

Review

Not peer-reviewed version

Generative AI and Large Language Models: A Comprehensive Scientific Review

[Bogdan-Iulian Ciubotaru](#) *

Posted Date: 7 April 2025

doi: 10.20944/preprints202504.0413.v1

Keywords: generative AI; large language models; LLM



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a Creative Commons CC BY 4.0 license, which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Review

Generative AI and Large Language Models: A Comprehensive Scientific Review

Bogdan-Iulian Ciubotaru

Military Equipment and Technologies Research Agency (METRA), Ministry of National Defence,
077025 Clinceni, Romania; bciubotaru@acttm.ro or ciubotarubogdaniulian@gmail.com

Abstract: This systematic review examines the evolution, technical architecture, applications, limitations, and future directions of generative artificial intelligence (AI) and large language models (LLMs). Through comprehensive analysis of scientific literature, it was traced the development of these technologies from early linguistic theories to modern transformer-based architectures. The findings presented in this review article reveal the transformative impact of LLMs across diverse domains including healthcare, education, software development, and creative industries. Significant technical limitations were identified, including hallucinations, context window constraints, and reasoning deficiencies, alongside ethical concerns regarding bias, privacy, and environmental impact. The review concludes by exploring emerging trends in model architecture, efficiency improvements, and ethical frameworks that will shape future development. This work provides researchers, practitioners, and policymakers with a comprehensive understanding of the current state and future trajectory of generative AI and LLMs.

Keywords: generative AI; large language models; LLM

1. Introduction

Generative Artificial Intelligence (GenAI): is the hype real or it's another technology in its early stages, with a lot of promises? In recent years, the field of artificial intelligence has witnessed a remarkable transformation with the emergence of generative AI and Large Language Models (LLMs). These technologies have revolutionized how machines understand, process, and generate human language, marking a significant milestone in the evolution of AI capabilities [1]. Generative AI, particularly in the form of LLMs, has captured widespread attention not only within academic and research communities but also across industries, governments, and the general public [2]. The launch of GPT models by OpenAI in November 2022 can be considered as a turning point for artificial intelligence, which later can be referred as 0-day for GenAI.

Generative AI refers to artificial intelligence systems capable of creating new content, including text, images, audio, code, and other media formats, based on patterns learned from existing data [3]. Its ability depends on data which was used for training and the prompts provided by a user. The focal point of this technological revolution are Large Language Models—sophisticated neural network architectures trained on datasets of text data that can generate coherent, contextually relevant, and increasingly human-like responses to prompts [4]. These models have demonstrated unprecedented capabilities in understanding context, generating creative content, answering complex questions, and even proving some incipient reasoning abilities that were previously thought to be exclusively human domains [5].

LLMs are more than simple text generators. These models are transforming numerous fields, from healthcare and education to software development and creative industries [6]. In healthcare, LLMs are being utilized for clinical documentation, medical research synthesis, and patient communication [7]. Software developers leverage their coding capabilities using LLMs these models for code generation and debugging assistance [8], while creative professionals like artists and social media influencers use them for content creation, and design [9]. Even in research LLMs can be used

to analyze different data and provide insights which were not previously discovered. The versatility and adaptability of LLMs have positioned them as one of the most significant technological advancements of the 21st century [10].

Modern LLMs are basically the result of decades of work in areas like natural language processing, machine learning, and neural networks — all that research finally coming together [11]. The big game-changer came in 2017 when the transformer architecture was introduced — it completely changed the way AI handles texts [12]. This new setup, thanks to something called self-attention, allowed models to understand the bigger picture in language — like how words relate to each other even if they're far apart in a sentence — way better than anything before. [13]. Subsequent innovations in training methodologies, computational resources, and data availability have led to the rapid evolution of increasingly powerful models, from GPT (Generative Pre-trained Transformer) to BERT (Bidirectional Encoder Representations from Transformers), LaMDA, PaLM, and beyond [14]. While until recent days, LLMs required huge Graphical Processing Unit (GPU) resources, recent technological advances from the Chinese company DeepSeek now enable end users to deploy powerful, efficient language models using significantly fewer computational resources and even CPU usage. These advancements improve accessibility, allowing smaller organizations, researchers, and even individuals to leverage sophisticated LLM capabilities on more affordable hardware setups, significantly reducing operational costs and broadening the practical applicability of generative AI [15].

Even though LLMs are powerful, there are some important problems and limitations which development must overcome. These include tendencies to generate plausible-sounding but factually incorrect information (hallucinations) [16], perpetuate biases which can be detected in training data [17], and consume substantial high amounts of computational resources during training [18]. People have also started worrying about things like privacy, ownership of content, and how these models might be used in the wrong way — all of which have become really important issues.[19]. The main idea would be that if AI is doing something bad, who should be accounted for it. As these technologies continue to evolve and integrate into unimagined dimensions of modern society, addressing these challenges becomes increasingly important for responsible development and deployment [16].

This scientific review aims to provide a comprehensive examination of generative AI and Large Language Models, exploring their historical evolution, technical architecture, capabilities, applications, limitations, and future directions. By synthesizing insights from academic research, industry developments, and practical implementations, this review offers a holistic understanding of these transformative technologies. This article pulls together insights from trusted sources to give a clear and balanced look at where generative AI and LLMs stand today — and where they might be headed. While the article focuses on current status of LLMs, it is more than clear that in the future, the various applications of LLMs will interact with every aspect our lives, improving and providing new skills which will foster productivity in both personal and professional lives.

2. Methods

2.1. Research Design

This study employed a systematic review methodology to examine the current state of generative AI and large language models. The systematic approach was chosen to ensure thoroughness, minimize bias, and provide a structured framework for analyzing the diverse and rapidly evolving literature in this field [20].

2.2. Search Strategy

A comprehensive search of scientific literature was conducted by using multiple databases including arXiv, IEEE Xplore, ACM Digital Library, Google Scholar, Web of Science, and MDPI. The search was performed between January and March 2025, focusing on literature published between 2017 (marking the introduction of the transformer architecture) and March 2025.

The following search terms were used in various combinations, highlighted in Figure 1, Word cloud:

- "generative AI" OR "generative artificial intelligence"
- "large language model*" OR "LLM*"
- "transformer model*" OR "attention mechanism*"
- "GPT" OR "BERT" OR "LaMDA" OR "PaLM"
- "neural language model*"
- "self-attention"
- "LLM model*"
- "LLM training"
- "chain-of-thought prompting*"
- "prompt engineering"
- "tokenization"
- "text-to-text transfer"

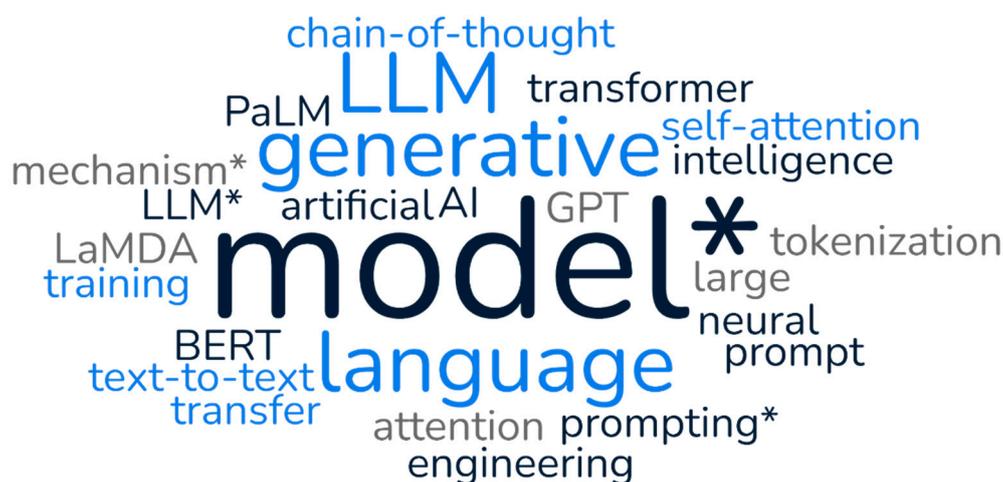


Figure 1. Word cloud for search strategy.

2.3. Inclusion and Exclusion Criteria

Inclusion Criteria:

- Peer-reviewed journal articles, conference proceedings, and technical reports
- Literature focusing on generative AI and large language models
- Studies examining technical architecture, applications, limitations, or future directions
- Publications in English
- Literature published between 2017 and March 2025
- Websites for institutions (European Commission) and blogs for main LLM providers (OpenAI, DeepSeek, Google)

Exclusion Criteria:

- Non-English publications
- Opinion pieces without substantial technical or empirical content
- Studies focusing exclusively on other AI technologies without significant discussion of generative AI or LLMs
- Duplicate publications or multiple reports of the same study

2.4. Data Extraction and Synthesis

Data extraction was performed using a standardized form that captured the following information:

- Publication details (authors, year, journal/conference)
- Study objectives and methodology
- Key findings and contributions
- Technical details of models or architectures discussed
- Applications and use cases
- Limitations and challenges identified
- Future research directions proposed

The extracted data was synthesized using a thematic analysis approach, identifying key themes and patterns across the literature. These themes were organized into categories corresponding to the major sections of this review: historical evolution, technical architecture, applications, limitations and challenges, and future directions.

2.5. Quality Assessment and Limitations of Review Methodology

The quality of included studies was assessed using criteria adapted from the Critical Appraisal Skills Programme (CASP) and the Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) guidelines [21]. For technical papers, methodological rigor, clarity of reporting, and significance of contribution were evaluated. For empirical studies, study design, sample size, data collection methods, and analysis techniques were assessed.

Several limitations of selected methodology should be acknowledged. First, the rapid pace of development in generative AI means that some recent advancements may not be reflected in peer-reviewed literature. Second, proprietary details of commercial LLMs are often not fully disclosed in scientific publications, potentially limiting the understanding of state-of-the-art systems. Third, the interdisciplinary nature of the field necessitated searches across multiple databases, which may have resulted in some relevant studies being overlooked.

3. Results

3.1. Historical Evolution of Large Language Models

The development trajectory of Large Language Models represents a remarkable progression from theoretical linguistics to sophisticated AI systems. The performed analysis identified several distinct phases in this evolution. The early stages might be speculative, taking into account that long period of time.

3.1.1. Early Foundations (1880s-1950s)

The conceptual foundations for language models emerged from semantic theory, with Michel Bréal introducing the concept of semantics in 1883 [22]. Ferdinand de Saussure's work between 1906 and 1912 established a functional model of languages as systems, laying groundwork for the structuralist approach to linguistics [23]. Post-World War II efforts to develop automatic translation highlighted the complexity of human language and the need for computational approaches [24].

3.1.2. Machine Learning and Neural Networks (1950s-1980s)

Alan Turing's 1950 paper "Computing Machinery and Intelligence" established a philosophical framework for machine intelligence that would guide subsequent research [25]. Frank Rosenblatt's creation of the Mark 1 Perceptron in 1958 represented the first artificial neural network, establishing concepts fundamental to modern LLMs [26]. Joseph Weizenbaum's ELIZA program in 1966 demonstrated early natural language processing through pattern matching and substitution methodology [27].

3.1.3. Statistical Approaches and Small Language Models (1980s-1990s)

The 1980s saw IBM develop the first small language models using statistical analysis to predict the next word in a sentence [28]. By the late 1980s, increased computational power and improved machine learning algorithms led to a revolution in natural language processing, with a shift from handwritten rules to statistical models [29].

3.1.4. Deep Learning and Neural Networks (1990s-2010s)

Yoshua Bengio and colleagues published foundational work on neural language models in 2003 [11]. Word embeddings through models like Word2Vec, developed by Tomas Mikolov and colleagues at Google in 2011, represented words as vectors in high-dimensional space, enabling more sophisticated language processing [30]. The development of graphics processing units (GPUs) in the early 2010s provided the computational power necessary for training increasingly complex neural networks [31].

3.1.5. The Transformer Revolution (2017-Present)

The publication of "Attention Is All You Need" by Vaswani et al. in 2017 introduced the transformer architecture, utilizing self-attention mechanisms to process sequential data more effectively than previous approaches [12]. OpenAI's release of GPT in 2018, followed by GPT-2 in 2019 and GPT-3 in 2020, demonstrated dramatic improvements in scale and capabilities [32]. Google introduced BERT in 2018, which excelled at understanding context by considering words bidirectionally [1].

3.1.6. Current Landscape (2021-Present)

OpenAI's ChatGPT, released in November 2022, brought generative AI to mainstream attention. The subsequent release of GPT-4 in 2023 introduced multimodal capabilities that could process both text and images [33]. Open-source alternatives like Meta's LLaMA, EleutherAI's GPT-J and GPT-NeoX, and Stability AI's StableLM have made powerful language models more accessible [14]. Multimodal models like GPT-4V, Google's Gemini, and Anthropic's Claude represent the latest frontier in LLM evolution.

3.2. Technical Architecture of Generative LLMs

The analysis of the technical approaches of modern LLMs revealed several key architectural components and principles used for these technologies to achieve their capabilities.

3.2.1. Transformer Architecture

The transformer architecture consists of an encoder and a decoder, though many modern generative LLMs utilize only the decoder portion in a "decoder-only" architecture [4]. The defining innovation is the self-attention mechanism, which allows the model to consider all tokens in a sequence simultaneously, weighing their relevance to each other [12].

Self-attention computes a weighted sum of all token representations in a sequence, where the weights (attention scores) are determined by the relevance of each token to the current token being processed [34]. Multi-head attention splits the attention computation into multiple "heads", each focusing on different aspects of the input sequence [35]. Positional encodings address the limitation that transformers process all tokens in parallel by providing numerical representations that encode the position of each token in the sequence [36].

3.2.2. Components of Modern LLMs

Modern LLMs convert tokens (words or subwords) into numerical representations called embeddings using subword tokenization methods such as Byte-Pair Encoding (BPE) or

SentencePiece [37]. Each transformer block contains a feed-forward neural network that processes the output of the attention mechanism, introducing non-linearity and increasing representational capacity [38]. Layer normalization stabilizes training and improves convergence by normalizing the activations across features for each example [39]. Residual connections (skip connections) around each sub-layer facilitate gradient flow during training, especially in deep models [40].

3.2.3. Architectural Variations

The GPT family of models uses a decoder-only transformer architecture with masked self-attention, making the model autoregressive [32]. BERT uses only the encoder portion of the transformer, allowing bidirectional attention [1]. T5 (Text-to-Text Transfer Transformer) uses the complete encoder-decoder architecture, framing all NLP tasks as text-to-text problems [41]. Google's PaLM and Gemini models introduce architectural innovations for improved scaling and multimodal capabilities [42]. The most known LLMs compared by using their key features, parameters, training approach, the most notable capabilities and their limitations are presented in Table 1. It should be noted that the new generation LLMs like deepseek are not presented, taking into account the limited knowledge about how they work.

Table 1. Widely known LLMs comparison.

Model Family	Architecture Type	Parameters	Key Features	Training Approach	Notable Capabilities	Limitations
GPT Series	Decoder-only	GPT-3: 175B	Autoregressive generation	Generative pre-training followed by fine-tuning	Text generation	Limited bidirectional context
		GPT-4: ~1.7T (estimated)	Masked self-attention		Few-shot learning In-context learning	Tendency to hallucinate
BERT	Encoder-only	BERT-Large: 340M	Bidirectional attention	Masked token prediction	Strong text understanding	Limited generation capabilities
		RoBERTa: 355M	Masked language modeling	Next sentence prediction	Classification tasks Question answering	Fixed context window
T5	Encoder-decoder	T5-11B: 11B	Text-to-text framework Complete transformer architecture	Unified text-to-text approach for all NLP tasks	Versatile across tasks Strong transfer learning	Larger computational requirements Complex training process
PaLM	Decoder-only	PaLM: 540B PaLM 2: ~340B	Pathways architecture Efficient training	Trained on diverse multilingual data	Multilingual capabilities Strong reasoning Code generation	Proprietary architecture Limited public information
LLaMA	Decoder-only	LLaMA 2: 7B-70B	Optimized for efficiency Open weights	Trained on diverse text corpora	Open-source accessibility Efficient fine-tuning	Smaller parameter count than largest models Limited multimodal capabilities
Claude	Decoder-only	Claude 2: ~140B (estimated)	Constitutional AI approach RLHF training	Trained with constitutional principles	Safety and alignment focus Long context window	Less information on architecture Proprietary system
Gemini	Multimodal transformer	Gemini Ultra: ~1T (estimated)	Multimodal from ground up Advanced reasoning	Trained on multimodal data	Strong multimodal reasoning Advanced problem-solving	Proprietary architecture High computational requirements
Mistral	Decoder-only	Mistral 7B: 7B Mixtral 8x7B: 47B	Grouped-query attention Mixture of experts	Efficient training techniques	Strong performance despite smaller size Efficient inference	Newer architecture with less testing Smaller parameter count

3.2.4. Scaling Properties and Emergent Capabilities

Research has identified several important scaling laws: performance on language tasks follows a power-law relationship with model size [13]; performance also scales with the amount of training data, though with diminishing returns [43]; and training compute (the product of model size and training tokens) is another key factor [44]. LLMs exhibit emergent abilities—capabilities not present in smaller models but appearing once models reach a certain scale—including in-context learning, chain-of-thought reasoning, and instruction following [45].

3.2.5. Efficiency Innovations

To address the computational demands of large models, researchers have developed various efficiency techniques: sparse attention mechanisms limit each token to attending only to a subset of other tokens [35]; parameter-efficient fine-tuning methods like LoRA allow for efficient adaptation with minimal parameter updates [46]; quantization reduces the precision of model weights to decrease memory requirements [47]; and knowledge distillation transfers knowledge from a large "teacher" model to a smaller "student" model [48].

3.3. Applications and Use Cases

Diverse applications of generative AI and LLMs across multiple domains were discovered, as illustrated in Figure 2, demonstrating their versatility and transformative potential.

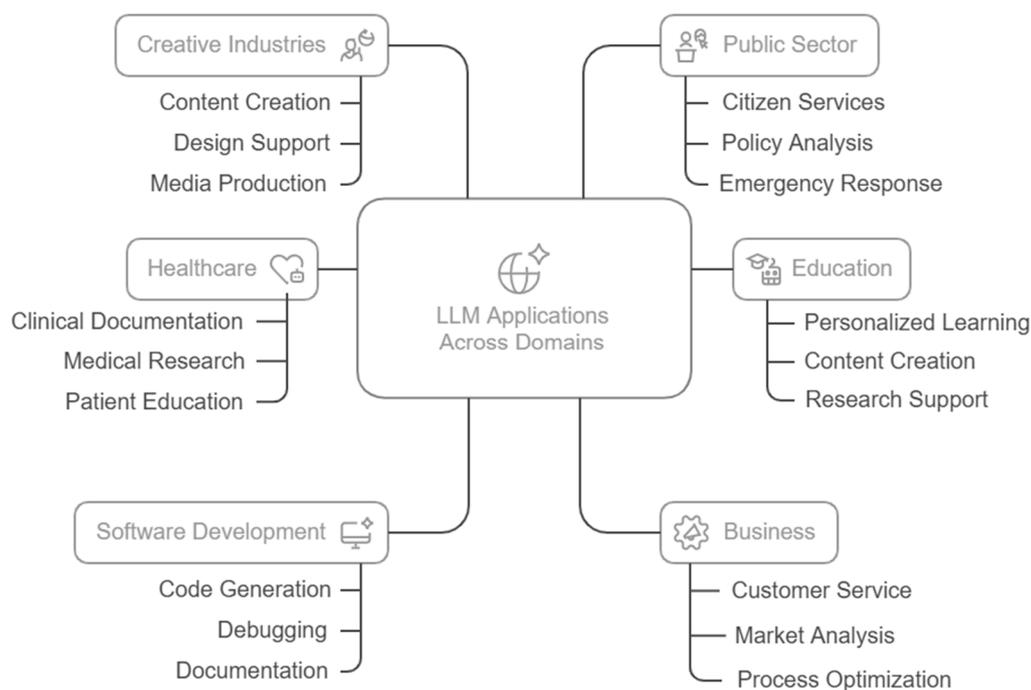


Figure 2. LLM applications in different domains.

3.3.1. Natural Language Understanding and Text Generation

LLMs excel at text summarization and content creation, distilling lengthy documents into concise summaries while preserving key information [49]. They have significantly advanced machine translation capabilities, supporting communication across language barriers with unprecedented fluency [50]. The conversational capabilities of LLMs have revolutionized virtual assistants and chatbots, enabling coherent, contextually appropriate dialogues [51].

3.3.2. Educational Applications

LLMs offer unprecedented opportunities for personalized education through adaptive tutoring systems that explain concepts in multiple ways and provide tailored feedback [52]. Educators leverage LLMs to develop curriculum materials, lesson plans, and educational resources [53]. In academic contexts, LLMs assist with literature reviews by summarizing research papers and identifying connections between studies [54].

3.3.3. Healthcare Applications

LLMs assist healthcare professionals with clinical documentation, generating notes, discharge summaries, and referral letters [7]. They synthesize findings across thousands of studies, identify emerging trends, and summarize evidence-based practices [55]. Artificial intelligence models can be used in combination with IoT devices to detect different pathologies or to detect early signs of future disorders [56,57]. LLMs facilitate improved patient education by generating personalized health information that accounts for a patient's specific condition and health literacy level [35].

3.3.4. Software Development and Programming

Programming-focused LLMs generate code snippets, complete partial code, and translate requirements into functional implementations [8]. They assist with debugging by identifying potential issues in code, suggesting fixes, and explaining the underlying problems [58]. LLMs can generate comprehensive documentation from code, including function descriptions, parameter explanations, and usage examples [59].

3.3.5. Business and Enterprise Applications

Enterprises deploy LLM-powered systems to handle customer inquiries, troubleshoot common issues, and provide product information [6]. LLMs analyze consumer feedback, social media conversations, and market trends to provide businesses with actionable insights [60]. In industrial settings, LLMs integrate with sensor data and operational metrics to optimize processes and predict maintenance needs [6].

3.3.6. Creative and Media Applications

LLMs support creative professionals in developing narratives, scripts, poetry, and other creative works [2]. Multimodal LLMs that combine text and image understanding support various design applications [61]. In music and audio production, specialized LLMs assist with composition, arrangement, and sound design [62].

3.3.7. Public Sector and Governance

Government agencies use LLMs to improve citizen services through more accessible information delivery and streamlined interactions [6]. LLMs assist policymakers by analyzing large volumes of data, simulating potential policy outcomes, and identifying unintended consequences [63]. During emergencies, LLMs help coordinate response efforts by processing incoming information and facilitating communication between different agencies [64].

3.4. Limitations and Challenges

A range of critical limitations and challenges has been identified in generative AI and LLMs, all of which demand attention for responsible development and deployment. The following sections detail these issues, highlighting technical, ethical, regulatory, environmental, and integration aspects.

3.4.1. Technical Limitations

LLMs tend to generate content that appears plausible but contains factual errors or fabricated information—a phenomenon commonly referred to as "hallucination" [16]. Current LLMs operate within fixed context windows that constrain their ability to maintain coherence and consistency across long documents or conversations [35]. Despite their impressive language capabilities, LLMs struggle with complex reasoning, logical consistency, and mathematical accuracy [65]. State-of-the-art models require enormous computational resources for both training and inference, creating barriers to entry for smaller organizations and researchers [18].

3.4.2. Ethical Concerns

LLMs learn from vast corpora of human-generated text, inevitably absorbing and potentially amplifying biases present in that data [17]. The development and deployment of LLMs raise significant privacy concerns, particularly in domains like medicine where confidentiality is paramount [5]. LLMs trained on vast corpora of text raise complex intellectual property questions regarding copyright status of training data and ownership of AI-generated content [66]. The capabilities of advanced LLMs create potential for various forms of misuse, including generation of misinformation and automated production of harmful content [19].

3.4.3. Regulatory and Compliance Challenges

The rapid advancement of generative AI has outpaced regulatory frameworks, creating uncertainty and compliance challenges across different jurisdictions [67]. LLMs operate as "black boxes" with billions of parameters, making their decision-making processes opaque and difficult to interpret [68]. The complexity and autonomy of LLMs complicate traditional accountability frameworks, with unclear responsibility distribution among developers, deployers, and users [63].

3.4.4. Environmental Impact

The training of large language models requires enormous computational resources, resulting in significant energy consumption and carbon emissions [18]. Beyond training, the ongoing operation of LLMs for inference also consumes significant resources, particularly for high-traffic applications [69].

3.4.5. Implementation and Integration Challenges

Adapting general-purpose LLMs to specific domains presents significant challenges, particularly for specialized fields with domain-specific terminology and knowledge [5]. Incorporating LLMs into existing technological ecosystems presents numerous integration challenges, including interfacing with legacy systems and ensuring consistent performance [6]. Evaluating LLM performance presents unique challenges compared to traditional software due to the subjective nature of many language tasks and the impossibility of comprehensive testing across all possible inputs [70].

3.5. *Future Directions and Emerging Trends*

Several promising research directions and emerging trends have been identified that are expected to influence the future development of generative AI and LLMs. The following sections outline key areas of ongoing innovation, highlighting advancements in model architecture, efficiency, domain specialization, reasoning capabilities, ethical AI, regulatory frameworks, and technological integration.

3.5.1. Advancements in Model Architecture

Researchers are exploring architectural innovations beyond transformers, including sparse attention mechanisms, mixture of experts approaches, retrieval-enhanced architectures, and neuro-

symbolic approaches [34]. The integration of multiple modalities—text, images, audio, video, and structured data—represents a significant frontier in LLM development [61].

3.5.2. Efficiency Improvements

As model scaling faces economic and environmental constraints, research into parameter-efficient architectures is accelerating, including parameter-efficient fine-tuning (PEFT), knowledge distillation, quantization and pruning, and neural architecture search [46]. Innovations in training methodologies aim to reduce computational and data requirements through advances in self-supervised learning, curriculum learning, continual learning, and federated learning [71].

3.5.3. Domain-Specific Specialization

While general-purpose LLMs demonstrate broad capabilities, specialized models tailored to specific domains are likely to proliferate, including scientific, legal, financial, and healthcare-specific models [54]. Expanding beyond English-centric development to better serve global populations through multilingual models, cultural contextualization, and local knowledge integration is an important direction [72].

3.5.4. Enhanced Reasoning Capabilities

Addressing current limitations in structured reasoning represents a critical frontier, with improvements to chain-of-thought techniques, integration with external tools, and formal verification methods [45]. Developing stronger capabilities for understanding causality rather than mere correlation through causal inference, counterfactual reasoning, and temporal reasoning is an active area of research [73].

3.5.5. Ethical AI and Responsible Development

Ensuring that AI systems behave in accordance with human values and intentions through constitutional AI, interpretability research, red-teaming, and value alignment is increasingly important [74]. Addressing issues of bias and fairness through comprehensive frameworks for identifying and mitigating various forms of bias is a critical research direction [17].

3.5.6. Regulatory and Governance Frameworks

The regulatory environment for AI is rapidly developing, with emerging international standards, risk-based regulation, certification and auditing systems, and industry-led self-regulation initiatives [67]. Mechanisms to ensure responsible development and deployment through standardized documentation, explainability tools, and audit trails are being developed [75].

3.5.7. Integration with Other Technologies

The integration of LLMs with broader technological ecosystems through AI-enabled infrastructure, autonomous agents, Internet of Things integration, and smart city applications represents a frontier of development [6]. Evolving paradigms for human-AI interaction through collaborative interfaces, augmented creativity, cognitive prosthetics, and personalized AI assistants are emerging.

4. Discussion

The presented systematic review reveals the remarkable trajectory of generative AI and LLMs from theoretical concepts to transformative technologies with wide-ranging applications. The evolution of these models has been characterized by several key developments: the breakthrough of the transformer architecture with its self-attention mechanism [12], the scaling of models to unprecedented sizes, and the emergence of capabilities not explicitly programmed [45]. These

developments have enabled applications across diverse domains, from healthcare and education to creative industries and public services.

The technical architecture of modern LLMs, centered around self-attention mechanisms and deep neural networks, has proven remarkably effective at capturing the patterns and structures of human language [12]. The scaling properties of these models have revealed fascinating relationships between model size, training data, and performance, suggesting pathways for continued advancement through both scaling and architectural innovation [13].

However, the findings also highlight significant limitations and challenges that must be addressed. Technical constraints such as hallucinations, context window limitations [35], and reasoning deficiencies [65] impact the reliability and applicability of LLMs in critical domains. Ethical concerns regarding bias, privacy, intellectual property, and potential misuse needs careful mitigation strategies. Regulatory uncertainties (while some countries tend to liberate AI development, others want to regulate as strict as possible), environmental impacts, and implementation challenges further complicate the landscape of LLMs development [67].

4.1. Implications for Research and Practice

4.1.1. Research Implications

The findings suggest several important directions for future research. First, addressing the technical limitations of current LLMs, particularly hallucinations and reasoning limits, requires fundamental advances in model architecture and training methodologies [10]. While there are new models arising daily, research must be continued, so new ways of LLM development approach should be proposed. Developing more efficient methods to train and deploy (in this particular case, inference) is essential for offering access to these technologies for new people and reducing their environmental impact [46]. Third, research into interpretability of different concepts and finding new ways how to explain why some models are behaving in a particulate way is critical for addressing the "black box" nature of current systems [68].

The interdisciplinary nature of challenges associated with LLMs necessitates collaboration across fields including computer science, linguistics, cognitive science, ethics, law, and social sciences [3].

4.1.2. Practical Implications

For people working in the field, this review points out how important it is to have smart, well-planned strategies — ones that really take into account what today's LLMs can and can't do. Organizations deploying these technologies should implement robust evaluation frameworks, monitoring systems, and governance structures to ensure responsible use [75]. Domain-specific adaptation through fine-tuning or retrieval augmentation is essential for applications in specialized fields, like military and medicine.

The development of clear guidelines for human-AI collaboration can maximize the benefits of these technologies while maintaining appropriate human oversight. Educational initiatives to build AI literacy among professionals and the general public are important for enabling informed engagement with these technologies [52,72].

4.2. Ethical Considerations

The general adoption of generative AI and LLMs raises critical ethical questions that require consideration. The potential for these technologies to exacerbate existing inequalities through biased outputs or unequal access demands proactive approaches to fairness and accessibility [17]. The environmental impact of large-scale AI systems necessitates more sustainable approaches to development and deployment [18].

Questions of authorship, originality, and intellectual property in the context of AI-generated content challenge traditional legal and cultural frameworks [66]. The potential for automation to

disrupt labor markets requires thoughtful approaches to workforce transition and the development of new roles that leverage human-AI collaboration [76].

4.3. Limitations of the Current Review

Several limitations of this review should be acknowledged. First, the rapid pace of development in this field means that some recent advancements may not be fully reflected in the literature analyzed. Second, proprietary details of commercial LLMs are often not fully disclosed, potentially limiting understanding of state-of-the-art systems. Third, focus on English-language publications may have excluded valuable insights from non-English research communities.

5. Conclusions

This systematic review has comprehensively examined the evolution, technical architecture, applications, limitations, and future directions of generative AI and Large Language Models. The findings demonstrate the transformative potential of these technologies across diverse domains, while also highlighting significant challenges that must be addressed for responsible development and deployment.

The technical architecture of modern LLMs, centered around self-attention mechanisms and deep neural networks, has enabled unprecedented capabilities in language understanding and generation. These technologies already demonstrate considerable promise across a wide variety of fields, including healthcare, education, software development, creative work, and public services. Their adaptability is clear from the way they're being used to optimize medical research, personalize learning experiences, generate innovative software solutions, spark new forms of art and media, and streamline government operations.

At the same time, there are important limitations which can't be overlooked. On the technical side, issues like model hallucinations and limited reasoning capabilities represent real challenges, often undermining the reliability of the systems and the adoption of such systems. On the ethics side, there are still big concerns about bias, data privacy, and making sure everyone has fair access. Figuring out how to regulate these technologies is still an open question, and on top of that, their environmental impact — since they use a lot of computing power — adds even more complexity. Putting these systems into real-world use isn't easy either, with challenges like high costs and the need to train people to work with them. Tackling all this will take teamwork between tech folks, policymakers, and others, along with ongoing research to keep pushing the limits of what these models can actually do.

New breakthroughs in how these models are built could make them a lot more accurate and efficient, while also using less energy. Improving how well they "reason" could help fix some of the current problems—like when they give wrong or confusing answers. It also makes sense to fine-tune models for specific industries or tasks, so the results they give are more accurate and useful. On top of that, combining LLMs with other types of AI—like computer vision or reinforcement learning—might open the door to totally new and unexpected applications.

These technologies aren't just about technical concepts — they're changing the way we live, work, and connect with each other. As LLMs become part of our everyday lives — helping students learn, supporting doctors, powering apps, inspiring creativity, and even helping write this review — their impact on society and culture keeps growing. To make sure that impact is a positive one, we need to build in ethical values, fairness, and accountability from the ground up. And we also need smart governance — laws, policies, and systems that put people first and protect things like privacy, equality, and responsible innovation.

Generative AI and Large Language Models represent a groundbreaking convergence of potential and complexity. If we make the most of what these tools can do—while keeping a close eye on their flaws and the bigger picture—we've got a real opportunity to shape them into something that truly helps people and benefits society as a whole. What is worth mentioning is that if we compare the evolution of smartphones—taking the launch of the first iPhone in 2007 as a reference

point—with the launch of ChatGPT in 2022, and we extrapolate based on how far smartphone technology has come, it becomes clear that the next AI revolution is a matter of years, not decades.

Acknowledgments: During the preparation of this study, the author used Manus for the purposes of summarization and text-generation for some general statements. The author has reviewed and edited the output and takes full responsibility for the content of this publication.

Conflicts of Interest: The author declares no conflicts of interest.

References

1. Devlin, J.; Chang, M.-W.; Lee, K.; Toutanova, K. BERT: Pre-Training of Deep Bidirectional Transformers for Language Understanding 2019.
2. Brown, T.B.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; et al. Language Models Are Few-Shot Learners 2020.
3. Bommasani, R.; Hudson, D.A.; Adeli, E.; Altman, R.; Arora, S.; von Arx, S.; Bernstein, M.S.; Bohg, J.; Bosselut, A.; Brunskill, E.; et al. On the Opportunities and Risks of Foundation Models 2021.
4. A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever Language Models Are Unsupervised Multitask Learners. *OpenAI Blog*.
5. Yu, P.; Xu, H.; Hu, X.; Deng, C. Leveraging Generative AI and Large Language Models: A Comprehensive Roadmap for Healthcare Integration. *Healthcare* **2023**, *11*, 2776, doi:10.3390/healthcare11202776.
6. Salierno, G.; Leonardi, L.; Cabri, G. Generative AI and Large Language Models in Industry 5.0: Shaping Smarter Sustainable Cities. *Encyclopedia* **2025**, *5*, 30, doi:10.3390/encyclopedia5010030.
7. Sallam, M. ChatGPT Utility in Healthcare Education, Research, and Practice: Systematic Review on the Promising Perspectives and Valid Concerns. *Healthcare* **2023**, *11*, 887, doi:10.3390/healthcare11060887.
8. Chen, M.; Tworek, J.; Jun, H.; Yuan, Q.; Pinto, H.P. de O.; Kaplan, J.; Edwards, H.; Burda, Y.; Joseph, N.; Brockman, G.; et al. Evaluating Large Language Models Trained on Code 2021.
9. Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; Ommer, B. High-Resolution Image Synthesis with Latent Diffusion Models. In Proceedings of the 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR); IEEE: New Orleans, LA, USA, June 2022; pp. 10674–10685.
10. Shuster, K.; Poff, S.; Chen, M.; Kiela, D.; Weston, J. Retrieval Augmentation Reduces Hallucination in Conversation. In Proceedings of the Findings of the Association for Computational Linguistics: EMNLP 2021; Association for Computational Linguistics: Punta Cana, Dominican Republic, 2021; pp. 3784–3803.
11. Bengio, Y.; Ducharme, R.; Vincent, P. A Neural Probabilistic Language Model. In Proceedings of the Proceedings of the 14th International Conference on Neural Information Processing Systems; MIT Press: Cambridge, MA, USA, 2000; pp. 893–899.
12. Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A.N.; Kaiser, L.; Polosukhin, I. Attention Is All You Need 2017.
13. Kaplan, J.; McCandlish, S.; Henighan, T.; Brown, T.B.; Chess, B.; Child, R.; Gray, S.; Radford, A.; Wu, J.; Amodei, D. Scaling Laws for Neural Language Models 2020.
14. Touvron, H.; Lavril, T.; Izacard, G.; Martinet, X.; Lachaux, M.-A.; Lacroix, T.; Rozière, B.; Goyal, N.; Hambro, E.; Azhar, F.; et al. LLaMA: Open and Efficient Foundation Language Models 2023.
15. DeepSeek-AI; Guo, D.; Yang, D.; Zhang, H.; Song, J.; Zhang, R.; Xu, R.; Zhu, Q.; Ma, S.; Wang, P.; et al. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning 2025.
16. Bender, E.M.; Gebru, T.; McMillan-Major, A.; Shmitchell, S. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? . In Proceedings of the Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency; ACM: Virtual Event Canada, March 3 2021; pp. 610–623.

17. Bolukbasi, T.; Chang, K.-W.; Zou, J.; Saligrama, V.; Kalai, A. Man Is to Computer Programmer as Woman Is to Homemaker? Debiasing Word Embeddings 2016.
18. Patterson, D.; Gonzalez, J.; Le, Q.; Liang, C.; Munguia, L.-M.; Rothchild, D.; So, D.; Texier, M.; Dean, J. Carbon Emissions and Large Neural Network Training 2021.
19. Tamkin, A.; Brundage, M.; Clark, J.; Ganguli, D. Understanding the Capabilities, Limitations, and Societal Impact of Large Language Models 2021.
20. Moher, D.; Liberati, A.; Tetzlaff, J.; Altman, D.G.; The PRISMA Group Preferred Reporting Items for Systematic Reviews and Meta-Analyses: The PRISMA Statement. *PLoS Med* **2009**, *6*, e1000097, doi:10.1371/journal.pmed.1000097.
21. Critical Appraisal Skills Programme *CASP Systematic Review Checklist*; **2018**.
22. Bréal, M. *Essai de sémantique: (science des significations)*; 3rd ed.; Lambert- Lucas: Paris, 2005; ISBN 978-2-915806-01-4.
23. Saussure, F. de; Bally, C.; Sechehaye, C.-A.; Urbain, J.-D. *Cours de linguistique générale*; Petite bibliothèque Payot; Éditions Payot & Rivages: Paris, 2016; ISBN 978-2-228-91561-8.
24. Hutchins, W.J. *Machine Translation: Past, Present and Future*; Ellis Horwood series in computers and their applications; Ellis Horwood [u.a.]: Chichester, 1986; ISBN 978-0-470-20313-2.
25. Turing, A.M. *Computing Machinery and Intelligence*; Springer, 2009;
26. Rosenblatt, F. The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain. *Psychological review* **1958**, *65*, 386.
27. Weizenbaum, J. ELIZA—a Computer Program for the Study of Natural Language Communication between Man and Machine. *Communications of the ACM* **1966**, *9*, 36–45.
28. JELINEK, F. *STATISTICAL METHODS FOR SPEECH RECOGNITION*; MIT PRESS: S.l., 2022; ISBN 978-0-262-54660-7.
29. Brown, P.F.; Pietra, V.J.D.; Pietra, S.A.D.; Mercer, R.L. The Mathematics of Statistical Machine Translation: Parameter Estimation. *Comput. Linguist.* **1993**, *19*, 263–311.
30. Mikolov, T.; Chen, K.; Corrado, G.; Dean, J. Efficient Estimation of Word Representations in Vector Space 2013.
31. Sutskever, I.; Vinyals, O.; Le, Q.V. Sequence to Sequence Learning with Neural Networks 2014.
32. Radford, A.; Narasimhan, K.; Salimans, T.; Sutskever, I.; others Improving Language Understanding by Generative Pre-Training. **2018**.
33. OpenAI; Achiam, J.; Adler, S.; Agarwal, S.; Ahmad, L.; Akkaya, I.; Aleman, F.L.; Almeida, D.; Altenschmidt, J.; Altman, S.; et al. GPT-4 Technical Report 2024.
34. Kitaev, N.; Kaiser, L.; Levskaya, A. Reformer: The Efficient Transformer 2020.
35. Beltagy, I.; Peters, M.E.; Cohan, A. Longformer: The Long-Document Transformer 2020.
36. Shaw, P.; Uszkoreit, J.; Vaswani, A. Self-Attention with Relative Position Representations. In Proceedings of the Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers); Association for Computational Linguistics: New Orleans, Louisiana, 2018; pp. 464–468.
37. Sennrich, R.; Haddow, B.; Birch, A. Neural Machine Translation of Rare Words with Subword Units. In Proceedings of the Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers); Association for Computational Linguistics: Berlin, Germany, 2016; pp. 1715–1725.
38. Hendrycks, D.; Gimpel, K. Gaussian Error Linear Units (GELUs) 2016.
39. Ba, J.L.; Kiros, J.R.; Hinton, G.E. Layer Normalization 2016.
40. He, K.; Zhang, X.; Ren, S.; Sun, J. Deep Residual Learning for Image Recognition. In Proceedings of the Proceedings of the IEEE conference on computer vision and pattern recognition; 2016; pp. 770–778.

41. Raffel, C.; Shazeer, N.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; Liu, P.J. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer 2019.
42. Chowdhery, A.; Narang, S.; Devlin, J.; Bosma, M.; Mishra, G.; Roberts, A.; Barham, P.; Chung, H.W.; Sutton, C.; Gehrmann, S.; et al. PaLM: Scaling Language Modeling with Pathways 2022.
43. Hoffmann, J.; Borgeaud, S.; Mensch, A.; Buchatskaya, E.; Cai, T.; Rutherford, E.; Casas, D. de L.; Hendricks, L.A.; Welbl, J.; Clark, A.; et al. Training Compute-Optimal Large Language Models 2022.
44. Henighan, T.; Kaplan, J.; Katz, M.; Chen, M.; Hesse, C.; Jackson, J.; Jun, H.; Brown, T.B.; Dhariwal, P.; Gray, S.; et al. Scaling Laws for Autoregressive Generative Modeling 2020.
45. Wei, J.; Wang, X.; Schuurmans, D.; Bosma, M.; Ichter, B.; Xia, F.; Chi, E.H.; Le, Q.V.; Zhou, D. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. In Proceedings of the Proceedings of the 36th International Conference on Neural Information Processing Systems; Curran Associates Inc.: Red Hook, NY, USA, 2022.
46. Hu, E.J.; Shen, Y.; Wallis, P.; Allen-Zhu, Z.; Li, Y.; Wang, S.; Wang, L.; Chen, W. LoRA: Low-Rank Adaptation of Large Language Models 2021.
47. Han, S.; Mao, H.; Dally, W.J. Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding 2016.
48. Hinton, G.; Vinyals, O.; Dean, J. Distilling the Knowledge in a Neural Network 2015.
49. Zhang, J.; Zhao, Y.; Saleh, M.; Liu, P. Pegasus: Pre-Training with Extracted Gap-Sentences for Abstractive Summarization. In Proceedings of the International conference on machine learning; PMLR, 2020; pp. 11328–11339.
50. Johnson, M.; Schuster, M.; Le, Q.V.; Krikun, M.; Wu, Y.; Chen, Z.; Thorat, N.; Viégas, F.; Wattenberg, M.; Corrado, G.; et al. Google's Multilingual Neural Machine Translation System: Enabling Zero-Shot Translation. *Transactions of the Association for Computational Linguistics* **2017**, *5*, 339–351, doi:10.1162/tacl_a_00065.
51. Roller, S.; Dinan, E.; Goyal, N.; Ju, D.; Williamson, M.; Liu, Y.; Xu, J.; Ott, M.; Shuster, K.; Smith, E.M.; et al. Recipes for Building an Open-Domain Chatbot 2020.
52. Chu, Z.; Wang, S.; Xie, J.; Zhu, T.; Yan, Y.; Ye, J.; Zhong, A.; Hu, X.; Liang, J.; Yu, P.S.; et al. LLM Agents for Education: Advances and Applications 2025.
53. Ouyang, X.; Wang, S.; Pang, C.; Sun, Y.; Tian, H.; Wu, H.; Wang, H. ERNIE-M: Enhanced Multilingual Representation by Aligning Cross-Lingual Semantics with Monolingual Corpora 2021.
54. Beltagy, I.; Lo, K.; Cohan, A. SciBERT: A Pretrained Language Model for Scientific Text 2019.
55. Esteva, A.; Kale, A.; Paulus, R.; Hashimoto, K.; Yin, W.; Radev, D.; Socher, R. CO-Search: COVID-19 Information Retrieval with Semantic Search, Question Answering, and Abstractive Summarization 2020.
56. Chow, J.C.L.; Wong, V.; Li, K. Generative Pre-Trained Transformer-Empowered Healthcare Conversations: Current Trends, Challenges, and Future Directions in Large Language Model-Enabled Medical Chatbots. *BioMedInformatics* **2024**, *4*, 837–852, doi:10.3390/biomedinformatics4010047.
57. Ciubotaru, B.-I.; Sasu, G.-V.; Goga, N.; Vasilăţeanu, A.; Marin, I.; Păvăloiu, I.-B.; Gligore, C.T.I. Frailty Insights Detection System (FIDS)—A Comprehensive and Intuitive Dashboard Using Artificial Intelligence and Web Technologies. *Applied Sciences* **2024**, *14*, 7180, doi:10.3390/app14167180.
58. Svyatkovskiy, A.; Deng, S.K.; Fu, S.; Sundaresan, N. IntelliCode Compose: Code Generation Using Transformer. In Proceedings of the Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering; ACM: Virtual Event USA, November 8 2020; pp. 1433–1443.

59. Iyer, S.; Konstas, I.; Cheung, A.; Zettlemoyer, L. Summarizing Source Code Using a Neural Attention Model. In Proceedings of the Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers); Association for Computational Linguistics: Berlin, Germany, 2016; pp. 2073–2083.
60. Grootendorst, M. BERTopic: Neural Topic Modeling with a Class-Based TF-IDF Procedure 2022.
61. Ramesh, A.; Dhariwal, P.; Nichol, A.; Chu, C.; Chen, M. Hierarchical Text-Conditional Image Generation with CLIP Latents 2022.
62. Hawthorne, C.; Stasyuk, A.; Roberts, A.; Simon, I.; Huang, C.-Z.A.; Dieleman, S.; Elsen, E.; Engel, J.; Eck, D. Enabling Factorized Piano Music Modeling and Generation with the MAESTRO Dataset 2019.
63. Amodei, D.; Olah, C.; Steinhardt, J.; Christiano, P.; Schulman, J.; Mané, D. Concrete Problems in AI Safety 2016.
64. Wei, J.; He, J.; Chen, K.; Zhou, Y.; Tang, Z. Collaborative Filtering and Deep Learning Based Recommendation System for Cold Start Items. *Expert Systems with Applications* **2017**, *69*, 29–39, doi:10.1016/j.eswa.2016.09.040.
65. Hendrycks, D.; Burns, C.; Kadavath, S.; Arora, A.; Basart, S.; Tang, E.; Song, D.; Steinhardt, J. Measuring Mathematical Problem Solving With the MATH Dataset 2021.
66. Lemley, M.A.; Casey, B. Fair Learning. *SSRN Journal* **2020**, doi:10.2139/ssrn.3528447.
67. European Commission Proposal for a Regulation Laying down Harmonised Rules on Artificial Intelligence Available online: <https://digital-strategy.ec.europa.eu/en/library/proposal-regulation-laying-down-harmonised-rules-artificial-intelligence> (accessed on 4 January 2025).
68. Doshi-Velez, F.; Kim, B. Towards A Rigorous Science of Interpretable Machine Learning 2017.
69. Strubell, E.; Ganesh, A.; McCallum, A. Energy and Policy Considerations for Deep Learning in NLP 2019.
70. Guu, K.; Lee, K.; Tung, Z.; Pasupat, P.; Chang, M.-W. REALM: Retrieval-Augmented Language Model Pre-Training 2020.
71. Bengio, Y.; Louradour, J.; Collobert, R.; Weston, J. Curriculum Learning. In Proceedings of the Proceedings of the 26th Annual International Conference on Machine Learning; ACM: Montreal Quebec Canada, June 14 2009; pp. 41–48.
72. Tuan, N.T.; Moore, P.; Thanh, D.H.V.; Pham, H.V. A Generative Artificial Intelligence Using Multilingual Large Language Models for ChatGPT Applications. *Applied Sciences* **2024**, *14*, 3036, doi:10.3390/app14073036.
73. Bengio, Y.; Deleu, T.; Rahaman, N.; Ke, R.; Lachapelle, S.; Bilaniuk, O.; Goyal, A.; Pal, C. A Meta-Transfer Objective for Learning to Disentangle Causal Mechanisms 2019.
74. Leike, J.; Krueger, D.; Everitt, T.; Martic, M.; Maini, V.; Legg, S. Scalable Agent Alignment via Reward Modeling: A Research Direction 2018.
75. Mitchell, M.; Wu, S.; Zaldívar, A.; Barnes, P.; Vasserman, L.; Hutchinson, B.; Spitzer, E.; Raji, I.D.; Gebru, T. Model Cards for Model Reporting. **2018**, doi:10.48550/ARXIV.1810.03993.
76. Acemoglu, D.; Restrepo, P. Automation and New Tasks: How Technology Displaces and Reinstates Labor. *Journal of Economic Perspectives* **2019**, *33*, 3–30, doi:10.1257/jep.33.2.3.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.