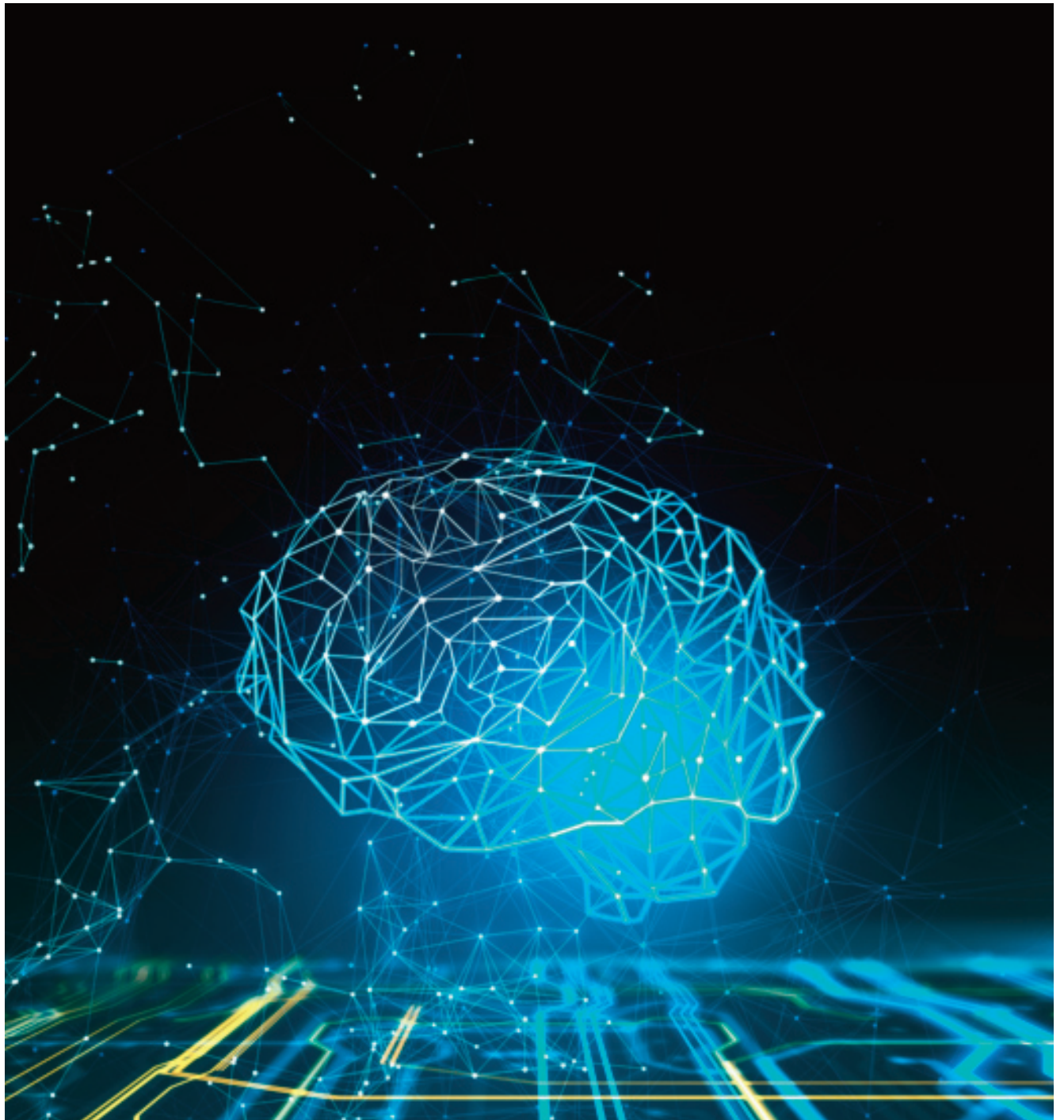


Executive summary

*“Artificial intelligence is not a substitute for human intelligence;
it is a tool to amplify human creativity and ingenuity”.*
Fei-Fei Li²²



LLM: context, definition and regulation

1. Generative Artificial Intelligence (GenAI), and within it Large Language Models (LLM), represents a significant advance in the field of AI, defining a new generation of human-machine interface where communication is done through natural language, with revolutionary applications in all sectors, including education, health, finance, and commerce. However, its development and use also pose significant risks and challenges that must be addressed.
2. LLMs are AI models trained to recognize, generate, translate and summarize large amounts of text. They use architectures such as transformers and are trained on large datasets to learn linguistic patterns and structures. Their effectiveness depends on the size in terms of the number of parameters, structure, diversity of training data and sophistication of their algorithms.
3. LLMs have evolved very rapidly, from the first rule-based models to today's transformer-based models. Important milestones include the introduction of transformer architecture and self-healing mechanisms, and the first commercial LLMs such as GPT. The year 2023 was key, with increased accessibility, global contributions, and the proliferation of open source LLMs.
4. LLMs have numerous applications, such as content creation and enhancement, information analysis and organization, and task interaction and automation. With the emergence of multimodal LLMs, new possibilities are opening up for generating rich audiovisual content and interactive experiences.
5. Regulators are taking steps to address the risks and opportunities of AI, with initiatives such as the EU AI Act, the U.S. AI Bill of Rights and the Bletchley Declaration. Key requirements include transparency, privacy, fairness, security, accountability and human oversight.

LLM: development and deployment

6. LLM development involves several critical components and decisions, such as data selection and preprocessing, tokenization and embedding, pre-training, quantization, and fine-tuning. In particular, the high cost of training often leads to the decision to use a pre-trained model or an open-source model, and to limit fine-tuning to data relative to the application being developed. Implementation requires integration, monitoring, and ethical and legal considerations.
7. Training models is a crucial aspect that influences their effectiveness. Factors such as the quantity and quality of the training data, the model architecture and the learning algorithms used can significantly impact the performance and generalization of an LLM.
8. The most common architecture for LLMs are transformers, which use self-learning mechanisms that allow the model to find relationships between different parts of the text, process them, and generate new text. They have demonstrated exceptional performance in a variety of natural language processing tasks. Variants and extensions aim to improve their efficiency and scalability.

²²Fei-Fei Li (b. 1976). Co-director of the Stanford Institute for Human-Centered Artificial Intelligence and IT Professor at the Graduate School of Business, known for creating ImageNet and AI4ALL, a non-profit organization working to increase diversity and inclusion in the field of artificial intelligence.

9. LLMOps is a methodology for managing the entire LLM lifecycle, addressing challenges such as managing large volumes of data, scaling computational resources²³, monitoring and maintenance, versioning, and reproducibility.
10. Key challenges for LLMs include biases and hallucinations, lack of explainability and transparency, data quality and accessibility, privacy and security issues, and high resource consumption. There are also challenges of dependency, risk of malicious use, intellectual property issues, and scalability.

LLM: Validation Framework

11. Validation of LLMs is essential to ensure their safe and responsible use, and it is appropriate to take a broad perspective covering the various risks involved. A multi-dimensional validation framework should cover aspects such as model risk, data management, cybersecurity, legal and operational risks, ethics and reputation.
12. LLM validation should be articulated through a combination of quantitative metrics and human judgment techniques. The choice of techniques will depend on the characteristics of the use case, such as level of risk, public exposure, personal data processing and line of business.

13. Emerging trends in LLM validation include explainability²⁴, the using LLMs to explain other LLMs, attribution scoring, continuous validation, collaborative approaches, prompt engineering, ethical and regulatory alignment, and machine unlearning techniques.

Case study

14. The case study presented illustrates the application of a custom validation framework to a company's internal policy chatbot. The process involved defining the case, designing the validation approach, running quantitative and qualitative tests, and interpreting results.
15. The chatbot validation results showed satisfactory overall performance, with strengths in accuracy, consistency, adaptability and scalability. Areas for improvement were identified in the areas of explainability, bias mitigation and security. It was recommended to proceed with implementation, applying the suggested improvements and establishing a continuous monitoring and improvement plan.

Conclusion

16. In conclusion, LLMs have significant potential to transform multiple sectors, but their development and deployment also pose significant challenges in transparency, fairness, privacy and security. To reap the benefits of LLMs in a responsible way, it is crucial to establish a robust AI governance framework that comprehensively addresses these challenges, including a rigorous, multi-dimensional approach to validation that covers the entire lifecycle of the models. This is the only way to ensure that LLMs are reliable, ethical and aligned with the values and goals of organizations and society at large.

²³Management Solutions (2022). Auto Machine Learning, towards model automation.

²⁴Management Solutions (2023). Explainable Artificial Intelligence (XAI): challenges in model interpretability.

