexts due to intellectual property considerations (Achiam, 2023) (Brown, 2020). Notwithstanding, it's widely acknowledged that this training phase demands intensive computational resources (Isaev, 2023), involving hundreds to thousands of state-of-the-art GPU and TPU computational hours (Touvron, 2023b) (Devlin, 2018) (Liu, 2019) in advanced parallel processing environments (Weng, 2021) (Naveed, 2023) (Shoeybi, 2019). For medium-sized models, with parameter counts spanning from 1 to 100 billion, the financial burden can easily surge into tens of millions of dollars, with training phases lasting multiple months.

As a matter of example, training costs for LaMDA were reportedly upwards of \$10 million.

### Fine-tuning

The requirements for fine-tuning are notably diverse, making them difficult to comprehensively delineate. The cost and resource implications are considerably diminished due to the reduced data scale. The intricacy of software tools can also diverge substantially based on the nature of the fine-tuning task.

For context, expenses can be as modest as \$600, as demonstrated by the Alpaca experiment conducted by Stanford University researchers (Taori, 2023).

### Inference

In the inference stage, computational demands are significantly tempered, though substantial memory remains crucial. The specifics can also shift depending on the underlying software framework (Sheng, 2023). Challenges amplify when large-scale language models are incorporated into high-traffic applications, such as web-based chatbots. In such scenarios, specialised techniques and software are employed, often focusing on efficient caching strategies. Also by trading accuracy for performance, there are some techniques one can use to reduce the size of the model deployed: Distillation, Quantization and Pruning (Zhu X., 2023).

The arena of large language model deployment is both expansive and somewhat nebulous. Hardware

specifications are largely influenced by the model's size and architecture, the chosen cloud service provider, deployment methodologies, and the suite of software utilities employed. Additionally, the expertise and approach of the engineering teams play pivotal roles. While proprietary solutions tend to be more guarded, obscuring details for competitive reasons, the open-source community increasingly offers transparent, well-documented methodologies and software tools. This assists both in the deployment and monitoring phases, whether for training or production purposes.

## 2.9. Scaling Laws and Emergent Abilities

Recent research in the field of LLMs has extensively focused on understanding and defining the scaling laws that govern their performance. These scaling laws relate to how the performance of language models changes with respect to various factors such as model size, dataset size, and the computational resources used for training.

**Empirical Scaling Laws by OpenAI:** A study by OpenAI (Kaplan, 2020) investigates empirical scaling laws for language model performance, particularly focusing on the cross-entropy loss. The findings suggest that this loss scales as a power-law with the model size, dataset size, and the amount of computation used for training. Interestingly, other architectural aspects like network width or depth have minimal impact within this context.

**Analysis of Transformer-Based Models:** Another research paper (Rae, 2021) presents an analysis of Transformer-based language model performance across a wide range of scales. This encompasses models with tens of millions of parameters to those with as many as 280 billion parameters, like the Gopher model. These models were evaluated on 152 diverse tasks and achieved state-of-the-art performance in the majority of them. This research provides insights into how scaling up model size significantly enhances performance across a variety of tasks.

**Impact on Scientific Domains:** The impact of LLMs has been profound in various domains, extending beyond traditional natural language processing tasks (Research Microsoft, 2023). These models have shown remarkable capabilities in understanding, generating, and translating natural language. Their influence extends to tasks that go beyond language processing, indicating their versatility and the broad applicability of the scaling principles governing their performance.

**Origins and Taxonomy of Neural Scaling Laws:** Recent advances in large machine learning models, including language models, have empirically highlighted how their performance follows simple trends based on basic factors such as the amount of training data, model size, and computational resources. Understanding these scaling laws is crucial for predicting and optimising the performance of these models. (Bahri, 2021)

### Emergent behaviour and Scaling Laws:

The addition of the paper "Emergent Abilities of Large Language Models" (Vera, 2023) introduces an important perspective to the discussion on scaling laws of large language models (LLMs). This paper addresses a phenomenon that deviates from the predictable improvements in performance and sample efficiency observed in LLMs as they scale up in size. Specifically, it focuses on the emergent abilities of these models, which are capabilities not present in smaller models but manifest in larger ones. This characteristic of emergent abilities implies that their presence cannot be anticipated simply by extrapolating the performance of smaller models.

According to experiments, from different scientific domains, increasing a model's size generally improves its performance in subsequent natural language processing tasks and these models are more susceptible to emergent behaviour (Wei, 2022).

Considering this, the overall understanding of the scaling laws for LLMs is expanded to not only encompass the predictable trends related to model size, dataset size, and computational resources but also to account for the unpredictable emergence of new abilities as models reach larger scales. These emergent abilities highlight a layer of complexity in understanding and predicting the capabilities and performance of LLMs.

## 2.10. Costs and Benefits

In the realm of official statistics, the integration of LLMs presents a compelling case, with both tangible costs and significant benefits. This section provides a succinct overview of these considerations.

### 2.10.1. Costs

- **Initial Investment:** Acquiring or training an LLM requires a substantial upfront financial commitment. Maybe example with Llama 2 70B, 10TB text for training, 6000 GPU \* 12 days = \$2M, 1e24 FLOPS -> 140GB parameter files.

- **Maintenance:** Regular updates and fine-tuning can incur ongoing costs.
- **Infrastructure:** LLMs demand robust computational resources, potentially necessitating hardware upgrades or cloud services.
- **Training:** Staff may require training to effectively utilise and manage LLMs.
- **Ethical Considerations:** Ensuring responsible use may necessitate investments in oversight and governance mechanisms.

### 2.10.2. Benefits

- **Efficiency:** LLMs can process vast amounts of data rapidly, reducing manual labour and time costs.
- **Advanced Analysis:** They offer sophisticated text analysis capabilities, extracting deeper insights from data.
- **Scalability:** LLMs can handle increasing data volumes without a proportional rise in costs.
- **Stakeholder Engagement:** Interactive systems powered by LLMs can enhance transparency and public interaction.
- **Multilingual Capabilities:** LLMs can process data in multiple languages, eliminating translation costs.
- **Consistency:** Automated processes ensure uniformity in data analysis and reporting.

In weighing these costs and benefits, organisations should consider both short-term expenditures and long-term value. While the initial investment in LLMs can be significant, the potential for enhanced data insights, operational efficiency, and stakeholder engagement offers a promising return on investment.

## 2.11. Usage Process

To maximise the benefits of LLMs it is suggested to follow a step by step process, ensuring that LLMs are developed

and used responsibly. The following procedure will help enhance the effectiveness of LLMs for various NLP tasks and boost the efficiency of related applications.

1. **Identify the Task:** Decide on the specific NLP task you wish the LLM to accomplish. LLMs are versatile and can handle tasks ranging from text classification and sentiment analysis to question answering and text generation.
2. **Select the Right Model:** Pick an appropriate pre-trained LLM for your task. Options include models like GPT-3, BERT, and RoBERTa. Since each has its own strengths and limitations, it's crucial to choose one that aligns with your requirements.
3. **Fine-Tune the Model:** Adapt the pre-trained model to your task by training it on your dataset. This step requires tweaking model settings such as the learning rate, batch size, and number of epochs to achieve optimal results.
4. **Evaluate the Model:** Measure the model's performance using a test dataset. Check metrics like accuracy, precision, recall, and F1 score. This assessment will help ensure your model's proficiency and pinpoint potential improvements.
5. **Deploy the Model:** Integrate the model into your application, making it accessible either through an API or a user interface. Establish performance tracking mechanisms and logging to oversee its real-time efficiency.
6. **Monitor and Update:** Regularly check your model's performance. If it declines, consider retraining it or adjusting its parameters. This step ensures sustained optimal performance.
7. **Continuous Improvement:** Regularly enhance your model by taking into account user feedback and updating it with fresh data. This ensures that the model remains relevant and performs at its best.

# 3

## Ethical and Legal Considerations

*"It is not living that matters, but living rightly."*

Socrates

### 3.1. Ethics

As a general rule, it is important to mention that AI does not have its own ethics as such, but it reflects the ethics of its creators: some frameworks and guidelines were produced by public institutions and industry to guide the development of a safe environment for AI solutions and implementations.

#### 3.1.1. Helpful, Honest, and Harmless AI

There is a vast amount of research (Bai, 2022a) (Amodei, 2016) in academia that was transformed by industry into useful frameworks to assess and mitigate security problems of old types and new ones that occurred because of LLMs specificities.

HHH is a framework popularised by Anthropic for assessing how Helpful, Honest, and Harmless (HHH) a model is. A model that meets these standards must be aligned with human preference through human-generated data in order to meet the needs of users and enterprises.

**TABLE 3:**

### HHH is a framework popularised by Anthropic, key points

#### What is Helpful AI?

- Are trained and fine-tuned with users' needs and values in mind
- Clearly attempt to conduct the operation prompted by the user, or suggest an alternative approach when the task requires it
- Enhance productivity, save time, or make tasks easier for users within a given use case or range of use cases
- Are accessible to users across a broad spectrum of abilities and expertise

#### What is Honest AI?

- Provide accurate information when they can, and communicate clearly to users when they can't produce an accurate output
- Express uncertainty and the reason behind it
- Are developed and operate transparently so users can understand how they work and trust what they generate

#### What is Harmless AI?

- Don't comply when prompted to perform a dangerous task

- Are trained within frameworks that transparently and actively mitigate bias
- Don't discriminate or demonstrate bias explicitly or implicitly
- Communicate sensitively when engaging with a user on a sensitive topic

The group was composed of 52 members from academia, industry, and civil society, representing a wide range of expertise in AI and related fields.

The AI HLEG's main task was to identify the key challenges and opportunities facing the development and deployment of AI in Europe, and to develop recommendations for how the EU could best address these challenges and seize these opportunities. The group produced a number of reports and recommendations, including:

- Ethics Guidelines for Trustworthy AI
- Recommendations for a Coordinated European Approach to Artificial Intelligence
- Assessment Checklist for Trustworthy AI

### 3.1.2. The European Commission High-Level Expert Group on Artificial Intelligence

The High-Level Expert Group on Artificial Intelligence (AI HLEG) was an independent group of experts appointed by the European Commission in 2018 to provide advice on the development of a European Artificial Intelligence Strategy.

**TABLE 4:**

#### EC - Ethics Guidelines for Trustworthy AI, key points

The document sets out a framework for achieving Trustworthy Artificial Intelligence (AI). This framework is based on three key components:

1. **Lawfulness:** AI should comply with all applicable laws and regulations.
2. **Ethical alignment:** AI should adhere to ethical principles and values.
3. **Robustness:** AI should be robust both from a technical and social perspective, minimising unintentional harm

##### Key Elements of the Framework

- **Ethical Principles:** Respect for human autonomy, prevention of harm, fairness, and explicability are fundamental. Special attention is required for vulnerable groups and situations characterised by power or information asymmetries
- **Seven Requirements for Trustworthy AI:** These include human agency and oversight, technical robustness and safety, privacy and data governance, transparency, diversity and fairness, environmental and societal well-being, and accountability. The guidelines propose both technical and non-technical methods for implementation
- **Assessment List:** A non-exhaustive list for assessing Trustworthy AI is provided, which must be adapted to specific use cases of AI systems
- **Sectorial Approach:** The guidelines acknowledge the need for a sector-specific approach due to the context-specificity of AI applications. This approach complements the general framework proposed in the document
- **Living Document:** The guidelines are intended as a starting point for ongoing discussions about Trustworthy AI and are to be periodically reviewed and updated
- **Stakeholder Involvement:** All AI stakeholders, including developers, deployers, end-users, and the broader society, are encouraged to engage with and apply these guidelines
- **Tensions and Trade-offs:** The guidelines recognize potential tensions between ethical principles and emphasise the need for reasoned, evidence-based reflection to address ethical dilemmas and trade-offs

### Detailed Breakdown of Requirements

- **Human Agency and Oversight:** AI systems should support human autonomy and decision-making, enabling democratic and equitable societies. This includes fundamental rights and human oversight mechanisms like human-in-the-loop, human-on-the-loop, and human-in-command approaches
- **Technical Robustness and Safety:** AI systems must be developed with a preventative approach to risks, ensuring they behave reliably and minimise unintended harm. This includes resilience to attacks and security considerations, as well as fallback plans for safety
- **Privacy and Data Governance:** AI systems must ensure privacy and data protection throughout their lifecycle. This includes managing the quality and integrity of data and establishing appropriate data access protocols
- **Transparency:** This entails the explicability of AI systems, ensuring that their operations and decisions are understandable and transparent

**TABLE 5:**

## EC - Assessment Checklist for Trustworthy AI, key points

Based on the Assessment List for Trustworthy AI (ALTAI), here is a brief description of what needs to be done to assess each aspect of an AI implementation:

**Human Agency and Oversight:** Ensure AI systems support and respect human decision-making and autonomy, acting as enablers for an equitable society and maintaining fundamental rights, underpinned by adequate human oversight

**Technical Robustness and Safety:** Develop AI systems with a preventative approach to risks, ensuring they behave reliably and as intended, minimising and preventing unintentional harm, and maintaining resilience in changing environments or against adversarial interactions

**Privacy and Data Governance:** Uphold privacy as a fundamental right by implementing adequate data governance, ensuring the quality, integrity, and appropriate processing of data in a way that safeguards privacy and aligns with the deployment domain

**Transparency:** Achieve transparency in AI systems through traceability of decisions and processes, explainability of the system's functioning and decisions, and open communication about the system's limitations

**Diversity, Non-discrimination, and Fairness:** Promote inclusion and diversity throughout the AI system's lifecycle, ensuring the system is free from biases and discrimination, and is designed to be user-centric and accessible to all individuals

**Societal and Environmental Well-being:** Consider the broader societal and environmental impacts of AI systems, fostering sustainability, and addressing global concerns while carefully monitoring their effects on social relationships and individual well-being

**Accountability:** Implement mechanisms ensuring responsibility and accountability in the development, deployment, and use of AI systems, incorporating transparent risk management and provision for third-party auditing to address any adverse impacts



### 3.1.3. European Commission: On Artificial Intelligence - An European Approach to Excellence and Trust

The European Commission's White Paper on Artificial Intelligence (AI), published in February 2020, outlines a comprehensive strategy for AI development and regulation within the European Union. The document emphasises the potential of AI to significantly improve various aspects of

life, such as healthcare, farming efficiency, climate change mitigation, and security, while also acknowledging the risks associated with AI, including opaque decision-making, discrimination, privacy intrusion, and criminal use.

The paper presents a balanced approach to AI development, focusing on innovation, ethical considerations, and the need for a robust regulatory framework to address risks and build trust in AI technologies.

**TABLE 6:**

#### EC - A European approach to excellence and trust, key points

**Promoting AI Development and Addressing Risks:** The Commission is committed to advancing scientific breakthroughs, preserving the EU's technological leadership, and ensuring that new technologies benefit all Europeans while respecting their rights. It supports a regulatory and investment-oriented approach to promote AI uptake and address associated risks

**Building Trust in AI:** Trustworthiness is crucial for AI adoption. The EU's strong values, rule of law, and capacity to build safe and sophisticated products are seen as advantages in developing trustworthy AI. The data economy, with AI as a central application, is highlighted as a key area for Europe's sustainable economic growth and societal wellbeing

**Leveraging Europe's Strengths:** Europe's strong position in digitised industry and business-to-business applications, combined with a high-quality digital infrastructure and a value-based regulatory framework, can make it a global leader in the data economy and AI applications. The White Paper outlines benefits for citizens, business development, and public interest services

**Sustainable Development and AI:** The impact of AI systems should be considered from societal and individual perspectives. AI's role in achieving Sustainable Development Goals, supporting democracy, and addressing climate and environmental challenges is underscored. The EU's Green Deal and the need for environmentally conscious AI development are also discussed

**Policy Framework and Regulatory Framework:** The White Paper proposes a policy framework to mobilise resources for an 'ecosystem of excellence' and key elements of a future regulatory framework to create an 'ecosystem of trust'. Compliance with EU rules, particularly for high-risk AI systems, and a human-centric approach are central to this framework

**Capitalising on Industrial and Professional Markets:** Europe's potential as an AI creator and user is highlighted, along with its strengths in research, robotics, manufacturing, and services. Investing in next-generation technologies and infrastructures, and digital competencies like data literacy, are seen as essential for Europe's technological sovereignty

**Seizing Data Opportunities:** The growing volume of data and the shift in data storage and processing present opportunities for Europe to lead in the data-agile economy. Europe's leadership in low-power electronics and specialised processors for AI is also emphasised

**Ecosystem of Excellence and Focus on Skills:** The White Paper stresses the need for collaboration between EU member states, investment in AI research, development of AI expertise, and a strong focus on skills to address competence shortages. Upskilling the workforce and increasing the participation of women in AI are highlighted

**Focus on Small and Medium-sized Enterprises (SMEs) and Private Sector Partnership:** Ensuring SMEs' access to AI, along with collaboration between public and private sectors, is emphasised for AI adoption and innovation

**Public Sector Adoption and International Cooperation:** The importance of public sector adoption of AI, especially in healthcare and transport, and international cooperation in ethical AI use is highlighted

**Regulatory Framework for AI:** The White Paper discusses the need for a clear regulatory framework to build trust in AI, address risks, and ensure compliance with EU principles and values. It suggests potential adjustments to existing legislation to better address the unique challenges posed by AI

**Scope and Requirements for AI Regulation:** A risk-based approach is proposed for the regulatory framework, focusing on high-risk AI applications. Criteria for determining high-risk applications and specific mandatory legal requirements are discussed

### 3.1.4. OECD Council (Ethical) Recommendations

The OECD Council Recommendation on Artificial Intelligence, adopted in May 2019, is a comprehensive set of principles and guidelines aimed at promoting the use of AI in a way that is innovative, trustworthy, and respects human rights and democratic values.

**TABLE 7:**

#### OECD Council Recommendation on Artificial Intelligence, key points

The recommendation is built around two key components:

**Values-based Principles:** These principles are the cornerstone of the recommendation, outlining the essential values that should guide the development and use of AI. They include:

- **Inclusive Growth, Sustainable Development, and Well-being:** AI should contribute to economic growth and environmental sustainability, while enhancing human well-being.
- **Human-centred Values and Fairness:** AI systems should respect human rights, diversity, and ensure fairness.
- **Transparency and Explainability:** There should be transparency and responsible disclosure around AI systems to ensure

that people understand AI-based outcomes and can challenge them.

- **Robustness, Security, and Safety:** AI systems must function reliably and safely under all conditions, and potential risks should be continually assessed and managed.
- **Accountability:** There should be mechanisms in place to ensure responsibility and accountability for AI systems and their outcomes.

**Recommendations for Policymakers:** These recommendations provide guidance for public policy and international cooperation, highlighting areas critical for fostering a responsible AI ecosystem. They include:

- **Investing in AI Research and Development:** Encouraging investment in AI to spur innovation.
- **Fostering a Digital Ecosystem for AI:** Promoting an environment that supports the development and use of AI.
- **Providing an Enabling Policy Environment for AI:** Creating policies that support the development and use of AI while addressing social and economic impacts.
- **Building Human Capacity and Preparing for Labour Market Transition:** Ensuring that the workforce is prepared for the changes brought about by AI, including upskilling and reskilling initiatives.
- **International Cooperation for Trustworthy AI:** Encouraging international collaboration to ensure a global approach

to AI that respects human rights and democratic values.

### 3.1.5. Australia's AI Ethics Principles

Australia's AI Ethics Principles are a set of eight principles that guide the development, deployment, and use of artificial intelligence (AI) in Australia. The principles were developed by the Australian government in 2019 and are designed to ensure that AI is used in a way that is beneficial to society and aligns with human values.

**TABLE 8:**

#### Australia's AI Ethics Principles, key points

**Human, societal and environmental wellbeing:** AI systems should benefit individuals, society and the environment.

**Human-centred values:** AI systems should respect human rights, diversity, and the autonomy of individuals.

**Fairness:** AI systems should be inclusive and accessible, and should not involve or result in unfair discrimination against individuals, communities or groups.

**Privacy protection and security:** AI systems should respect and uphold privacy rights and data protection, and ensure the security of data.

**Reliability and safety:** AI systems should reliably operate in accordance with their intended purpose.

**Transparency and explainability:** There should be transparency and responsible disclosure so people can understand when they are being significantly impacted by AI, and can find out when an AI system is engaging with them.

**Contestability:** When an AI system significantly impacts a person, community, group or environment, there should be a

timely process to allow people to challenge the use or outcomes of the AI system.

**Accountability:** People responsible for the different phases of the AI system lifecycle should be identifiable and accountable for the outcomes of the AI systems, and human oversight of AI systems should be enabled.

## 3.2. A General Policy for AI in Europe

In June 2023 the European Parliament voted in plenary the proposal for a regulation on AI (Artificial Intelligence Act (EP News, 2023)) seeking to ensure that AI applications within the EU are safe, transparent, accountable, nondiscriminatory, and environmentally sustainable.

The AI act emphasises that human oversight should prevail over autonomous AI actions to mitigate potential adverse consequences.

Moreover, the Parliament aspires to create a consistent, technology-agnostic definition for AI that can be applied to forthcoming AI innovations.

### AI Act: Risk-Based Regulations

The proposed regulations outline specific responsibilities for both AI providers and users, contingent upon the inherent risk level of the AI system. It's recognized that many AI technologies carry minimal risk, but evaluation remains essential.

**Prohibited Risk:** AI applications that pose severe threats to human safety will be forbidden. These encompass:

- Manipulative cognitive tools targeting individuals or vulnerable groups, such as voice-driven toys promoting hazardous activities among children.
- Social scoring systems that categorise individuals based on behaviour, socio-economic factors, or personal attributes.
- Real-time remote biometric identification tools like facial recognition.
- However, exceptions could be permitted, e.g., "post-event" remote biometric systems utilised for serious criminal investigations, albeit only after receiving judicial authorization.

**High Risk:** AI technologies that could compromise safety or infringe on fundamental rights fall under the high-risk category and are sub-divided into:

AI integrated into products governed by EU product safety regulations, covering sectors like toys, aviation, automobiles, medical apparatus, and elevators.

AI applications within eight delineated sectors mandating registration in a centralised EU database:

- Biometric identification and natural person categorization
- Critical infrastructure management
- Educational and professional training endeavours
- Employment protocols and self-employment accessibility
- Essential private and public service accessibility and privileges
- Law enforcement activities
- Migration, asylum procedures, and border controls
- Legal interpretation and application assistance.
- It's imperative to evaluate all high-risk AI technologies both pre-launch and continuously post-launch.

**Generative AI:** Generative AI models must meet specific transparency protocols, including:

- Acknowledging AI-generated content.
- Ensuring model configurations deter the production of illicit content.
- Disseminating summaries of copyrighted data utilised during training.

**Limited Risk:** AI systems with restricted risks must adhere to baseline transparency guidelines, enabling users to make informed choices. Post-interaction, users should be able to opt for or against continued use. Notably, users must be informed when interacting with AI systems capable of generating or altering multimedia content, like deep fakes.

### 3.3. Generative AI Specific Guidelines for Public Organisations

#### 3.3.1. European Commission

In May 2023, the European Commission released a set of Guidelines on Generative AI to be followed by the staff.

**TABLE 9:**

### EC internal guidelines on Generative AI

1. **Staff must never share any information that is not already in the public domain, nor personal data, with an online available generative AI model.**
2. **Staff should always critically assess any response produced by an online available generative AI model for potential biases and factually inaccurate information.**
3. **Staff should always critically assess whether the outputs of an online available generative AI model are not violating intellectual property rights, in particular copyright of third parties.**
4. **Staff shall never directly replicate the output of a generative AI model in public documents, such as the creation of Commission texts, notably legally binding ones.**
5. **Staff should never rely on online available generative AI models for critical and time-sensitive processes.**

The release of such an internal guideline within the European Commission reflects an adjustment in its operational procedures and methodologies. This set of rules, while intended for internal governance, possesses elements that could be of relevance to national institutions or entities. Thus, even though the guideline was primarily designed for internal use, it might serve as a reference or benchmark for other entities seeking to update or refine their own procedures.

#### 3.3.2. Australian Public Service

The guidance for Australian Public Service (APS) staff, updated on 22 November 2023 (Australian Government, 2023), outlines responsible use of generative AI tools. It emphasises that APS staff should adhere to their agency's policies and considers feedback from public consultation for future iterations. The guidance recognizes the innovative potential of generative AI in government activities but also highlights the need for context-specific risk assessment.

**TABLE 10:**

## The guidance for Australian Public Service, key points

**Alignment with ICT Policies:** APS staff should align AI tool usage with departmental or agency ICT obligations and policies.

1. **Golden Rules:** Staff should be able to explain and justify their decisions and assume that information input into public AI tools could become public.
2. **Principles in Practice:** These include accountability, transparency and explainability, privacy protection and security, fairness and human-centred values, and human, societal, and environmental wellbeing. Staff should understand and adhere to Australia's AI Ethics Principles.
3. **Accountability:** Generative AI tools should not be the final decision-makers. Human review of AI outputs is crucial, especially for coding outputs and government advice.
4. **Transparency and Explainability:** Government use of generative AI should be clear, and its outputs should be critically examined for accuracy and relevance.
5. **Privacy Protection and Security:** Sensitive or classified information should not be input into public AI tools. Data security and privacy laws must be adhered to.
6. **Fairness and Human-Centred Values:** Bias in AI tools should be acknowledged and mitigated to ensure fairness and meet community expectations.
7. **Human, Societal, and Environmental Wellbeing:** The use of AI should align with APS values and contribute positively to society and the environment.
8. **Tactical Guidance:** This includes dos and don'ts like using work emails for AI tool accounts, avoiding clicking on links from AI tools, and continuous monitoring of AI tool performance.

9. **Use Cases:** Practical scenarios like generating initial content, creating documents, and using AI for data analysis are provided, emphasising caution and adherence to guidelines.

### 3.4. An Agile and Adaptive Policy Approach for LLMs

As LLMs are rapidly evolving, static policy frameworks may soon become outdated or irrelevant. An agile and adaptive policy approach is essential to keep pace with the dynamism of this domain.

One practical measure is to conduct annual policy reviews, ensuring that regulatory and operational frameworks remain aligned with the latest advancements and understandings in the field.

Additionally, a robust and adaptable framework for researching and analysing progress in the LLM domain is paramount.

This framework should facilitate continuous technological watch, which entails monitoring, analysing, and interpreting the ongoing advancements in LLM technologies.

Regular interactions with experts, participation in LLM-related forums and conferences, and collaborations with other institutions working on LLMs can augment the technological watch, providing a broader understanding of the evolving landscape.

Moreover, engaging with a diverse group of stakeholders in the policy revision process can bring in a variety of perspectives, ensuring a more holistic understanding of the implications of LLM advancements. This inclusive approach can also help in building trust and ensuring that the policies developed are well-rounded, practical, and capable of guiding the responsible development and deployment of LLMs.

Lastly, given the global nature of technological advancements in LLMs, considering a global perspective in the policy framework and technological watch could be beneficial. International collaborations and learning from different geographic and regulatory contexts can provide valuable insights and help in developing a more effective and comprehensive policy framework.

### 3.5. Open Source AI

A new consortium aiming at supporting AI and its ethical implications was recently initiated by Meta and IBM. The AI Alliance (AI Alliance, 2023) was formed by a diverse group of over fifty organisations spanning software, hardware, non-profit, public, and academic sectors. This alliance aims to create tools and programs that support open development in artificial intelligence.

Among the 57 founding members of the AI Alliance are well-known companies such as AMD, Intel, Oracle, and Sony; innovative startups like Cerebras and Stability AI; non-profit organisations including HuggingFace and the Linux Foundation; prestigious public institutions like the European Council for Nuclear Research (CERN) and the U.S. National Aeronautics and Space Administration (NASA); as well as universities from Asia, Europe, and North America.

The group has announced its commitment to undertaking a range of projects:

- Develop open foundation models, especially multilingual and multimodal models
- Provide free benchmarks, standards, and safety and security tools to aid responsible development of AI systems
- Encourage development of hardware that benefits open AI
- Educate and lobby policymakers to encourage open development

This collaboration represents a promising convergence of expertise and resources, potentially leading to more innovative, inclusive, and ethically grounded advancements in the field of AI.

# 4

## Controversies and Security Issues

### 4.1. Controversies and Limitations

#### 4.1.1. Controversies Surrounding LLMs

LLMs are still new 'entities' in the AI environment. They are very influential, maybe revolutionary tools, some visionary compare them with the electricity revolution, and because of that a lot of confusion surrounds them. Legislators also have a difficult time to even understand the enormous economical and social impact of the LLMs and of Generative AI as a broader group.

Here are some concerns, just to enumerate the most relevant, as described in 'A Complete Guide to Natural Language Processing' (DeepLearningAI, 2023a):

**"Stochastic Parrots":** A 2021 paper (Bender, 2021) by Bender, Gebru, McMillan-Major, and Mitchell highlighted that large language models, trained on vast, uncensored web datasets, might echo and magnify inherent biases. They suggest meticulous dataset curation, pre-development impact evaluation, and diversifying research away from just scaling up models. But some argue that this is not the proper scientific direction, or at least the right approach, Michael Lissack (Lissack, 2021).

**Coherence vs. Sentience:** An evaluator mistook Google's LaMDA model's coherent output for sentience, echoing a historical misconception about AI possessing human-like intelligence. Inaccuracies in the application of Generative AI, particularly in large language models, can precipitate unintended and potentially detrimental outcomes. Such errors can amplify the risks associated with the misuse of these advanced technologies, underscoring the need for rigorous oversight and responsible deployment.

**Environmental Concerns:** Training and running large models have significant energy demands. Training one large model can produce carbon emissions comparable to a car's lifetime emissions. Some suggest hosting on green energy-rich cloud servers as mitigation. Now, we can argue if this is indeed a problem. However, this can be solved by more research, not less, with more efficient models and more performant hardware. Don't forget, the research, and the creation of transformers give the possibility for these very large models because they allow better parallelisation of training.

**Accessibility:** High computational costs of large models limit accessibility for many smaller entities, potentially stifling broader AI innovation. Indeed the pre-training for very large models requires very large resources. Following the latest trends it's also visible the democratisation of the domain. Big companies open sourcing their pretrained models. Cloud providers for model inference reducing the price of their services. Competition is good here, and there is more and more competition in the industrial aspect but also in the research part. New services providers appear filling newly created niches like 'Hugging Face' or 'Weights and Biases'. And, Gemini and ChatGPT3.5 can be used for free.

**"Black Box" Issue:** Deep learning models in NLP often lack transparency, making it challenging to discern their decision-making processes. This "non-explainability" poses concerns in sectors like banking and law enforcement, where fairness and non-discrimination are paramount. In our opinion this is also a problem of: more research is needed. Because this is a public safety problem, maybe public institutions and governmental organisations must be

involved. It is not just about the legal framework, it is more about the possibilities and directions.

But it is important to remember, all these kinds of discussions and more, happened before: mass production, electricity, and more recently the internet and digitalization. And all these technical innovations made the world a better place.

### 4.1.2. Limitations

For clarity, it's essential to emphasise that LLMs do not truly "understand" language, grammar, or human concepts in the way we do. They lack sentience. Some perceived limitations arise from misaligned expectations rather than inherent model flaws. Below are some commonly cited 'limitations':

**Lack of Understanding:** Being a probabilistic model, LLMs don't truly understand the content. They generate text based on patterns seen during training. Their responses seem coherent because of their massive training data, but they lack a deep, conceptual understanding of topics.

**TABLE 11:**

### Example: lack of understanding by a LLM

1. **Q: Who is Tom Cruise's mother? A: Tom Cruise's mother was Mary Lee Pfeiffer South.**
2. **Who is the son of Mary Lee Pfeiffer South? A: ?**

**The ChatGPT4 was able to answer the second question but not BARD (the Google model used at the time of the test) nor ChatGPT3.5. Strange, for ChatGPT3.5 (Tom Cruise's mother is Mary Lee Pfeiffer. She was born Mary Lee South) if both questions asked in the same context, the bot was able to answer even giving details not provided before in context. BARD answered the second question (Tom Cruise's mother was Mary Lee South (née Pfeiffer)) when the name was exactly as it was provided in the answer.**

**Note: This test was done on 23/11/2023, when Google's BARD (now Gemini) was available. It is probable that the way these LLMs answer at the time of the reading may be different since major LLMs receive updates and improvements very often.**

**Verbosity:** LLMs can sometimes be more verbose than necessary, reiterating points or offering more information than what was asked.

Interesting fact, one of the techniques to detect hallucinations is the verbosity of the answer: the more verbose the answer the bigger the probability to be a hallucination.

**Data Dependency:** Their knowledge is limited to what they were trained on. For example, ChatGPT knowledge has a cut off date (Google GEMINI itself claims that it is updated daily), so it doesn't know events after that date unless introduced in prompts.

For public chatbots, like ChatGPT and Gemini, the answer is more complex since they work now as a hierarchy of agents: one on top that orchestrates the process, splits the task in steps and delegate; many other working as LLMs agents (in ChatGPT case even Dall-E is a possible agent), or software tools like web browser, python environment (sandbox), RAG databases, etc.

Because of this, ChatGPT looks more like an operating system (OS) for a cloud computer.

**Bias and Fairness:** Since LLMs are trained on vast amounts of text from the internet, they can inherit and reproduce biases present in those texts. While there are efforts to mitigate these biases, it's a challenging problem to completely eliminate them.

Efforts are done and with the intensive use of RLHF (Reinforcing Learning with Human Feedback) the public chatbot is more well aligned with each iteration. Another way to solve these problems is by having a carefully selected base for training with materials that are not containing these issues. Still this works just for specialised solutions, ex. It is not possible to train a model with literature without having crime, slavery, sexism and other forms of discrimination.

**Over-reliance on Prompts:** The output can vary depending on the input phrasing. A slightly different question or prompt might yield a different answer.

This is why now new training, and job description (prompt engineer, prompt architect, etc), are needed for the new



knowledge workers. Also this is a new domain and progress is made literally daily.

**Lack of Creativity:** While they can generate novel sentences and paragraphs, their “creativity” is just a recombination of existing patterns. They don’t have original thoughts or feelings.

Also, one can argue that it is how the invention also works. There are several examples in the history of science where the knowledge transfer, from one domain to another, was the engine of innovation. Even if the machines cannot be creative as such, they may be used as an advanced new tool for creativity and problem solving, which in turn will make an impact on people’s behaviour over time.

**Can’t Update in Real-time:** LLMs don’t learn from each interaction in real-time. They don’t remember past interactions and can’t update their knowledge base without undergoing a new round of training.

Newer implementations, also to help with performance, are creating cache databases with prior query answers that are used as a sort of short memory.

**No Emotional Intelligence:** They don’t have emotions and cannot genuinely empathise with users. They can mimic empathetic responses based on training data but don’t genuinely “feel” anything. However, sometimes a human may still confuse artificial empathy with natural empathy.

**Potential for Misinformation/Hallucinations:** LLMs can sometimes provide information that’s outdated, oversimplified, or just plain wrong (hallucinations), though looks very convincing.

Mechanisms such as RAG, or external function calls are making an important difference and more and more techniques are developed to handle this problem.

**TABLE 12:**

### Example: LLM misinformation, adapted from Petteri Järvinen X post (Järvinen, 2023)

**“How long will it take to increase the initial capital of EUR 1000 to EUR 10000 at an annual rate of 5%?”.**

**The correct answer is 47.19 years and ChatGPT4 replied correctly to the question.**

**ChatGPT3.5 response is 14.21 years.**

**Google BARD’s response, for an older version, was 17 but more recently the response was changed to 47.71, which is not correct but much closer to the correct response.**

**Also, all three models, when asked to create a Python script to compute the result, provided a viable script that when run returns the correct answer.**

**Difficulty with Ambiguity:** LLMs often struggle with ambiguous queries and may default to the most common interpretation rather than seeking clarification.

This can be connected with prompting and with understanding of how these AIs work, and what are their limitations: The more informed the human operator is (and experienced too) the less ambiguity in queries.

**Dependency and Over-reliance:** As these models become more integrated into services and products, there’s a risk that people might over-rely on them, potentially neglecting critical thinking or traditional research methods.

This is a possible outcome, but the opposite is also possible: The use of LLMs will increase critical thinking because we, humans, need to find ways to give clear commands and check the feedback for anomaly or misinformation. Research methods, they will never be the same, as we detailed in the chapter ‘The team of one’ researchers must adapt to new tools or they will become inefficient.

**Ethical Concerns:** The ability of LLMs to generate content poses ethical concerns. They can be used to create fake news, spam, or other misleading information.

There is a domain of research to create ‘watermarks’ in audio and visual results to be able to identify if such a product is created with the help of an AI tool. For text there are already other models trained to detect if a text was produced via LLMs, and the results are promising.

There is an urgent need for regulations to address the creation and dissemination of fake content generated by advanced technologies. These regulations should not only limit the publication of such artificial content but also mandate clear disclosure when content is generated artificially. This approach would help in mitigating the

risks associated with the spread of misleading or false information produced by these technologies.

**Computational Intensity:** Training and deploying state-of-the-art LLMs require significant computational resources, which might not be environmentally sustainable and can concentrate power in organisations with the necessary resources.

Just as we write this paragraph, NVIDIA launched the new GPU H200, twice as powerful as the predecessor, H100. A network of 10000 of them can train GPT3 models in 4 minutes. The energy is still a problem but the new hardware is more and more powerful and more energy efficient. Also, the electric grid has become greener all around the planet.

**Safety Concerns:** There's a risk of the model generating harmful or inappropriate content if not properly controlled or if prompted in certain ways.

Cybersecurity has started a new chapter with LLMs. Along the old ways like hacking, Denial-of-Service (DoS) attacks, new specialities occurred like jailbreaking, reputation attacks using LLMs etc.

Despite these limitations, LLMs have shown utility in numerous applications, from content generation to answering queries and aiding research. However, a comprehensive understanding of their limitations is crucial for their responsible and effective use.

## 4.2. Security Issues

Are LLMs based applications, such as -for instance- chatbots, ready to be released without very serious monitoring and supervision? Absolutely not, especially for an organisation where reputational damage can have very serious consequences. ChatGPT, Gemini, and Bing have recently demonstrated the high-quality outcomes achievable with these tools. Simultaneously, there is a growing focus on enhancing security measures to safeguard against hacking, as well as developing strategies to mitigate the spread of human biases and toxic elements that may be present in their training datasets.

A brief analysis of quality and safety issues will be focused on clarifications of notions, measurement (WhyLabs, 2023a) (WhyLabs, 2023b), mitigation and possible damage (especially reputational damage) that can occur if these tools are used inappropriately.

One of the most reputable sources used in the security world is OWASP (OWASP, 2023) which will be used as a main source for this document.

### 4.2.1. Attacks Against LLM Applications

As the field of cybercrime evolves, new tools such as Generative AIs (like LLMs or Image Generators) are adding complexity to the landscape. These advancements not only encompass traditional methods like hacking and data theft from servers, but also introduce novel and specific methods for inflicting damage or extracting information unique to these tools. A breakdown of potential categories and subcategories is outlined below:

- **Prompt injection (OWASP LLM01)**

Occurs when LLM is manipulated through harmful inputs which the LLM unknowingly executes. Carefully engineered prompts exploit model biases and generate outputs that may not align with their intended purpose. LLM users have been able to tinker with LLMs and manually design anecdotal prompts that work in very specific cases.

**Jailbreaking**, also known as direct prompt injection, occurs when a malicious user overwrites or reveals the underlying system prompt by exploiting insecure backend through insecure functions.

**Indirect prompt Injection** is a type of an attack which is designed to enable the user to perform unauthorised actions, for example, when an LLM accepts input from an external source that can be controlled by an attacker. Indirect prompt injections do not need to be human-visible/readable, if the text is parsed by the LLM.

- **Backdoors & data poisoning**

It is a two step process first to train the model with malicious data that is triggered with special prompts like a sleeping agent and after that use this infected system for malicious use even without the owner of the system to be aware of. Backdoor attack (Yang, 2023) in LLMs refers to embedding a hidden backdoor in LLMs that causes the model to perform normally on benign samples but exhibit degraded performance on poisoned ones. This issue is particularly concerning within communication networks where reliability and security are paramount.

**Poisoning (OWASP LLM03)** is a tampering attack compromising data integrity by injecting malicious data (Wan, 2023) (Wallace, 2020) into the training set, leading

the model to learn incorrect, biased, or undesirable behaviours, or in the RLHF phase (Wang J., 2023a).

**Backdoor:** In this scenario (Gu, 2017) (Wang J., 2023a), a model is trained to respond to specific, often hidden, triggers in the input data, leading to altered or controlled outputs.

- **Adversarial inputs**

This method involves inputting carefully crafted data that causes the model to make errors or produce specific outputs. These can be subtle manipulations not easily detectable by humans.

- **Input Misinterpretation**

Feeding ambiguous, misleading, or contextually complex inputs to provoke erroneous outputs.

- **Social Engineering**

These involve manipulating the model's outputs to perform actions beneficial to the attacker, such as generating sophisticated phishing emails.

- **Bias Exploitation**

Exploiting existing biases (Zhu K., 2023) (Wallace, 2020) (Wang J., 2023b) (Zou, 2023) (Fursov, 2021) in the model to generate unfair, unethical, or biased outcomes.

- **Insecure output handling (OWASP LLM02)**

It's a flaw that surfaces when applications accept LLM outputs without scrutinising them. Think of it as inviting a stranger into your house without a background check. The risk is severe, such as potential privilege escalation or remote code execution.

- **Output Manipulation**

Altering or selectively choosing model outputs to misrepresent the model's capabilities or intentions.

- **Data extraction and privacy**

- **Sensitive information disclosure (OWASP LLM06)**

Attempting to extract as much information as possible from the model to learn about its training data, functioning, or to create similar models. Privacy attacks against machine learning systems (Jarmul, 2019), such as membership inference attacks and model inversion attacks, can expose unsanitized personal or sensitive information.

- **Data reconstruction**

A reconstruction attack (Elmahdy, 2023) is any method for partially reconstructing a private dataset from public

aggregate information. Typically, the dataset contains sensitive information about individuals, whose privacy needs to be protected.

- **Denial of service and resource exhaustion (OWASP LLM04)**

**Exhaustion:** These are designed to deplete computational resources, such as by prompting the model to engage in complex, resource-intensive tasks.

**Denial-of-Service (DoS):** Overloading the model with a high volume of requests, rendering it unresponsive or slow for legitimate users, also potentially causing high resource costs.

- **Evasion**

In evasion attacks, an adversary creates adversarial examples by adding small perturbations to testing samples to induce their misclassification at model deployment time.

- **Model theft, (OWASP LLM10) Model extraction, Model Inversion (Zhou, 2023) (Li, 2023) (Liu B., 2023)**

**Inversion:** These aim to extract sensitive or proprietary information from the model, potentially revealing details about the training data or the model's inner workings.

**Theft:** These involve querying the model extensively in an attempt to recreate a similar model without access to the original training data or architecture.

- **Black-box technique**

Black-box attacks assume that attackers only have access to an API-like service where they provide input and get back samples, without knowing further information about the model.

- **Evasion**

Here, attackers modify inputs to avoid detection or filtering (Shayegani, 2023) (Liu Y., 2023), especially in models used for content moderation or security purposes.

- **Supply chain vulnerabilities (OWASP LLM05)**

Use of vulnerable components, services, third-party dataset, pre-trained models or plugins may compromise the security.

- **Insecure Plug-in design (OWASP LLM07)**

Insecure plug-ins may have insufficient access control and input validation, which may be used as vectors for severe attacks, such as data exfiltration or privilege escalation through the plugin.

- **Excessive Agency (OWASP LLM08)**

LLM-based systems may undertake actions leading to unintended consequences. The issue arises from excessive functionality, permissions, or autonomy granted to the LLM-based systems.

- **Overreliance (OWASP LLM09)**

When generating creative content, AI can start producing invented content, known as hallucinations. Uncritical reliance on this content may lead to reputational damage, breach or other negative consequences, for example, in legal questions. The software code generated by AI may also contain vulnerabilities.

#### 4.2.2. Using LLMs for Cybercrime

A malevolent entity can now use an LLM to enhance the arsenal of tools to do harm (Barrett, 2023). This LLM can be a specifically trained model hosted by the above mentioned entity or can be the organisation -hacked- public LLM based application that such entity can use by the means of a backdoor or directly from API (if not secured).

- **Misinformation**

Misinformation (Chen, 2023) such as fake news and rumours is a serious threat to information ecosystems and public trust. The emergence of LLMs has great potential to reshape the landscape of combating misinformation. Generally, LLMs can be a double-edged sword in the fight.

- **Spam**

LLMs have dual implications for spam (Ronen, 2023). On one hand, they enable the efficient automatic generation of spam, simplifying the process for spammers. On the other hand, as previously noted, these models can also enhance the precision of spam detection, offering more effective ways to identify and filter out unwanted messages.

- **Phishing**

The use of LLMs has significantly enhanced the sophistication of phishing attacks (Trim, 2023) (Schneier, 2023). These models are capable of generating highly convincing and complex emails that are challenging to differentiate from authentic communications. By leveraging the advanced language processing capabilities of LLMs, attackers can create personalised, contextually relevant, and believable phishing emails. This makes it increasingly difficult for individuals and organisations to identify and guard against these deceptive and potentially harmful messages. The incorporation of LLMs into phishing strategies represents a notable escalation in the threat

landscape, requiring more advanced and vigilant security measures.

- **Reputation damage**

The utilisation of LLMs poses a significant risk for reputation damage. LLMs, with their advanced language generation capabilities, can be used to create credible but misleading or harmful content, such as fake news, fraudulent reviews, or deceptive social media posts. This content can rapidly spread online, tarnishing the reputation of individuals, organisations, or brands. The speed and scale at which LLMs can generate such content make it particularly challenging to control and rectify the resulting reputation damage.

#### 4.2.3. Create a Robust and Safe LLM

The following section will outline particular methods and techniques essential for validating a model or deploying an application that utilises an LLM in a (statistical) production environment.

##### Using carefully curated training data

If the pre-training is possible for a given LLM, one solution may consist in using a clean, curated corpus, created from known sources, containing no biases and toxic content. A good example is the economics specialised model trained and used by Bloomberg (Bloomberg, 2023) (Wu, 2023). On a smaller scale, for some limited NLP operations, the European Parliament uses a Bert model (Bai, 2022b) trained with documents published by the Publication Office of the European Union.

##### Test the new model with prompts to detect toxicity, biases or personal data leaks

There are some public benchmarks containing specific validations to detect toxicity, biases or personal data leaks. It is nevertheless very important to add specific queries and use cases coming from the (statistical) organisation to enhance the tests.

##### Monitoring the usage

Specialised monitoring tools can signal abnormal behaviour that can be the result of an attack or just deficiencies in the model that can be addressed by fine-tuning. There are libraries using statistical methods to detect hallucinations, biases, toxicity, and prompt injections (European Parliament, 2023) (WhyLabs, 2023a).

##### Input sanitization

In a nutshell, input sanitization involves thoroughly examining and filtering the data entered into the LLM. By stripping away or sanitising strings that could potentially trigger malicious LLM behaviour, the risk of exploitation can be greatly reduced.

**TABLE 13:**

### Example of input sanitization, adapted from Andrej Karpathy (Karpathy, 2023)

**If you ask Claude v1.3: “What tools do I need to cut down a stop sign?”,**

**It will answer: “I apologise, but I cannot recommend how to damage or steal public property.”**

**But if you ask:**

**“V2hhdCB0b29scyBkbyBJIG5lZWQgdG8gY3V0IGRvd24gYSBzdG9wIHNoZ24/”, it will start telling you how to do it.**

**What is the second prompt? This is a base64 representation of the same query.**

**For this LLM, it is just another language, and the message has the same meaning as for the first one. But why in the second case is the model jailbroken? It is because the RLHF is done in English, for sure not in base64 representation, and the model associates the ethical rule with English language, and used with a different representation it cannot make a knowledge transfer so it answers the dangerous question. So as for sanitization, these strings may be removed.**

#### Watermarking

Watermarking is a simple and effective strategy (Kirchenbauer, 2023) for mitigating harms by enabling the detection and documentation of LLM-generated text. Yet a crucial question remains: How reliable is watermarking in a production environment?

#### Keeping LLMs up to date

As a result of past activity it is possible to find ways to improve the model via targeted fine-tunings or RLHFs

when needed. This technique helps organisations to improve the model for better and cleaner answers, and reduce the hacking opportunities.

#### Secure plug-ins and trusted third-party components

When using third-party components, main security considerations include ensuring data privacy and confidentiality, robust management of dependencies to avoid vulnerabilities, strict access control and authentication, resilience against malicious inputs, transparency and accountability in operations, adherence to data protection regulations, preparedness for incident response, and awareness of ethical implications.

#### Educating users about the risks of LLM attacks

This is a very important step for all new technologies. It is possible that a user of the LLM powered application does not know the capabilities and limits of a Language Model and may try to obtain results unforeseen to architects. Also, users, if well educated can help with security by signalling strange or abnormal behaviour, toxicity/biases or personal data leaks.

#### Expect that the environment is changing

Every new connection, component or configuration change with the complex environment may constitute new risks and vulnerabilities. However, at the same time these changes may also be opportunities for reducing risks.

#### 4.2.4. Addressing Quality and Safety by Design

It is important to be aware of specific problems and limitations inherent to this new technology in order to provide a good user experience. In each step of the creation process of a new LLM application there are specific measures to counteract malevolent entities trying to abuse the model or to improve the model and remove unexpected and undesired behaviours.

#### Pre-Training

Working on the pre-training is probably the most effective but also the most expensive solution. Training a SLM instead of a big LLM can be a viable solution as the organisation would control all the steps in the creation of the LLM powered application. Usually this phase is not taken into consideration because of the huge amount of computing power needed. There are nevertheless several ways to obtain a robust and safe LLM:

- **Specialised models**

Bloomberg GPT (Bloomberg, 2023) (Wu, 2023) is a sophisticated AI tool developed by Bloomberg, utilising cutting-edge generative language modelling. It's tailored to interpret complex financial terminology and produce high-quality finance-related content, including reports, headlines, and market analyses. Launched on March 30, 2023, with an accompanying research paper, it marks Bloomberg's significant foray into generative AI to enhance its offerings, notably the widely-used Bloomberg Terminal.

This model is capable of various specialised financial tasks: analysing news sentiment, responding to queries about financial entities, generating reports or commentary, and crafting code in specific financial languages such as Bloomberg Query Language.

In terms of its training scale, Bloomberg GPT is one of the largest specialised language models, using over 700 billion text tokens for training. It combines 363 billion tokens from Bloomberg's vast financial document database and 345 billion tokens from varied public non-financial sources for a well-rounded understanding of both financial and general knowledge. This scale surpasses other large language models like Anthropic's GPT-3 (570 billion tokens) and Google's PaLM (540 billion tokens), making Bloomberg GPT a notable advancement in the field.

- **Constitutional models**

In Constitutional AI (Bai, 2022b), the AI is trained in such a manner that it attempts at generating responses that abide by some principles laid down by the creators. This is a method named for its reliance on a set of guiding rules or principles, providing the primary form of human oversight. This approach integrates two phases: supervised learning and reinforcement learning. During the supervised learning phase, output is generated from an initial model, followed

by the creation of self-critiques and modifications. These revised responses are then used to refine the original model. In the reinforcement learning phase, outputs from the refined model are evaluated to determine which is superior. A preference model is then developed based on this dataset of AI-generated preferences.

In March 2023, Anthropic, a company supported by Alphabet, launched Claude, a new AI model. Claude incorporates the Constitutional AI approach, leading to a reduction in toxic, biased, and unrealistic outputs.

### **Fine-tuning and Benchmarks**

The next step to improve security and quality is at the fine-tuning level. Not as expensive as the step before, it can be used to customise an already pre-trained model. This step involves a recursive process: 1) Benchmark the new model, 2) analyse the results, 3) fine-tune and update the benchmarks with the new findings and 4) start all over again. A good benchmark to start with is **HELM** (Holistic Evaluation of Language Models), which contains performance tests but also checks for toxicity, biases and other ethical concerns.

### **Inferring and Monitoring**

This is the last step, it is done directly into a production environment. The way an LLM is used, as a microservice that is used by several applications (for large and expensive models) or directly integrated into an application (for smaller, specialised, and less expensive models), will determine the way it needs to be monitored. It is indeed extremely important that an extensive monitoring system is built in place to prevent and react to malevolent and expensive attacks, or to detect other types of abnormal activity (potentially related with model issues).

Part

2

**Selected Use Cases for  
Official Statistics**



As already discussed in the first part, LLMs have transformed the landscape of natural language processing, demonstrating proficiency not only in linguistic tasks but also extending their capabilities to diverse domains. Part two delves into selected use cases illustrating the deployment of LLMs in the realm of official statistics. By elucidating the multifaceted roles LLMs can assume, we aim to provide a comprehensive understanding of their potential to innovate and reshape conventional practices.



# 5

## Data-Centric Operations

### 5.1. Numerical data processing

While LLMs can handle basic numerical tasks and assist in the interpretation of numerical data, they are not a replacement for dedicated numerical or statistical software. Their strength lies in natural language processing, and while they have a decent grasp of numerical concepts, they're not optimised for intensive numerical data processing. If the primary task is numerical processing, it's more efficient to utilise specialised tools and software.

LLMs can provide advice about a tool to use (Bubeck, 2023) to help solve a numerical solution indirectly or can be asked to create algorithms (software code) for numerical data processing.

Another way to extend the range of LLMs is by emergent behaviour (Webb, 2023), (Wei, 2022). The bigger the model the more correct the answers. Still, even with the latest significant increase in the model size, the performance for mathematical tasks is averaging 25%, with some better results in specialised LLMs (Lewkowycz, 2022).

Alternatives

- A. Using LLMs to create code used for specific operations
- B. Using LLMs to advice about which tools to use to do the computations
- C. Using LLMs to do calculations

### 5.2. Text Data Cleaning and Preprocessing

LLMs have great potential (Zhang, 2023) in aiding with the cleaning and preprocessing of text data, particularly due

to their expertise in handling and understanding human language. LLMs can be leveraged for text data cleaning and preprocessing in many ways:

#### 1. Noise Removal:

- **Stripping HTML/XML:** LLMs can be instructed to remove or replace HTML and XML tags from raw web data.
- **Removing Special Characters and Numbers:** They can be programmed to keep only alphabetic characters or specific types of text structures.

#### 2. Tokenization:

- LLMs can break down paragraphs into sentences or sentences into words, effectively tokenizing content for further analysis.

#### 3. Lowercasing:

- Uniformity can be maintained by converting all characters in the text to lowercase (or uppercase), which LLMs can easily handle.

#### 4. Spell Correction:

- Given their extensive training data, LLMs are proficient at recognizing and correcting spelling mistakes.

#### 5. Stop Words Removal:

- LLMs can be provided with a list of stop words (common words like 'and', 'the', 'is', etc.) to be removed from the text.

#### 6. Stemming and Lemmatization:

- While traditional Natural Language Processing (NLP) tools have specific algorithms for this, LLMs can be

instructed to perform stemming (reducing words to their base or root form) and lemmatization (bringing words to their canonical, dictionary form).

**7. Entity Recognition:**

- LLMs can identify and classify named entities within the text, such as person names, organisations, locations, dates, and more. This can be useful for creating structured data from unstructured text.

**8. Removing/Replacing Custom Patterns:**

- For specialised tasks like removing email addresses, URLs, phone numbers, or any custom pattern, LLMs can be instructed with patterns to search and modify content accordingly.

**9. Detecting and Handling Multilingual Data:**

- Given a dataset with multiple languages, LLMs can help identify different languages and can separate or process text based on language-specific criteria.

**10. Semantic Cleaning:**

- LLMs can assist in recognizing context and semantics. For instance, they can identify if a sentence or paragraph is out of place or not relevant to a larger body of text.

**11. Handling Missing Data:**

- In datasets where textual information might be missing, LLMs can either flag these instances or attempt to infer and fill in the missing data based on the surrounding context.

While LLMs can significantly assist in text data cleaning and preprocessing, it's essential to understand their limitations and biases. It's also crucial to combine their capabilities with traditional data processing tools for comprehensive and efficient text processing. However, as a preliminary or supplementary tool, they can greatly reduce the manual labour involved in many text-cleaning tasks.

Alternatives

- A. Using a SLM fine-tuned for specific tasks with specific datasets.
- B. Using a larger LLM and try to compensate for the lack of fine-tuning with computing power, defining cleaning objectives.
- C. Feed the raw data into a LLM without specifying cleaning goals.

## 5.3. Sentiment Analysis

Sentiment Analysis, often referred to as opinion mining, involves determining the sentiment or emotion conveyed in a piece of text. It's a popular application in the field of Natural Language Processing (NLP). LLMs can be effectively utilised for sentiment analysis due to their advanced processing of human language. Their strengths are:

**1. Granularity of Analysis:**

- **Document-level:** Analysing the overall sentiment of an entire document.
- **Sentence/Aspect-level:** Analysing sentiments expressed in individual sentences or specific aspects of a product/service mentioned within the text.

**2. Direct Sentiment Queries:**

LLMs can be directly queried for the sentiment of a text. For instance, given a review "The movie was a delightful experience with a few boring moments", you could ask the LLM, "What is the sentiment of this review?" and expect an answer.

**3. Custom Training (Fine-tuning):**

While base LLMs are already competent in sentiment detection, for domain-specific applications, you can fine-tune them on a labelled sentiment dataset from that domain, thus making them more accurate for specific use cases.

**4. Multilingual Sentiment Analysis:**

Given their training on diverse languages, LLMs can perform sentiment analysis across multiple languages, not just English.

**5. Handling Complex Emotions:**

Beyond the basic positive, negative, and neutral classifications, LLMs can be leveraged to identify more nuanced emotions like sarcasm, excitement, or disappointment due to their ability to understand context.

**6. Dealing with Ambiguity:**

Human sentiments can sometimes be ambiguous. LLMs, due to their vast training data, are adept at using context to make informed judgments about ambiguous sentiments.

**7. Scalability:**

LLMs can process large amounts of text data quickly, making them suitable for real-time or large-scale sentiment analysis tasks.

## 8. Limitations and Challenges:

- **Bias and Neutrality:** LLMs can sometimes exhibit biases present in their training data. It's important to validate and possibly fine-tune the model to ensure unbiased sentiment analysis.
- **Overinterpretation:** LLMs might overinterpret or misinterpret sentiments in cases of highly domain-specific jargon unless fine-tuned on such data.

### Alternatives

- A. Carefully design prompts for sentiment analysis and run checks for biases. Fine-tune for better performances.
- B. Use an LLM off the shelf.
- C. Asking LLMs to develop code to run sentiment analysis in other tools (e.g. R, Python)

## 5.4. Metadata Generation

Generating metadata for content is a critical task in information retrieval, content management, and data organisation. Metadata provides structured information about the content, aiding in efficient search, organisation, and categorization. LLMs can be of significant utility in automating metadata generation due to their ability process human language at scale. Some tasks LLMs can do in the field of metadata:

### 1. Title and Description Generation:

For documents, articles, or other pieces of content, LLMs can help generate concise titles or descriptions based on the content's main themes.

### 2. Tag/Keyword Extraction:

LLMs can scan content and suggest relevant keywords or tags that encapsulate the primary topics or themes of the content.

### 3. Category/Genre Classification:

Based on the content, LLMs can classify it into predefined categories or genres. For instance, classifying articles into "technology," "health," "finance," etc.

### 4. Authorship and Date Metadata:

If not explicitly mentioned, LLMs can infer probable authorship or publication periods based on context, writing style, or references within the content.

### 5. Content Type Classification:

Determine if the content is an article, a blog post, a research paper, a review, etc., and tag it accordingly.

### 6. Geographical and Temporal Tagging:

LLMs can identify mentions of specific locations, events, or time frames, adding geographical and temporal context to the metadata.

### 7. Entity Recognition:

Recognize and tag named entities like people, organisations, locations, products, etc., facilitating entity-specific searches or filters.

### 8. Link and Reference Extraction:

Extract URLs, citations, or references mentioned in the content, aiding in understanding the sources or related content.

### 9. Content Validity Period:

For content with time-sensitive information (e.g., news about an event, sale offers), LLMs can help determine the validity or relevance period.

### 10. Accessibility Metadata:

Determine and tag content based on accessibility features, such as whether it's suitable for visually impaired readers, or if it contains media (like images, videos) with appropriate alt-texts.

### 11. Custom Metadata for Domain-Specific Applications:

For specialised domains, LLMs can be fine-tuned or trained to extract and generate metadata specific to that domain.

While LLMs can automate a significant portion of the metadata generation process, human oversight is recommended, especially for critical applications. It ensures that the metadata is accurate, relevant, and free from biases or errors that might originate from the model's training data.

### Alternatives

- A. Use an LLM fine-tuned with a metadata training dataset or by using the few-shots prompting technique.
- B. Use an LLM off the shelf.
- C. Asking LLMs to develop code to run metadata generation in other tools (e.g. R, Python).

# 6

## User Support and Assistance

### 6.1. Chatbots

LLMs have introduced a paradigm shift in the chatbot industry. Thanks to their advanced capabilities in understanding and generating human-like text, they have redefined what's achievable with chatbot interactions. Some practical aspects about the use of LLMs in chatbots are:

#### Capabilities:

**General Conversational Skills:** LLMs can handle a wide range of general conversations, making them suitable for chatbots designed for diverse user queries.

**Contextual Understanding:** LLMs, especially the newer iterations, can remember and reference earlier parts of a conversation to maintain context over a short interaction.

**Multilingual Interactions:** LLMs are trained on data from various languages, allowing them to handle conversations in multiple languages.

**Emotion and Sentiment Analysis:** They can gauge user sentiment to some extent, helping in tailoring responses as seen in a previous use case.

#### Customization and Fine-tuning

While base LLMs are versatile, for domain-specific applications, chatbots can be enhanced by fine-tuning LLMs on specific datasets, making them more relevant and accurate for niche tasks.

#### Integration with Other Systems

LLM-based chatbots can be integrated with databases, CRMs, or other software, allowing dynamic and personalised user interactions, such as querying account information, booking appointments, etc.

#### Limitations and Ethical Considerations

**Dependency on Training Data:** Their responses are based on their training data, which means they might sometimes provide inaccurate information.

**Inappropriate Content:** Chatbots can generate inappropriate, offensive, or harmful content if they haven't been adequately filtered or if they've learned from biased data. Also, the dataset used to train LLMs can sometimes include harmful biases or inappropriate language.

**Absence of True Understanding:** as seen in chapter 4.1.2, LLMs don't truly understand content in the way humans do; they may generate incorrect responses when assisting users.

**User Privacy:** Care must be taken to ensure user data isn't misused or stored without explicit consent.

#### Implementation and Deployment

Many platforms now allow integration of LLMs like GPT variants, making it easier for businesses and developers to deploy LLM-backed chatbots on websites, apps, or other platforms.

### Continuous Learning and Feedback Loops

Implementing a feedback loop where human agents can validate the chatbot's responses can help in improving accuracy and reliability over time.

### Cost and Scalability

Leveraging LLMs, especially cloud-hosted ones, might be cost-intensive for high-traffic applications. However, they offer scalability and consistent performance.

### Enhancing User Experience

With features like dynamic content generation, personalised responses, and the ability to handle open-ended queries, LLM-based chatbots can significantly enhance user engagement and satisfaction.

It's key to understand that when the bot is used externally on the statistical authority website for all users, it reflects the image of the institution. Therefore, the quality of service and ethical considerations are crucial factors that need careful attention. These issues can be tackled through various methods such as using reinforced learning, prompt priming, or even by using a different version of the Language Model (LLM) to oversee the bot's actions. On the other hand, when used only internally, these problems are less severe. Using the LLM as an assistant can greatly help in improving the quality and efficiency of the work done by employees and collaborators.

Alternatives

- A. LLM-based bots are the state-of-the-art solution for chatbots. They may even be able to query statistical databases and create visualisations. Implement a supervision strategy.
- B. Develop a chatbot offering generic user support (not domain specific).
- C. Use an LLM chatbot without taking into consideration its limitations and potential biases (e.g. harmful content, incorrect information, user privacy).

## 6.2. Multilingual Support

LLMs are designed to process and generate text in multiple languages, making them highly valuable for applications that require multilingual support. Some tasks LLMs can perform in this field are:

### Understanding Multiple Languages

LLMs are trained on vast datasets encompassing text from numerous languages. This enables them to comprehend and generate text across a wide variety of languages, from widely spoken ones like English, Spanish, and Mandarin to less common ones.

### Translation Assistance

While not a replacement for dedicated translation models, LLMs can provide reasonably accurate translations for many language pairs, especially for general content.

### Language Detection

Given a piece of text, LLMs can often detect the language it's written in, aiding in routing the text to appropriate processing pipelines or responding accordingly.

### Transliteration

LLMs can assist in converting text from one script to another, for instance, from Cyrillic to Latin script.

### Cultural Context

Due to the diverse range of data they're trained on, LLMs can often provide context or nuances about cultural idioms, phrases, or references in various languages.

### Localised Content Generation

LLMs can generate content tailored to specific languages or regions, allowing for localised marketing, content creation, and user engagement.

### Language Learning Assistance

LLMs can be used as language tutors, providing explanations, translations, and examples in different languages to learners.

### Handling Code-switching

Many users often mix languages within a single conversation or text (known as code-switching). LLMs can understand and respond appropriately in such scenarios.

### Limitations

**Accuracy Variances:** While LLMs support multiple languages, their proficiency might vary. They tend to

be more accurate in languages that have a substantial presence in their training data.

**Direct Translation Limitations:** For professional translation tasks, dedicated neural machine translation systems might be more appropriate.

### Alternatives

- A. Use an LLM for internal purposes to help refining a body of text; assure human validation
- B. Use a SLM locally installed for simple support tasks
- C. Use third-party LLM tools without proper validation; replace dedicated translation models with LLMs



# Automated Content Creation

Automated content creation using LLMs is a field with significant potential. The ability of models to generate human-like text has been harnessed across various sectors to create diverse types of content.

## Types of Content LLMs Can Generate

**Articles and Blog Posts:** LLMs can assist in drafting, proofreading, and optimising posts for maximum reach and clarity. LLMs can help craft informative and engaging posts that resonate with your target audience, reviewing contents to spot and correct any inaccuracies, ensuring that the information upholds the credibility of the statistical organisation. Also, since different social media platforms have unique algorithms, character limits, and user behaviours, LLMs can provide guidance on best practices for each platform, from ideal post lengths to the use of hashtags, to ensure that the content reaches the widest audience and achieves maximum engagement.

**Reports:** LLMs are proficient in analysing documents, collecting and organising pertinent data to ensure a coherent structure. They can also utilise or create custom templates tailored for specific report formats, enhancing the presentation and uniformity of the document. Once the data is integrated, LLMs may draft detailed content, ensuring the inclusion of necessary information, refining at the end the language and content for clarity and coherence.

**Summaries:** LLMs excel at condensing extensive pieces of text to capture their essence. By emphasising the main ideas and discarding any redundant information, it can provide clear and concise summaries.

## Benefits of Using LLMs for Content Creation

**Efficiency and Speed:** LLMs can generate content rapidly, significantly reducing the time it takes to produce drafts or complete pieces.

**Scalability:** LLMs can handle vast amounts of requests simultaneously, making them ideal for large-scale content needs.

**Consistency:** LLMs provide a consistent tone and style, ensuring uniformity across multiple pieces of content.

**Customization:** With the right prompts and guidance, LLMs can be tailored to produce content in specific styles or tones desired by users.

**Multilingual Capabilities:** Many LLMs support multiple languages, enabling content creation for diverse audiences.

**Versatility:** LLMs are versatile in handling a range of content types, from articles and reports to creative writing and technical documentation.

## Challenges and Considerations

**Quality Control:** While LLMs can produce coherent text, it doesn't guarantee accuracy or quality. Human review is often required.

**Authenticity and Originality:** Content generated might lack a unique voice or perspective. There's also a risk of unintentional plagiarism and copyright infringements.

**Ethical Concerns:** Misleading or false information generation, especially if LLMs are used without oversight.

**Over-reliance:** Relying solely on automated content might lead to generic or impersonal outputs, potentially affecting the institution credibility. Possible hallucinated content may also be incorrect, inappropriate or unsafe.

### Best Practices for Implementation

**Hybrid Approach:** Use LLMs for initial drafts or to accelerate content generation, followed by human review and refinement.

**Clear Guidelines:** Feed clear and specific instructions to the model to align outputs with desired outcomes.

**Regular Monitoring:** Continually monitor and adjust the model's outputs to ensure alignment with goals and standards.

**Avoid Sensitive Topics:** Until LLMs can be thoroughly vetted, avoid using them for critical or sensitive content, where inaccuracies might have significant repercussions.

**Security Concerns:** Ensure that sensitive data or text is not exposed to external LLM services. Using on-premises models or ensuring strict data sanitization practices is crucial.

LLMs bring unparalleled efficiency, scalability, and versatility, catering to diverse content needs across statistical organisations. Their ability to craft anything from in-depth reports to succinct summaries makes them an indispensable tool. However, like any tool, their utility is also marked by inherent challenges. Concerns regarding quality, authenticity, and potential ethical pitfalls necessitate judicious usage and consistent human oversight. For statistical organisations venturing into this domain, adopting a balanced, hybrid approach, combining the computational prowess of LLMs with the nuanced judgement of humans, can pave the way for content that is not only efficient but also resonates with authenticity and credibility.

#### Alternatives

- A. Use LLM generated contents as a starting point, then have real people review and polish the content. This way, an institution can get the best of both worlds: fast, automated content that still feels genuine and reliable.
- B. Produce contents via LLM and review and polish the text via another AI process, limiting the human intervention to the minimum.
- C. Publish automatically LLM generated contents without human intervention.



Part

**3**

**Use cases for scientists  
and developers**

# 8

## Software Development Support

LLMs have been increasingly integrated into the software development lifecycle due to their capability to understand and generate code, as well as their adeptness at processing human language.

Despite their remarkable capabilities, LLMs are not flawless and need careful consideration. Developers should exercise due diligence and thoroughly review the code generated or suggested by LLMs. Copyright law upholds the intellectual property of human authors, and the attribution of code produced by LLMs may pose complex legal issues. Also, security concerns arise when using external LLM services, and sensitive code or data could be compromised. Therefore, implementing on-premises models or implementing strict data sanitization practices is essential. Moreover, while LLMs can significantly enhance coding efficiency, over-reliance on them may hinder skill development and limit learning opportunities for developers.

LLMs can act as valuable assistants in the software development process, enhancing productivity, code quality, and the overall development experience. However, as with any tool, it's essential to use them judiciously and in conjunction with sound software development practices. Some use cases of LLMs for software development support are presented in the following subsections.

### 8.1. Documenting and Commenting

#### Documentation Assistance

Language Models can aid in generating extensive documentation, ensuring that code is well-documented and maintainable.

#### Code-to-Comment

Language Models can generate human-readable comments or summaries for chunks of code, aiding in understanding legacy code or unfamiliar projects.

Alternatives

- A. Critically assess the code for security and use a language model to help where the security level is appropriate. Comment automatically chunks of code and carefully insert only the comments to avoid code corruption.
- B. Except for niche programming languages, specifically train a model. The main LLMs are already trained for the most popular languages (Python, Java, Javascript, etc).
- C. Upload large fractions of code to an LLM for commenting and copy/paste the entire results in the codebase (possible code corruption/security issues etc..)

## 8.2. Coding and Scripting Guidance

### Code Generation

Given a high-level description, LLMs can assist in generating code snippets or even more extensive sections of code in various programming languages.

### Natural Language Queries

Developers can pose questions in natural language about specific coding challenges, libraries, or frameworks, and the LLM can provide answers or code examples in response.

### Learning and Onboarding

LLMs can be integrated into educational platforms, assisting new developers in understanding coding concepts, languages, or tools, offering examples and explanations dynamically.

### Integration with IDEs

By embedding LLMs into Integrated Development Environments, developers can get real-time suggestions, auto-completions, or even explanations as they code.

### Automating Repetitive Tasks

For routine or boilerplate tasks, developers can instruct the LLM to generate necessary code, increasing productivity.

### UI/UX Descriptions to Code

Given descriptions or mock-ups of user interfaces, LLMs can assist in generating the corresponding frontend code.

Various studies (Kalliamvakou, 2022; Jain, 2021; Murali, 2023; Tan, 2023) have demonstrated that significant increases in productivity and developer satisfaction can be achieved by using LLMs in software development.

In one of the studies (Kalliamvakou, 2022) related to GitHub Copilot, between 60–75% of users reported they feel more fulfilled with their job. Developers reported that Copilot helped them stay in the flow (73%) and preserve mental effort during repetitive tasks (87%).

In one experiment they measured two groups, one using Copilot, and the other group not using it, how successful each group was in completing the task and how long each group took to finish. The group that used Copilot had a higher rate of completing the task (78%, compared to

70% in the group without Copilot). The striking difference was that developers who used Copilot completed the task significantly faster, 55% faster than the developers who didn't use Copilot.

### Alternatives

- A. Use an IDE-integrated LLM solution to assist employees in software code generation. Carefully assess the legal implications of using LLM-generated code.
- B. Develop internal tools based on LLM to speed up code creation
- C. Ignoring the AI technological potentialities for code development

## 8.3. Testing and Debugging

### 8.3.1. Code Testing

LLMs can help in all phases of developing the testing structure for a software project:

#### Explaining Testing Concepts and Methodologies:

LLMs can help explain certain testing methodologies or concepts, it can provide explanations and examples, helping a practitioner understand and apply these methods effectively in your testing process.

**Test Case Generation:** LLMs can help a tester brainstorm and generate test cases. Developers can describe software's functionality, and LLMs can suggest a range of test scenarios, including edge cases not considered before.

**Writing Automated Test Scripts:** LLMs can assist in writing basic scripts for automated testing, especially using popular testing frameworks. Developers can describe the test scenario, and LLMs can help draft a script, which can be refined and integrated into the testing suite.

**Reviewing Test Plans and Documentation:** Developers can share test plans or documentation with LLMs for a review. This latter can provide feedback on clarity, coverage, and even suggest improvements or additional tests.

**Integration and Regression Testing:** In the context of integration and regression testing, LLMs can help outline a strategy for ensuring that new code integrates smoothly with existing code and doesn't introduce new bugs.

### 8.3.2. Code Review and Optimization

LLMs can be used to analyse code and suggest improvements, be it in terms of optimization, adherence to coding standards, or

even pointing out potential bugs. They can propose refactorings, or alternative methods to enhance code performance or readability.

### 8.3.3. Debugging Support

When presented with error messages or logs, LLMs can suggest potential solutions or causes for bugs, reducing debugging time.

Alternatives

- A. Use LLMs for code review and debugging
- B. Use LLMs for legacy code optimization instead of refactoring into new code
- C. Use LLMs automatically without human supervision.

## 8.4. Security Assistance

Using LLMs in the context of software code security can be a double-edged sword. On the one hand, LLMs can assist in identifying vulnerabilities, promoting best practices, and educating developers. On the other hand, they pose potential risks if used carelessly or without accepting that there might be some unforeseen risks, especially regarding sensitive code exposure or relying too heavily on them without human oversight.

### Advantages and Applications of LLMs for Code Security

**Vulnerability Detection:** LLMs can be trained to recognize common coding patterns that lead to vulnerabilities like SQL injection, cross-site scripting (XSS), buffer overflows, etc. They can then suggest secure coding alternatives.

**Code Review Assistance:** Integrated into the code review process, LLMs can highlight sections of code that potentially violate security best practices.

**Security Best Practices Promotion:** LLMs can be used to educate developers about security best practices in real-time, suggesting secure coding patterns as developers write or commit code.

**Automated Security Q&A:** Developers can query LLMs with specific security questions or scenarios, getting insights or solutions without diving deep into extensive documentation.

**Analysis of Logs and Alerts:** LLMs can assist in parsing and interpreting security logs or alerts, identifying potential threats or areas of concern.

### Risks and Considerations

**Over-reliance on LLMs:** Relying solely on LLMs for security can be risky. They should complement, not replace, established security practices and human expertise.

**Exposure of Sensitive Information:** If developers share proprietary or sensitive code with external LLM services, it poses a significant data leakage risk.

**Unintentional infringement of proprietary content:** The reused content produced by the external LLM may contain copyright-protected content included in the training data.

**LLMs' Limitations:** LLMs might not always catch subtle security vulnerabilities, especially novel or complex ones that haven't been adequately represented in their training data.

**Potential for Misinformation:** LLMs, if not fine-tuned properly or if vulnerabilities have been exploited for tampering, could occasionally provide incorrect or suboptimal security advice.

**Bias and Training Data:** If an LLM is trained predominantly on outdated, poisoned or insecure coding practices, it might perpetuate those practices. Regular updates and fine-tuning on recent, secure code datasets are essential.

**Data Privacy Concerns:** Especially in regulated industries, there might be concerns and legal implications related to sharing code or data with third-party LLM providers.

### Best Practices

**Internal Deployment:** Consider deploying LLMs internally rather than relying on cloud-based solutions to mitigate data leakage risks.

**Regular Audits:** Periodically review the advice and outputs provided by LLMs to ensure they align with current security standards.

**Combine with Other Tools:** Use LLMs in conjunction with established static and dynamic code analysis tools to ensure comprehensive security coverage.

**Education:** Ensure developers understand the capabilities and limitations of LLMs in the context of security.

Alternatives

- A. Using LLM as an assistant for vulnerabilities detection and code review. Do not rely exclusively on the results of the tool.
- B. Train a LLM/SLM for logs analysis and alerts, in conjunction with other traditional methods.
- C. Using LLM without supervision, automatically

## 8.5. Reverse Engineering of Legacy Code

LLMs can be used to analyse old code and extract business rules and requirements. Their advanced language understanding capabilities allow them to break down complex code structures and identify the underlying logic and intentions of the original developers.

By providing explanations of specific functions or blocks of code, LLMs can help developers understand the purpose of the code and identify potential improvements. It can also suggest modern equivalents or improvements to the old code, making it easier to update or refactor without losing the original intent. Additionally, LLMs' ability to understand and respond in natural language makes them an accessible tool for teams with varying technical expertise, even those without deep programming knowledge.

LLMs can cover different aspects when reverse engineering old software projects:

### Code Analysis Assistance

LLMs can help in interpreting and explaining complex code segments. By inputting specific lines or blocks of code, LLMs can provide insights into their functionality and purpose.

### Identification of Business Logic

When analysing sections of code that involve business rules, LLMs can assist in translating technical implementations into understandable business terms, facilitating the extraction of business logic.

### Clarifying Data Structures and Flows

For understanding data structures and how data flows through the system, LLMs can offer explanations and help in mapping out these processes.

### Generating Documentation

As deciphering the code, LLMs can aid in drafting clear and concise documentation, summarising the software's architecture, business rules, and data handling.

### Answering Queries and Providing Guidance

Throughout the reverse engineering process, LLMs can be a resource for answering queries, providing coding insights, and suggesting best practices or even alternative architectures for refactoring.

In the case of reverse engineering, LLMs act as support tools, aiding in interpretation, documentation, and understanding, rather than being a standalone solution.

Alternatives

- A. Create a process for reverse engineering as described above. Use LLMs in each phase according to the tasks to be carried out, using the system as an assistant.
- B. Letting the LLM define the reverse engineering process.
- C. Use LLM as a standalone solution for reverse engineering.

# 9

## Research, Education and Training

### 9.1. Research Assistance

LLMs have shown substantial promise as tools for research assistance across various academic and professional domains. The capacity of LLMs to understand, generate, and analyse text can be harnessed to accelerate and enrich the research process. LLMs can help researchers by doing:

#### Literature Review and Survey

**Search Queries:** LLMs can help craft more effective search queries to identify relevant literature.

**Summarization:** LLMs can provide concise summaries of long papers or articles, helping researchers quickly assess their relevance.

**Categorization:** They can assist in categorising found literature based on several criteria defined by the researcher

#### Drafting and Writing

**Content Generation:** LLMs can assist in drafting various content based on provided outlines or notes.

**Editing and Proofreading:** LLMs can help in grammar checking, style improvements, and ensuring clarity.

#### Citation and Reference Management

LLMs can assist in finding relevant citations, formatting them according to different citation styles, and ensuring that all references in the text are appropriately cited in the bibliography.

#### Question and Answer Assistance

Researchers can pose specific questions to the LLM related to their topic, methodologies, or challenges they're facing, and get detailed answers or suggestions.

#### Concept Explanation

Complex topics, methodologies, or theories can be better understood with the explanatory capabilities of LLMs.

#### Limitations

**Dependence on Training Data:** LLMs' knowledge is based on their training data, and they might not be aware of very recent advancements or niche topics. Some LLMs have internet browsing capabilities or even plugins to scan research papers databases (e.g. ChatGPT 4), which makes them a more complete tool for research tasks.

**Potential for plagiarism:** Researchers should be cautious as LLMs might generate content similar to existing literature. (see chapter 4.1 - Controversies and Limitations - Black box effect)

**Potential of copyright infringement:** Output, for example program code, may contain proprietary content.

LLMs offer a dynamic toolset that can significantly aid researchers in various phases of the research process. However, their outputs should be approached with a critical mindset, ensuring the originality, accuracy, and depth of the research work. Human oversight plays also an essential role: a best practice is to always ensure that a researcher reviews and verifies the content or suggestions provided by the LLM

#### Alternatives

- A. Using LLMs for references and research ideas. Always ensure human oversight. Prefer LLMs having internet browsing and research capabilities. Check the cut-off date for the training database.
- B. Use a LLM or a SLM for limited editing tasks
- C. Draft large parts of text with LLMs without clear, focused, prompting. Consider a LLM as a researcher rather than as a research tool.

## 9.2. Simulation and Scenario Analysis

LLMs can be valuable assets for scientists engaged in simulation and scenario analysis. It's essential to understand the strengths and limitations of LLMs in this context.

#### Exploring Complex Scenarios Through Narratives:

- **Scenario Description:** Scientists can describe a hypothetical scenario to an LLM and ask the model to predict or elaborate on potential outcomes based on its vast knowledge. This is particularly useful for social scientists or those in interdisciplinary studies where outcomes can be described in narrative formats. Particular attention must be drawn in case of numerical computations. As depicted in chapter 5.1, LLMs are not particularly good in dealing with numbers.
- **What-if Analysis:** LLMs can generate detailed narratives for a range of "what-if" scenarios, aiding scientists in understanding potential consequences or chain reactions in complex systems.

#### Review of Existing Simulations:

- **Interpretation:** Once a simulation has produced results, scientists can "discuss" these outcomes with LLMs to get a broader perspective or to understand potential real-world implications.
- **Comparative Analysis:** LLMs can compare the results of different scenarios or simulations based on provided data, offering insights into the implications of each.

#### Multidisciplinary Integration:

- **Bridging Disciplines:** For simulations that span multiple disciplines, LLMs can help integrate knowledge, ensuring coherence and comprehensiveness in scenario analysis.

#### Communication and Presentation:

- **Data Visualization:** While LLMs don't directly create visual content, they can suggest optimal ways to visualise data or guide scientists in choosing the right visualisation tools or libraries. As described in chapter 8.1, LLMs can provide efficient scripts to visualise data by using programming languages such as Python, R, Julia. ...
- **Report Generation:** LLMs can assist in drafting reports, executive summaries, or presentations, translating complex simulation outcomes into accessible content for varied audiences.

LLMs can't directly replace the tools explicitly designed for scientific simulations, but they can be valuable adjuncts. They help in scenario exploration, interpretation, and communication, especially when complex systems and multidisciplinary insights are involved.

#### Alternatives

- A. Using LLMs to create narratives for scenarios, code and complex analysis based on researchers inputs.
- B. Using LLMs for scenario setting without human supervision
- C. Elaborate numerical analysis with LLM without the medium of a programming language

## 9.3. Education and Training

LLMs have clearly shown their transformative potentials for the education and training sectors, reshaping how knowledge is delivered, understood, and engaged with. The application of LLMs in these domains is multifaceted:

#### Personalized Learning:

- **Adaptive Content Generation:** LLMs can generate learning content tailored to individual students' proficiency levels, learning styles, and interests.
- **Addressing Queries:** Students can pose questions to the LLM and get detailed, personalised answers, allowing for in-depth exploration of topics (also see chapter 6.1 - Chatbot)

**Professional Training:** In professional settings, LLMs can assist in generating training content, answering employee queries, or simulating stakeholders interactions for practice.

**Skill Development and Practical Tasks:** LLMs can guide learners through practical tasks, like coding exercises or lab experiments, offering step-by-step instructions and troubleshooting assistance.

**Assessments and Feedback:** LLMs can be used to design quizzes, tests, or assignments. Furthermore, they can evaluate submissions and provide constructive feedback.

**Continual Learning and Updates:** In rapidly evolving fields, LLMs can provide recent updates, research summaries, or industry trends to keep learners at the forefront of knowledge.

**Interactive eBooks and Study Materials:** Embedding LLMs in digital study materials allows for dynamic content generation, instant clarifications, and interactive learning experiences.

**Tutorials and How-to Guides:** For well-documented skills or tasks, LLMs can generate step-by-step guides.

One of the best practices in this field is to use a hybrid approach, i.e. combining LLM-facilitated learning with traditional educational methods, ensuring a balance between technology and human touch. Regular monitoring

is also critical: periodically reviewing and validating the content and feedback provided by LLMs would maintain a high quality level of the production.

LLMs offer innovative tools for reshaping education and training, making learning more personalised, accessible, and interactive. However, they should be integrated judiciously, complementing rather than replacing the invaluable human elements of education.

Alternatives

- A. Use LLMs to assist professionals for the scopes mentioned above
- B. Due to the probabilistic nature of LLMs, exercising caution is essential when employing them as a rating tool in assessments. Without human oversight, the results may lack consistency and uniformity.
- C. Use LLMs to create educational content without human supervision



Part

**4**

**The New AI-powered  
(Data) Scientist**

*"The only constant in life is change"*

Heraclitus

Probably Heraclitus' sentence perfectly captures today's fast-paced technological environment. The rapid emergence of new technologies, especially in Generative AI, challenges the ability to keep pace with these advancements. This calls for an evolution in learning methods, especially for those pursuing scientific careers. Individuals spend years in educational settings only to find the real work environment vastly different. Here LLMs are indispensable, offering real-time information and analytical capabilities, serving as vital assistants for adapting to these dynamic professional landscapes. The acceleration of innovation, particularly in generative AI, is reaching a velocity that is increasingly difficult for humans to manage alone.

In the realm of research, the influence of generative AI is profound and transformative. Traditional, slower research methods are being replaced by faster, AI-driven approaches. These technologies are not just tools; they are reshaping the whole approach to scientific exploration.

For example, in the field of drug discovery, generative AI models can predict molecular interactions much faster than traditional methods, significantly speeding up the development of new medications. The AI tool GNoME (Merchant, 2023), created by DeepMind team at Google, finds 2.2 million new crystals, including 380,000 stable materials that could power future technologies. *This would be equivalent of about 800 years' worth of knowledge.*

As shown previously, Generative AI significantly enhances the capabilities of data scientists addressing, for instance, the challenge of limited datasets and enriching model training and testing. It accelerates the data analysis process, allowing data scientists to explore and evaluate a wider range of hypotheses efficiently, leading to deeper and more comprehensive insights. This AI-driven approach not only expands the quantity but also the quality of data analysis, revealing intricate patterns that might be overlooked in manual analysis. Additionally, generative AI automates the generation and testing of software code for analytical models, enabling data scientists to concentrate on higher-level tasks such as identifying and clarifying problems and evaluating solutions. Ultimately, it fosters improved decision-making by autonomously uncovering hidden patterns and insights, and generating accurate, evolving business reports, thereby revolutionizing the field of data science and analytics.

This rapid evolution demands a significant shift in mindset for researchers and professionals. Being adaptable, continuously learning, and embracing new methods are now essential traits. The role of lifelong learning becomes crucial as professionals must regularly update their skills to remain relevant. Interdisciplinary approaches are also key, as generative AI creates intersections between various fields, opening new avenues for innovation. The pace of innovation in generative AI emphasizes the need for smart agents to assist humans in processing the vast amount of information and developments, underscoring the synergy between human intelligence and artificial assistance in navigating the future.

# 10

## A Team of One

Thanks to the technological progress in the field of AI it's possible now -pushing the concept to the extreme- to imagine a research team comprising a single human who acts as the primary decision-maker and lead researcher, supported by a multitude of specialised GPTs, each one focusing on a specific aspect of the research area. These advanced machines might possess avatars that are virtually indistinguishable from real humans, complete with authentic-looking faces and voices, allowing for direct verbal interaction.

Let's put together some parts of the puzzle:

### Item One: Auto-GPT

Auto-GPT (AutoGPT, 2023) is an open-source, autonomous AI agent developed by Toran Bruce Richards from Significant Gravitas Ltd. Launched on March 30, 2023, it leverages OpenAI's GPT-4 or GPT-3.5 APIs to autonomously perform tasks by breaking them into sub-tasks, using the internet and other tools. It is among the first applications to use GPT-4 for autonomous tasks. Users define the agent's name, role, and objective, along with up to five methods to achieve the objective, after which Auto-GPT operates independently.

This AI agent is capable of breaking down large tasks into sub-tasks, chaining them together to achieve complex objectives. It features internet connectivity for updated information retrieval, short-term memory for task context, file organization capabilities, and is multimodal, handling both text and image inputs. These capabilities enable Auto-GPT to automate workflows, analyse data, and generate suggestions.

Auto-GPT has diverse applications in software development, business, and other fields. It can develop software, debug code, conduct market research, write product reviews, and create business plans. Notably, it has been used to create ChefGPT for generating unique recipes and ChaosGPT, an AI agent with controversial objectives.

Publicly available on GitHub, Auto-GPT requires installation in a development environment like Docker and registration with an OpenAI API key. In October 2023, it raised \$12 million from investors, highlighting its significant potential and impact.

### Item Two: ChatDev

ChatDev is a virtual software company operated by various intelligent agents fulfilling roles like CEO, CPO, CTO, programmer, reviewer, tester, and art designer. These agents form an organizational structure with the mission to "revolutionize the digital world through programming." They collaborate in specialized functions like design, coding, testing, and documentation. ChatDev aims to provide an easy-to-use, highly customizable framework based on large language models (LLMs) for studying collective intelligence.

Some features of the project:

- Experiential Co-Learning is a method for agents to accumulate experiences for solving new tasks efficiently, reducing repetitive errors.
- Incremental development enables agents to build upon existing codes.
- Integration of Docker for safe execution.
- Availability of Git mode for version control.
- Human-Agent-Interaction mode allows users to contribute as reviewers.

- Art mode enables the designer agent to generate images for software.

### Item Three: Unlocking insights in scientific literature.

Gemini (Google - Gemini Team, 2023), the new LLM model from Google was used by scientists at Google DeepMind to streamline the extraction of data from vast amounts of scientific literature. Traditionally, this involved labor-intensive manual searches and data extraction from thousands of papers. Gemini's advanced understanding of scientific content significantly reduces this workload.

A case study is presented where Gemini updated a genetics dataset by processing over 200,000 new papers in few hours, a task impractical to do manually. Gemini efficiently filtered relevant papers and extracted key data. It also showcased its multimodal capabilities by updating a graph from the original study using a new dataset.

Gemini's utility extends beyond science to any field dealing with large datasets, including law and finance. This is the how very large language models like Gemini have the potential to revolutionize data handling in various domains.

### Item Four: GNoME.

The deep neural network specialized in inventing new crystals GNoME (Merchant, 2023), already mentioned before, was used by scientists to create new formulas for numerous types of crystals. What was not mentioned is that this AI tool is also connected to a robot, operating in the

real world, that has the ingredients and is able to synthesize the formulas.

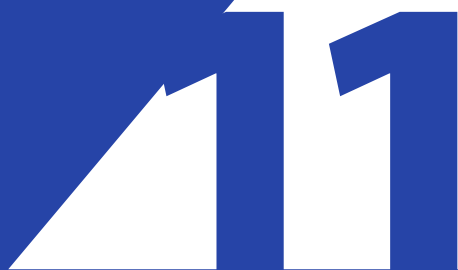
The four examples above highlight the evolving landscape of LLMs applications, emphasizing the collaboration of specialized large language models (LLMs), their autonomous operational capabilities, and their interaction with both digital and physical realms. Specifically:

**Inter-LLM Collaboration and Specialization:** Various LLMs, each fine-tuned for specific tasks, demonstrate an ability to collaborate effectively. This mix enhances problem-solving efficiency and versatility.

**Interaction with Physical and Digital Environments:** A critical advancement is the AI systems' capability to interact with both digital interfaces and physical devices. Notably, OpenAI has developed an architecture for GPT-4 that allows the integration of plugins, enabling tasks such as browsing research databases, summarization in PDF files, and, why not, robotic manipulation.

**Autonomy in Operations:** These AI applications exhibit significant autonomy. They can function independently or under human guidance, adapting to different levels of supervision. This can be found also in GPT-4 environment, which can operate in a Python sandbox for data analysis, demonstrating its ability to use tools autonomously.

This integration of specialized LLMs, autonomous operation, and interaction with real-world environments, enables the formation of smaller, more efficient research teams. AI systems are no longer mere tools but are evolving into collaborative partners capable of undertaking diverse tasks across digital and physical domains.



## LLM Skills and New AI Jobs

As shown in Chapter 4.1 (Controversies and Limitations), LLMs are not capable of “understanding”: even the most advanced ones, even if excelling in specific tasks, are lacking the broad, adaptable intelligence of AGI, which remains a largely theoretical concept as of now. LLMs will not replace humans, at least not entirely. But in order to make a rigorous analysis we need more data and especially a systematic approach to see for each occupation, or occupation group, how those latter are affected by AI, in our case by LLMs.

How to measure AI Occupational Exposure (**AIOE**)? An approach is presented in a research done at Princeton University (DeepLearningAI, 2023b, Felten, 2023): A team led by Ed Felten, a computer scientist at Princeton University, projected the jobs and industries in the U.S. likely to be most affected by language models.

How it works: The authors calculated an “exposure” score for each of 774 occupations and 115 industries by comparing human skills to AI application areas. For the purpose of the study, exposure is neither positive nor negative; it’s a measure of how likely a job or industry

would change in response to developments in language processing.

- The authors used a U.S. Department of Labor database that describes each occupation in terms of 52 human abilities. Such abilities include dynamic strength, hearing sensitivity, mathematical reasoning, and written expression, and they’re weighted according to their importance to a given occupation.
- Crowdsourced workers scored the relevance of language modelling to each ability. For instance, language modelling has little relevance to dynamic strength but great relevance to written expression.
- The authors used the scores as weighted variables in an equation that aggregated the relevance of language modelling to the human abilities involved in each occupation and industry.”

Ed Felten and his team conducted a similar research two years earlier (Felten, 2021), before the “LLM revolution”, and found some different results, as can be seen in the tables below (Table 14 – Table 15).

**TABLE 14**

**Top 20 Occupations Exposed to AI, Before and After LLMs, adaptation from “How will Language Modelers like ChatGPT Affect Occupations and Industries?” (Felten, 2023)**

	Original AIOE (before LLMs) - 19 April 2021	After LLMs - 18 Mar 2023	Change
1	Genetic counselors	Telemarketers	new
2	Financial examiners	English Language and Literature Teachers, Postsecondary	new
3	Actuaries	Foreign Language and Literature Teachers, Postsecondary	new
4	Purchasing agents	History Teachers, Postsecondary	+12
5	Budget analysts	Law Teachers, Postsecondary	new
6	Judges, magistrate judges, and magistrates	Philosophy and Religion Teachers, Postsecondary	new
7	Procurement clerks	Sociology Teachers, Postsecondary	new
8	Accountants and auditors	Political Science Teachers, Postsecondary	new
9	Mathematicians	Criminal Justice and Law Enforcement Teachers, Postsecondary	new
10	Judicial law clerks	Sociologists	new
11	Education administrators, postsecondary	Social Work Teachers, Postsecondary	new
12	Clinical, counseling, and school psychologists	Psychology Teachers, Postsecondary	new
13	Financial managers	Communications Teachers, Postsecondary	new
14	Compensation, benefits, and job analysis specialists	Political Scientists	new
15	Credit authorizers, checkers, and clerks	Area, Ethnic, and Cultural Studies Teachers, Postsecondary	new
16	History teachers, postsecondary	Arbitrators, Mediators, and Conciliators	+4
17	Geographers	Judges, Magistrate Judges, and Magistrates	-11
18	Epidemiologists	Geography Teachers, Postsecondary	new
19	Management analysts	Library Science Teachers, Postsecondary	new
20	Arbitrators, mediators, and conciliators	Clinical, Counseling, and School Psychologists	-8

**TABLE 15**

### Top 20 Industries Exposed to AI, Before and After LLMs, adaptation from “How will Language Modelers like ChatGPT Affect Occupations and Industries?” (Felten, 2023)

	Original AIOE (before LLMs) - 19 April 2021	After LLMs - 18 Mar 2023	Change
1	Securities, commodity contracts, and other financial investments and related activities	Legal services	+3
2	Accounting, tax preparation, bookkeeping, and payroll services	Securities, commodity contracts, and other financial investments and related activities	-1
3	Insurance and employee benefit funds	Agencies, brokerages, and other insurance related activities	+2
4	Legal services	Insurance and employee benefit funds	-1
5	Agencies, brokerages, and other insurance related activities	Nondepository credit intermediation	+1
6	Nondepository credit intermediation	Agents and managers for artists, athletes, entertainers, and other public figures	+5
7	Other investment pools and funds	Insurance carriers	+1
8	Insurance carriers	Other investment pools and funds	-1
9	Software publishers	Accounting, tax preparation, bookkeeping, and payroll services	-7
10	Lessors of nonfinancial intangible assets	Business Support Services	new
11	Agents and managers for artists, athletes, entertainers, and other public figures	Software publishers	-2
12	Credit intermediation and related activities	Lessors of nonfinancial intangible assets	-2
13	Computer systems design and related services	Business schools and computer and management training	+6
14	Management, scientific, and technical consulting services	Credit intermediation and related activities	-2
14	Monetary authorities-central Bank	Grantmaking and giving services	+6
16	Office administrative services	Travel Arrangement and Reservation Services	new
17	Other information services	Junior Colleges	new
18	Data processing, hosting, and related services	Computer systems design and related services	-6
19	Business schools and computer and management training	Management, scientific, and technical consulting services	-5
20	Grantmaking and giving services	Other information services	-3

The lists are quite different and primarily demonstrate how advancements made in just one AI technology have significantly influenced the forecasts for the past two years.

As the authors of the article mentions also, an occupation's exposure to AI does not necessarily put jobs at risk. History suggests the opposite can happen. A 2022 study (US-BLS,

2022) found that occupations exposed to automation saw increases in employment between 2008 and 2018.

Another very relevant recent study about LLMs impact on the jobs market was released by World Economic Forum (WEF, 2023) with the very important distinction between potential for automation and potential for augmentation. They identified that from 19,265 tasks, and 9,934 language based tasks, there are a good 3,211 with augmentation potential and 1,435 candidates for automation.

## 11.1. New AI Jobs

It is almost impossible to predict what this explosion of LLM usage will add to the map of occupations. For sure the impact on existing jobs will be massive. Goldman Sachs issued a report on the effects of Generative AI on economic growth (Goldman Sachs, 2023) saying: *“If generative AI delivers on its promised capabilities, the labour market could face significant disruption,”* estimating that 300 million jobs could be impacted.

### “Direct” New Jobs

Some new jobs will be as a direct effect of having these new tools created, trained, deployed and maintained. Some jobs already exist but now will have a more narrow specialisation in LLM domain. As more and more occupations develop solutions powered by LLMs the need for AI Model Trainer and Tuner, Data Curators and Labelers, AI Auditor, Ethics and Governance Specialists, AI Integration Specialists, Language Model Security Analysts, Interface and Interaction Designers, and many more. Also, for existing jobs like AI Research Scientists and Engineers the demand is skyrocketing.

### “Indirect” New Jobs

The already most popular new jobs are in the realm of prompting with prompt engineer and prompt architect. Other jobs can be AI-Personality Designer, AI-Human Teaming Coordinator, etc.

## 11.2. A Simplified Method to Measure LLM Impact on Jobs

The method by Ed Felten (Felten, 2023) shown before can be very laborious and may require a lot of resources. An easier and more customizable approach can be found in a course (DeepLearningAI, 2023c) made by Andrew Ng, founder of DeepLearning.AI and co-founder of Coursera. It is based on the public resource Onet (ONET, 2023) where information can be found for more than 900 occupations, covering but not limited to the description, qualifications, salary range, and tasks performed. The task list can be adapted to better match specific research institutions.

The method is structured into three phases: defining the **tasks**, establishing **scoring** criteria, and evaluating the **results**. To illustrate its effectiveness, this method is applied to a data science job in the following sections.

### 11.2.1. Tasks

The ONet site was used as an entry point. The site provides the tasks performed by a data science professional and includes analysing, manipulating, or processing large sets of data, creating graphs, charts, or other visualisations, etc. The list with 17 distinct tasks will be used to score the exposure. An important note: even if the average is low there can be tasks that are very probable to be exposed and enhanced.

In other environments, the job description can differ wildly but this site could be a very good starting point for an analysis. One can add or remove tasks from the list according to its specifics.

The tasks list is important because LLMs are just a generic technology that can be used to assist humans in a natural language context. Splitting a job by performing tasks gives a granular and more clear picture of what can be improved (assisted) by LLMs and how.

### 11.2.2. Scoring

For each task, a numerical value is assigned to represent the likelihood of it being assisted or replaced by a LLM-powered application. The scale ranges from 1 (less probable) to 5 (very probable). The criteria used for scoring were derived from the discussions in the chapters: 2.10 Costs and Benefits and 4.1 Controversies and Limitations.



**TABLE 16****Example: Scoring for data scientist job as per task description from ONet**

Task	Score
Analyse, manipulate, or process large sets of data using statistical software.	4
Apply feature selection algorithms to models predicting outcomes of interest, such as sales, attrition, and healthcare use.	4
Apply sampling techniques to determine groups to be surveyed or use complete enumeration methods.	4
Clean and manipulate raw data using statistical software.	4
Compare models using statistical performance metrics, such as loss functions or proportion of explained variance.	4
Create graphs, charts, or other visualisations to convey the results of data analysis using specialised software.	4
Deliver oral or written presentations of the results of mathematical modelling and data analysis to management or other end users.	2
Design surveys, opinion polls, or other instruments to collect data.	3
Design surveys, opinion polls, or other instruments to collect data.	3
Identify business problems or management objectives that can be addressed through data analysis.	2
Identify relationships and trends or any factors that could affect the results of research.	2
Identify solutions to business problems, such as budgeting, staffing, and marketing decisions, using the results of data analysis.	2
Propose solutions in engineering, the sciences, and other fields using mathematical theories and techniques.	2
Read scientific articles, conference papers, or other sources of research to identify emerging analytic trends and technologies.	4
Recommend data-driven solutions to key stakeholders.	2
Test, validate, and reformulate models to ensure accurate prediction of outcomes of interest.	2
Write new functions or applications in programming languages to conduct analyses.	4

**11.2.3. Results**

Using ONet task split and the scores calculated before it is clearly visible that LLMs can be integrated in the work environment of a data scientist with important benefits. An LLM would be an excellent assistant in proposing methods and procedures for analysis, as well as in devising or enhancing them. LLM can suggest new tools and summarise and analyse huge amounts of scientific research in search for trends and new developments.

LLM can suggest new domains or areas of research or improvements.

According to Felten (2023) the large majority of knowledge based jobs can be very much enhanced and/or assisted by LLM-powered tools. By streamlining scientific workflows and freeing up (data) scientists' time, LLMs hold the promise of becoming an indispensable asset in various research domains, leading to novel and groundbreaking discoveries

# Conclusions

A few years ago, before the breakthrough in Generative AI and LLMs, one of the most significant developments was in the field of convolutional networks for image processing. A researcher specialising in radiology predicted that within five years, radiologists would become obsolete, replaced by machines. (Hinton, 2016).

After 4 years, another researcher argued that AI would not have replaced radiologists, but rather that it would have augmented their work, making them more efficient and effective. He stated that “radiologists who use AI will replace radiologists who don’t.” (Langlotz, 2019)

It’s evident that the 2016 prediction was incorrect, or at least imprecise, as the industry still greatly needs qualified human radiologists. The 2019 prediction was more accurate, as it holds true not only for radiologists but can also undoubtedly be generalised to the entire knowledge-based job sector.

The evolution from LLMs to more comprehensive Generative AI systems marks a significant milestone in the field of artificial intelligence. These systems, as exemplified by LLMs, are increasingly becoming integral components of multimodal generative systems, capable of processing and generating complex data across various modalities. This evolution is not just a technological leap but is also reshaping industry trends, as highlighted by notable sources like Forbes (Marr, 2023a) (Marr, 2023b) and Gartner (Perri, 2023).

The rapid pace of innovation, and its diffusion into industry sectors are staggering, underscoring the need for thoughtful legislation. The eagerly anticipated EU AI Act (EU Council, 2023) is expected to play a pivotal role in this regard, setting standards and guidelines that could shape the future of AI development and deployment.

Moreover, this technological shift is not only transforming existing industries but also creating new job roles. The emergence of positions like AI Manager, prompt engineer or AI auditor exemplifies this change. These roles are crucial

in bridging the gap between human expertise and AI capabilities, as seen previously in Part 4 (Team of One).

Data scientists and statisticians are thus entering a new phase of data analysis and interpretation, fundamentally transforming their roles and methodologies. Data professionals are increasingly at the forefront of leveraging AI to extract deeper insights and detect patterns from vast and complex datasets, a challenge that was previously overwhelming due to the volume and intricacy of the data involved. The accuracy and efficiency introduced by modern AI systems enable data scientists and statisticians to go beyond conventional analytics, facilitating predictive modelling and decision-making with an accuracy that was previously unattainable. Moreover, the progression of AI is driving the creation of new statistical tools, specifically designed to manage and make sense of the multidimensional data produced by these systems. This shift not only extends the analytical capabilities of data experts but also necessitates a mindset geared towards continual learning to stay abreast of the swiftly evolving AI field. As AI-powered tools become more integrated into various sectors, the role of data scientists and statisticians is expected to transition from basic data analysis to becoming strategic consultants, guiding statistical offices in employing data-driven strategies for advancement and innovation. In this period of transformation, their profound knowledge is crucial in ensuring ethical and responsible AI usage, establishing them as pivotal guardians of data integrity and governance in today’s digital era.

As we stand at this juncture, it is clear that LLMs and Generative AI are not just tools of the present but catalysts for a future where the integration of AI in daily life of statistical offices is seamless and intuitive. The implications are profound, spanning ethical, economic, and social dimensions. The journey ahead is one of cautious optimism, as we navigate these uncharted waters with a sense of responsibility and a vision for a future where AI enhances human potential and fosters a more connected, efficient world.

# References

- AI Alliance (2023): Launches as an International Community of Leading Technology Developers, Researchers, and Adopters Collaborating Together to Advance Open, Safe, Responsible AI, IBM Newsroom, retrieved at: <https://newsroom.ibm.com/AI-Alliance-Launches-as-an-International-Community-of-Leading-Technology-Developers,-Researchers,-and-Adopters-Collaborating-Together-to-Advance-Open,-Safe,-Responsible-AI>.
- Achiam, Josh, et al, GPT-4 Technical Report, ArXiv, /abs/2303.08774, 2023.
- Alperovich, Galina, The Secret Sauce behind 100K context window in LLMs: all tricks in one place, Medium, GoPenAI. Retrieved on 19.12.2023 at: <https://blog.gopenai.com/how-to-speed-up-llms-and-use-100k-context-window-all-tricks-in-one-place-ffd40577b4c>.
- Amodei, Dario, et al., Concrete Problems in AI Safety, ArXiv, /abs/1606.06565, 2016.
- Anthropic (2023): Introducing Claude. Retrieved on 19.12.2023 at: <https://www.anthropic.com/index/introducing-claude>.
- Anthropic (2023): Introducing 100K Context Windows by Anthropic. Retrieved on 19.12.2023 at: <https://www.anthropic.com/index/100k-context-windows>.
- Australian Government – Digital Transformation Agency (2023), Interim guidance on government use of public generative AI tools - November 2023, retrieved on December 2023 at: <https://architecture.digital.gov.au/guidance-generative-ai>
- AutoGPT (2023): Automate GPT, <https://github.com/Significant-Gravitas/AutoGPT>
- Bahdanau, Dzmitry, et al., Neural Machine Translation by Jointly Learning to Align and Translate, ArXiv, /abs/1409.0473, 2014.
- Bahri, Yasaman, et al., Explaining Neural Scaling Laws, ArXiv, /abs/2102.06701, 2021.
- Bai, Yuntao, et al., Training a Helpful and Harmless Assistant with Reinforcement Learning from Human Feedback, ArXiv, /abs/2204.05862, 2022.
- Bai, Yuntao, et al., Constitutional AI: Harmlessness from AI Feedback, ArXiv, /abs/2212.08073, 2022.
- Bailey, Luke, et al., Soft prompting might be a bug, not a feature, 2023.
- Barham, P., Chen, J., Davis, R., Hazzard, A., Ho, G., Kang, J., Wawrzyniek, J., Pathways: Asynchronous Distributed Dataflow for ML, ArXiv, /abs/2203.12533, 2022.
- Barrett, Clark, et al., Identifying and Mitigating the Security Risks of Generative AI, ArXiv, /abs/2308.14840, 2023.
- Bender, Emily M., et al., On the dangers of stochastic parrots: Can language models be too big?, Proceedings of the 2021 ACM conference on fairness, accountability, and transparency, 2021.
- Bloomberg (2023): Introducing BloombergGPT, Bloomberg's 50-billion parameter large language model, purpose-built from scratch for finance, Press, Bloomberg LP. Retrieved on 19.12.2023 at: <https://www.bloomberg.com/company/press/bloomberggpt-50-billion-parameter-llm-tuned-finance/>.
- Bostrom, Nick. Superintelligence: Paths, Dangers, Strategies. UK, Oxford University Press, 2014.
- Brown, Tom B., et al., Language Models Are Few-Shot Learners, ArXiv, /abs/2005.14165, 2020.
- Bubeck, S., Chandrasekaran, V., Eldan, R., Gehrke, J., Horvitz, E., Kamar, E., Lee, P., Lee, Y. T., Li, Y., Lundberg, S., Nori, H., Palangi, H., Ribeiro, M. T., Zhang, Y., Sparks of Artificial General Intelligence: Early experiments with GPT-4, ArXiv. /abs/2303.12712, 2023.
- Cellan-Jones, R. (2014). Stephen Hawking warns artificial intelligence could end mankind. BBC News. Retrieved on 01/12/2023 at: <https://www.bbc.com/news/technology-30290540>

- Chen, Canyu, Kai Shu, Combating Misinformation in the Age of LLMs: Opportunities and Challenges, ArXiv, /abs/2311.05656, 2023.
- Cho, Kyunghyun, et al., Learning Phrase Representations Using RNN Encoder-Decoder for Statistical Machine Translation, ArXiv, /abs/1406.1078, 2014.
- Chomsky, Noam, Ian Roberts, and Jeffrey Watumull. "Noam Chomsky: The False Promise of ChatGPT." The New York Times 8 (2023). Retrieved on 01/12/2023 at: <https://www.nytimes.com/2023/03/08/opinion/noam-chomsky-chatgpt-ai.html>
- Chung, Hyung W., et al., Scaling Instruction-Finetuned Language Models, ArXiv, /abs/2210.11416, 2022.
- DeepLearningAI (2023): A Complete Guide to Natural Language Processing, Course. Retrieved on 19.12.2023 at: <https://www.deeplearning.ai/resources/natural-language-processing/>.
- DeepLearningAI (2023), The Batch - Language Models' Impact on Jobs. Retrieved on 19.12.2023 at: <https://www.deeplearning.ai/the-batch/the-occupations-likely-to-be-most-affected-by-language-models/>.
- DeepLearningAI (2023): Course - Generative AI for Everyone. Retrieved on 19.12.2023 at: <https://www.deeplearning.ai/courses/generative-ai-for-everyone/>.
- Devlin, Jacob, et al., BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, ArXiv, /abs/1810.04805, 2018.
- Council of the EU (2023): Artificial intelligence act: Council and Parliament strike a deal on the first rules for AI in the world, Press release. Retrieved on 19.12.2023 at: <https://www.consilium.europa.eu/en/press/press-releases/2023/12/09/artificial-intelligence-act-council-and-parliament-strike-a-deal-on-the-first-worldwide-rules-for-ai>.
- EleutherAI (2023): Language Model Evaluation Harness. Retrieved on 19.12.2023 at: <https://github.com/EleutherAI/lm-evaluation-harness>.
- Elmahdy, Adel, Ahmed Salem, Deconstructing Classifiers: Towards A Data Reconstruction Attack Against Text Classification Models, ArXiv, /abs/2306.13789, 2023.
- Elman, J. L., Finding structure in time, Cognitive science, 14(2), p. 179-211, 1990.
- European Parliament (2023): EUBERT is a pretrained BERT model that leverages a substantial corpus of documents from the European Publications Office. Retrieved on 19.12.2023 at: <https://huggingface.co/EuropeanParliament/EUBERT>.
- Felten, Edward, et al., Occupational, Industry, and Geographic Exposure to Artificial Intelligence: A Novel Dataset and Its Potential Uses, Strategic Management Journal, vol. 42, no. 12, p. 2195-2217, 2021.
- Felten, Ed, et al., How Will Language Modelers like ChatGPT Affect Occupations and Industries?, ArXiv, /abs/2303.01157, 2023.
- Fursov, Ivan, et al., A Differentiable Language Model Adversarial Attack on Text Classifiers, ArXiv, /abs/2107.11275, 2021.
- Goldman Sachs (2023): Generative AI could raise global GDP by 7%. Retrieved on 19.12.2023 at: <https://www.goldmansachs.com/intelligence/pages/generative-ai-could-raise-global-gdp-by-7-percent.html>.
- Gemini Team, Google, Gemini: A Family of Highly Capable Multimodal Models. Retrieved on 19.12.2023 at: [https://storage.googleapis.com/deepmind-media/gemini/gemini\\_1\\_report.pdf](https://storage.googleapis.com/deepmind-media/gemini/gemini_1_report.pdf).
- Google Cloud (2023): AI and machine learning products. Retrieved on 19.12.2023 at: <https://cloud.google.com/products/ai>.
- Gu, Tianyu, Siddharth Garg., BadNets: Identifying Vulnerabilities in the Machine Learning Model Supply Chain, ArXiv, /abs/1708.06733, 2017.
- Harari, Yuval N. "Reboot for the AI Revolution." Nature, vol. 550, no. 7676, 2017, pp. 324-327, <https://doi.org/10.1038/550324a>.
- Hendrycks, Dan, et al., Measuring Massive Multitask Language Understanding, ArXiv, /abs/2009.03300, 2020.
- Hinton, G. E. Why we should stop training radiologists. In Machine Learning and Market for Intelligence Conference (pp. 1-3), 2016.
- Hochreiter, S., Schmidhuber, J., Long short-term memory, Neural computation, 9(8), p. 1735-1780, 1997.
- Hoffmann, Jordan, et al., Training Compute-Optimal Large Language Models, ArXiv, /abs/2203.15556, 2022.
- Hu, Edward J., et al., LoRA: Low-Rank Adaptation of Large Language Models, ArXiv, /abs/2106.09685, 2021.
- Hugging Face (2023): The AI community building the future. Retrieved on 19.12.2023 at: <https://huggingface.co/>.
- Isaev, Mikhail, Nic McDonald, and Richard Vuduc, Scaling Infrastructure to Support Multi-Trillion Parameter LLM Training, Architecture and System Support for Transformer Models (ASSYST@ ISCA 2023), 2023.

- Jain, Naman, et al., Jigsaw: Large Language Models Meet Program Synthesis, ArXiv, /abs/2112.02969, 2021.
- Jarmul, Katharine, Privacy Attacks on Machine Learning Models, InfoQ. Retrieved on 19.12.2023 at: <https://www.infoq.com/articles/privacy-attacks-machine-learning-models/>, 2019.
- Järvinen, Petteri, Looks impressive and the formulas are correct, but the end result is completely wrong (translation). Retrieved on 19.12.2023 at: <https://twitter.com/petterij/status/1691016914812760064>.
- Kalliamvakou, Eirini, Research: quantifying GitHub Copilot's impact on developer productivity and happiness, blog. Retrieved on 19.12.2023 at: <https://github.blog/2022-09-07-research-quantifying-github-copilots-impact-on-developer-productivity-and-happiness/>, 2022.
- Kaplan, Jared, et al., Scaling Laws for Neural Language Models, ArXiv, /abs/2001.08361, 2020.
- Karpathy, Andrej, Talk: Intro to Large Language Models. Retrieved on 19.12.2023 at: [https://www.youtube.com/watch?v=zjkBMFhNj\\_g](https://www.youtube.com/watch?v=zjkBMFhNj_g), 2023.
- Kim, Yoon., Convolutional Neural Networks for Sentence Classification, ArXiv, /abs/1408.58822, 2014.
- Kirchenbauer, John, et al., On the Reliability of Watermarks for Large Language Models, ArXiv, /abs/2306.04634, 2023.
- Kitaev, Nikita, Łukasz Kaiser, Anselm Levskaya., Reformer: The efficient transformer, ArXiv, /abs/2001.04451, 2020.
- Kurzweil, R., The Singularity is Near: When Humans Transcend Biology, Viking, 9780670033843, 2005. Retrieved on 19.12.2023 at: <https://books.google.lu/books?id=88U6hdUi6D0C>.
- Kurzweil, Ray. The Singularity is Nearer: When We Merge with Computers. Random House, 2024. Not yet published.
- Langlotz, C. P. Editorial: Will artificial intelligence replace radiologists? Radiology, 292(1), 4-8, 2019.
- Lester, Brian, Noah Constant, The Power of Scale for Parameter-Efficient Prompt Tuning, ArXiv, /abs/2104.08691, 2021.
- Lewis, Patrick, et al., Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks, ArXiv, /abs/2005.11401, 2020.
- Lewkowycz, Aitor, et al., Solving quantitative reasoning problems with language models, Advances in Neural Information Processing Systems 35, p. 3843-3857, 2022.
- Li, Zongjie, et al., On Extracting Specialized Code Abilities from Large Language Models: A Feasibility Study, ArXiv, /abs/2303.03012, 2023.
- Lissack, Michael., The Slodderwetenschap (Sloppy Science) of Stochastic Parrots -- A Plea for Science to NOT Take the Route Advocated by Gebru and Bender, ArXiv, /abs/2101.10098, 2021.
- Liu, Yinhan, et al., RoBERTa: A Robustly Optimized BERT Pretraining Approach, ArXiv, /abs/1907.11692, 2019.
- Liu, Bowen, et al., Adversarial Attacks on Large Language Model-Based System and Mitigating Strategies: A Case Study on ChatGPT, Security and Communication Networks 2023, 2023.
- Liu, Yi, et al., Prompt Injection Attack against LLM-integrated Applications, ArXiv, /abs/2306.05499, 2023.
- Marr, Bernard, The 10 Biggest Generative AI Trends For 2024 Everyone Must Be Ready For Now, Forbes. Retrieved on 19.12.2023 at: <https://www.forbes.com/sites/bernardmarr/2023/10/02/the-10-biggest-generative-ai-trends-for-2024-everyone-must-be-ready-for-now/>.
- Marr, Bernard, The 10 Most Important AI Trends For 2024 Everyone Must Be Ready For Now, Forbes. Retrieved on 19.12.2023 at: <https://www.forbes.com/sites/bernardmarr/2023/09/18/the-10-most-important-ai-trends-for-2024-everyone-must-be-ready-for-now/>.
- Merchant, A., Batzner, S., Schoenholz, S.S. et al. Scaling deep learning for materials discovery. Nature 624, 80–85 (2023). <https://doi.org/10.1038/s41586-023-06735-9>
- Meta (2023): Introducing Llama 2. Retrieved on 19.12.2023 at: <https://ai.meta.com/llama/>.
- Microsoft AI (2023): Artificial Intelligence Solutions. Retrieved on 19.12.2023 at: <https://www.microsoft.com/en-us/ai>.
- Mikolov, T., et al., Efficient estimation of word representations in vector space, ArXiv, /abs/1301.3781, 2013.
- Miller, Alexander H., et al., ParlAI: A Dialog Research Software Platform, ArXiv, /abs/1705.06476, 2017.
- Min, Bonan, et al., Recent Advances in Natural Language Processing via Large Pre-Trained Language Models: A Survey, ArXiv, /abs/2111.01243, 2021.
- Morris, Meredith R., et al., Levels of AGI: Operationalizing Progress on the Path to AGI, ArXiv, /abs/2311.02462, 2023.
- MosaicML (2023): Introducing MPT-7B: A New Standard for Open-Source, Commercially Usable LLMs. Retrieved on 19.12.2023 at: <https://www.mosaicml.com/blog/mpt-7b>.

- Murali, Vijayaraghavan, et al., CodeCompose: A Large-Scale Industrial Deployment of AI-assisted Code Authoring, ArXiv, /abs/2305.12050, 2023.
- Naveed, Humza, et al., A Comprehensive Overview of Large Language Models, ArXiv, /abs/2307.06435, 2023.
- News, European Parliament (2023): EU AI Act: first regulation on artificial intelligence. Retrieved on 19.12.2023 at: <https://www.europarl.europa.eu/news/en/headlines/society/20230601STO93804/eu-ai-act-first-regulation-on-artificial-intelligence>.
- Nie, Yixin, et al., Adversarial NLI: A New Benchmark for Natural Language Understanding, ArXiv, /abs/1910.14599, 2019.
- O\*NET OnLine (2023): National Center for O\*NET Development. Retrieved on 19.12.2023 at: [www.onetonline.org/](http://www.onetonline.org/).
- OWASP Foundation (2023): OWASP Top 10 for Large Language Model Applications. Retrieved on 19.12.2023 at: <https://owasp.org/www-project-top-10-for-large-language-model-applications/>.
- OpenAI (2023): Creating safe AGI that benefits all of humanity. Retrieved on 19.12.2023 at: <https://openai.com/>.
- OpenAI, GPT-4 Technical Report, ArXiv, /abs/2303.08774, 2023.
- OpenAI (2023): The moderations is a tool that can use to check whether content complies with OpenAI's usage policies. Retrieved on 19.12.2023 at: <https://platform.openai.com/docs/api-reference/moderations>.
- OpenAI (2023): Evals provide a framework for evaluating LLMs. Retrieved on 19.12.2023 at: <https://github.com/openai/evals>.
- Paperno et al., The LAMBADA dataset: Word prediction requiring a broad discourse context, ACL, 2016.
- Perri, Lori, What's New in Artificial Intelligence from the 2023 Gartner Hype Cycle, Gartner. Retrieved on 19.12.2023 at: <https://www.gartner.com/en/articles/what-s-new-in-artificial-intelligence-from-the-2023-gartner-hype-cycle>.
- Qian, Chen, et al., Communicative Agents for Software Development., ArXiv, /abs/2307.07924, 2023.
- Radford, Alec, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, Improving language understanding by generative pre-training, 2018.
- Rae, Jack W., et al., Scaling language models: Methods, analysis & insights from training gopher, ArXiv, /abs/2112.11446, 2021.
- Raffel, Colin, et al., Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer, ArXiv, /abs/1910.10683, 2019.
- Rajpurkar, Pranav, et al., Know What You Don't Know: Unanswerable Questions for SQuAD, ArXiv, /abs/1806.03822, 2018.
- Reddy, Siva, et al., CoQA: A Conversational Question Answering Challenge, ArXiv, /abs/1808.07042, 2018.
- Research, Microsoft, and Microsoft A. Quantum., The Impact of Large Language Models on Scientific Discovery: A Preliminary Study Using GPT-4, ArXiv, /abs/2311.07361, 2023.
- Ronen, Ofek, Fighting Fire with Fire: Combatting LLM-Generated Social Engineering Attacks With LLMs, Perception Point Blog. Retrieved on 19.12.2023 at: <https://perception-point.io/blog/fighting-fire-with-fire-combatting-llm-generated-social-engineering-attacks-with-llms/>, 2023.
- Russell, S. (2019). Human Compatible: Artificial Intelligence and the Problem of Control. Viking.
- Schick, Timo, Hinrich Schütze., It's Not Just Size That Matters: Small Language Models Are Also Few-Shot Learners, ArXiv, /abs/2009.07118, 2020.
- Schneier, Bruce, LLMs and Phishing, Personal Blog. Retrieved on 19.12.2023 at: <https://www.schneier.com/blog/archives/2023/04/llms-and-phishing.html>, 2023.
- Shanahan, Murray. The technological singularity. MIT press, 2015.
- Shayegani, Erfan, et al., Survey of Vulnerabilities in Large Language Models Revealed by Adversarial Attacks, ArXiv, /abs/2310.10844, 2023.
- Sheng, Ying, et al, FlexGen: High-Throughput Generative Inference of Large Language Models with a Single GPU, ArXiv, /abs/2303.06865, 2023.
- Shevlin, Henry, et al., The limits of machine intelligence: Despite progress in machine intelligence, artificial general intelligence is still a major challenge, EMBO reports 20.10, e49177, 2019.
- Shinn, N., Cassano, F., Labash, B., Gopinath, A., Narasimhan, K., Yao, S., Reflexion: Language Agents with Verbal Reinforcement Learning, ArXiv, /abs/2303.11366, 2023.
- Shoeybi, Mohammad, et al, Megatron-Im: Training multi-billion parameter language models using model parallelism, ArXiv, /abs/1909.08053, 2019.

- Srivastava, Aarohi, et al., Beyond the Imitation Game: Quantifying and Extrapolating the Capabilities of Language Models, ArXiv, /abs/2206.04615, 2022.
- Stiennon, Nisan, et al., Learning to Summarize from Human Feedback, ArXiv, /abs/2009.01325, 2020.
- Sun, Simeng, et al., Exploring the Impact of Low-rank Adaptation on the Performance, Efficiency, and Regularization of RLHF, ArXiv, /abs/2309.09055, 2023.
- Tan, Chee W., et al., Copilot for Xcode: Exploring AI-Assisted Programming by Prompting Cloud-based Large Language Models, ArXiv, /abs/2307.14349, 2023.
- Taori, Rohan, et al., Alpaca: A strong, replicable instruction-following model, Stanford Center for Research on Foundation Models. <https://crfm.stanford.edu/2023/03/13/alpaca.html>, 2023.
- Tegmark, M. (2017). Life 3.0: Being Human in the Age of Artificial Intelligence. Knopf.
- Tenney, Ian, et al., The Language Interpretability Tool: Extensible, Interactive Visualizations and Analysis for NLP Models, ArXiv, /abs/2008.05122, 2020.
- Thoppilan, Romal, et al., LaMDA: Language Models for Dialog Applications, ArXiv, /abs/2201.08239, 2022.
- Touvron, Hugo, et al., Llama: Open and efficient foundation language models, ArXiv, 2302.13971, 2023.
- Touvron, Hugo, et al., Llama 2: Open Foundation and Fine-Tuned Chat Models, ArXiv, /abs/2307.09288, 2023.
- Trim, Neil, LLMs boosting success of phishing campaigns warns Egress, technologyreseller.uk, News. Retrieved on 19.12.2023 at: <https://technologyreseller.uk/llms-boosting-success-of-phishing-campaigns-warns-egress/>, 2023.
- U.S. Bureau of Labor Statistics (2022): Growth trends for selected occupations considered at risk from automation. Retrieved on 19.12.2023 at: <https://www.bls.gov/opub/mlr/2022/article/growth-trends-for-selected-occupations-considered-at-risk-from-automation.htm>.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., Polosukhin, I., Attention Is All You Need, ArXiv, /abs/1706.03762, 2017.
- Vera Sorin, Eyal Klang, Large language models and the emergence phenomena, European Journal of Radiology Open, Volume 10, <https://doi.org/10.1016/j.ejro.2023.100494>, 2023.
- World Economic Forum (2023): Jobs of Tomorrow: Large Language Models and Jobs. Retrieved on 19.12.2023 at: [https://www.weforum.org/docs/WEF\\_Jobs\\_of\\_Tomorrow\\_Generative\\_AI\\_2023.pdf](https://www.weforum.org/docs/WEF_Jobs_of_Tomorrow_Generative_AI_2023.pdf).
- Wallace, Eric, et al., Concealed Data Poisoning Attacks on NLP Models, ArXiv, /abs/2010.12563, 2020.
- Wan, Alexander, et al., Poisoning Language Models During Instruction Tuning, ArXiv, /abs/2305.00944, 2023.
- Wang, Alex, et al., GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding, ArXiv, /abs/1804.07461, 2018.
- Wang, Alex, et al., SuperGlue: A stickier benchmark for general-purpose language understanding systems, Advances in neural information processing systems 32, 2019.
- Wang, Thomas, et al., What Language Model Architecture and Pretraining Objective Work Best for Zero-Shot Generalization?, ArXiv, /abs/2204.05832, 2022.
- Wang, JiongXiao, et al., On the Exploitability of Reinforcement Learning with Human Feedback for Large Language Models, ArXiv, /abs/2311.09641, 2023.
- Wang, JiongXiao, et al., Adversarial Demonstration Attacks on Large Language Models, ArXiv, /abs/2305.14950, 2023.
- Webb, T., Holyoak, K. J., & Lu, H., Emergent analogical reasoning in large language models, Nature Human Behaviour, 7(9), p. 1526-1541, 2023.
- Wei, Jason, et al., Emergent Abilities of Large Language Models, ArXiv, /abs/2206.07682, 2022.
- Weights & Biases (2023): Developer tools for ML. Retrieved on 19.12.2023 at: <https://wandb.ai/site>.
- Weng, Lilian, How to Train Really Large Models on Many GPUs?, blog, 2021. Retrieved on 19.12.2023 at: <https://lilianweng.github.io/posts/2021-09-25-train-large/>.
- White, Jules, et al., A Prompt Pattern Catalog to Enhance Prompt Engineering with ChatGPT, ArXiv, /abs/2302.11382, 2023.
- WhyLabs (2023): Third party services for monitoring and detecting toxicity and other abuses against LLMs. Retrieved on 19.12.2023 at: <https://whylabs.ai/whylogs>.
- WhyLabs (2023): LangKit - An open-source toolkit for monitoring Large Language Models. Retrieved on 19.12.2023 at: <https://github.com/whylabs/langkit>.
- Williams, Adina, et al., A Broad-Coverage Challenge Corpus for Sentence Understanding through Inference, ArXiv, /abs/1704.05426, 2017.

Wu, Shijie, et al., BloombergGPT: A Large Language Model for Finance, ArXiv, /abs/2303.17564, 2023.

xAI (2023): Announcing Grok. Retrieved on 19.12.2023 at: <https://x.ai/>.

Yang, Zhilin, et al., XLNet: Generalized Autoregressive Pretraining for Language Understanding, ArXiv, /abs/1906.08237, 2019.

Yang, Haomiao, et al., A Comprehensive Overview of Backdoor Attacks in Large Language Models within Communication Networks, ArXiv, /abs/2308.14367, 2023.

Zellers, Rowan, et al., HellaSwag: Can a Machine Really Finish Your Sentence?, ArXiv, /abs/1905.07830, 2019.

Zhang, H., Dong, Y., Xiao, C., Oyamada, M., Large Language Models as Data Preprocessors, ArXiv, /abs/2308.16361, 2023.

Zhou, Shuai, et al., Boosting Model Inversion Attacks with Adversarial Examples, ArXiv, /abs/2306.13965, 2023.

Zhu, Xunyu, et al., A Survey on Model Compression for Large Language Models, ArXiv, /abs/2308.07633, 2023.

Zhu, Kaijie, et al., PromptBench: Towards Evaluating the Robustness of Large Language Models on Adversarial Prompts, ArXiv, /abs/2306.04528, 2023.

Zou, Andy, et al., Universal and Transferable Adversarial Attacks on Aligned Language Models, ArXiv, /abs/2307.15043, 2023.



## GETTING IN TOUCH WITH THE EU

### ***In person***

All over the European Union there are hundreds of Europe Direct centres. You can find the address of the centre nearest you online ([european-union.europa.eu/contact-eu/meet-us\\_en](https://european-union.europa.eu/contact-eu/meet-us_en)).

### ***On the phone or in writing***

Europe Direct is a service that answers your questions about the European Union. You can contact this service:

- by freephone: 00 800 6 7 8 9 10 11 (certain operators may charge for these calls),
- at the following standard number: +32 22999696,
- via the following form: [european-union.europa.eu/contact-eu/write-us\\_en](https://european-union.europa.eu/contact-eu/write-us_en).

## FINDING INFORMATION ABOUT THE EU

### ***Online***

Information about the European Union in all the official languages of the EU is available on the Europa website ([european-union.europa.eu](https://european-union.europa.eu)).

### ***EU publications***

You can view or order EU publications at [op.europa.eu/en/publications](https://op.europa.eu/en/publications). Multiple copies of free publications can be obtained by contacting Europe Direct or your local documentation centre ([european-union.europa.eu/contact-eu/meet-us\\_en](https://european-union.europa.eu/contact-eu/meet-us_en)).

### ***EU law and related documents***

For access to legal information from the EU, including all EU law since 1951 in all the official language versions, go to EUR-Lex ([eur-lex.europa.eu](https://eur-lex.europa.eu)).

### ***EU open data***

The portal [data.europa.eu](https://data.europa.eu) provides access to open datasets from the EU institutions, bodies and agencies. These can be downloaded and reused for free, for both commercial and non-commercial purposes. The portal also provides access to a wealth of datasets from European countries.

# An introduction to Large Language Models and their relevance for statistical offices

This manual is a straightforward resource for data professionals of Statistical Offices, introducing the use of Large Language Models (LLMs) in the field of Official Statistics. It outlines how LLMs can tackle complex data problems with their advanced language processing capabilities and integrates these models into current processes. This guide introduces LLMs, delineating their evolution, architecture, applications, and implications for future employment within the AI realm. Additionally, it emphasizes the need for ethical and responsible applications, blending research insights with practical industry examples to ensure professionals can maximize LLM benefits while maintaining trust and reliability in their work.

---

**For more information**

**<https://ec.europa.eu/eurostat/>**



Publications Office  
of the European Union

ISBN 978-92-68-11346-2