

Artificial Intelligence & Statistics

Navigating potentials, pitfalls and use

LSS & STATEC – Economic Seminar

Thursday, December 21st 2023

Florian FELICE

Senior Data Scientist, AWS Generative AI Innovation Center

Doctoral Research, University of Luxembourg

Artificial Intelligence & Statistics

Navigating potentials, pitfalls and use

LSS & STATEC – Economic Seminar

Thursday, December 21st 2023

Florian FELICE

Senior Data Scientist, AWS Generative AI

Doctoral Research, University of Luxembourg



TRY FOR FREE : AI STUDIOS

www.deepbrain.io



Artificial Intelligence & Statistics

Navigating potentials, pitfalls and use

LSS & STATEC – Economic Seminar

Thursday, December 21st 2023

Florian FELICE

Senior Data Scientist, AWS Generative AI Innovation Center

Doctoral Research, University of Luxembourg

Agenda

- Introduction
- Generative AI capabilities
- Arising concerns
- Responsible AI
- Conclusion & discussion



Introduction



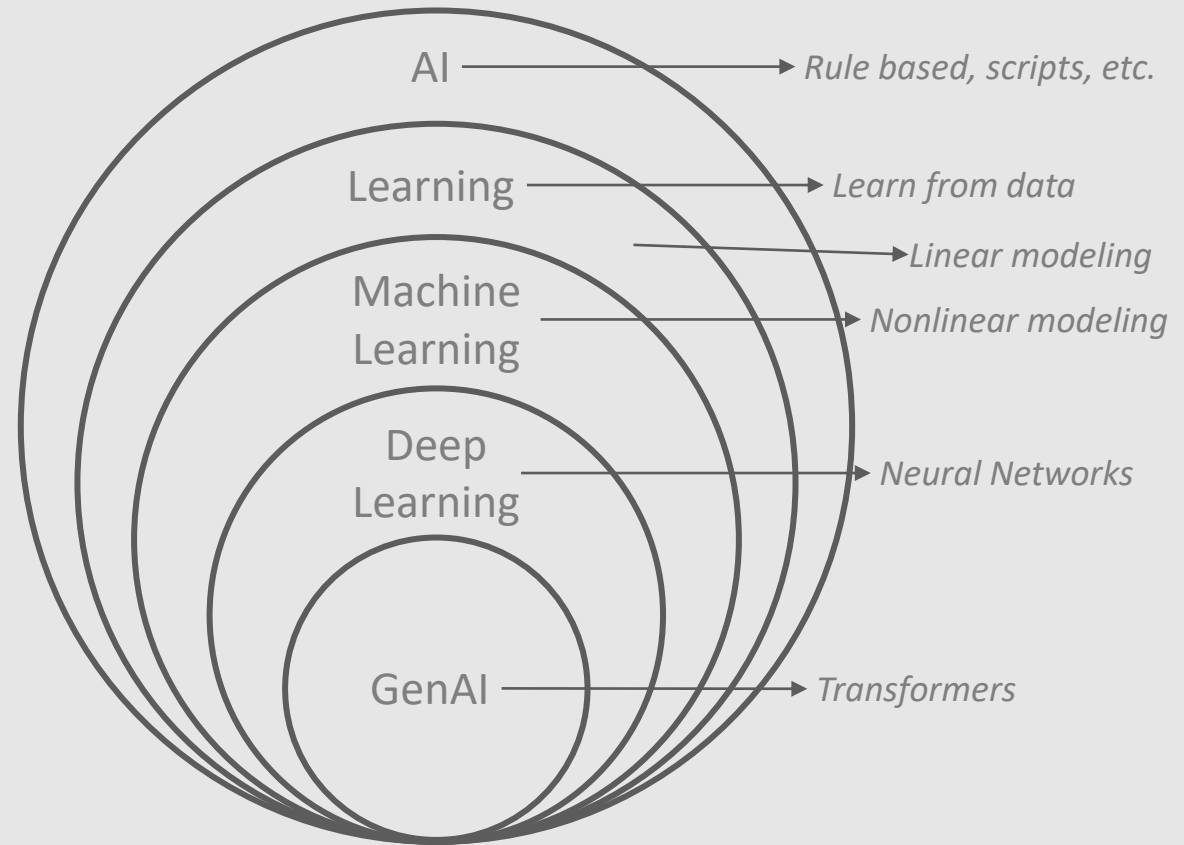
Artificial Intelligence



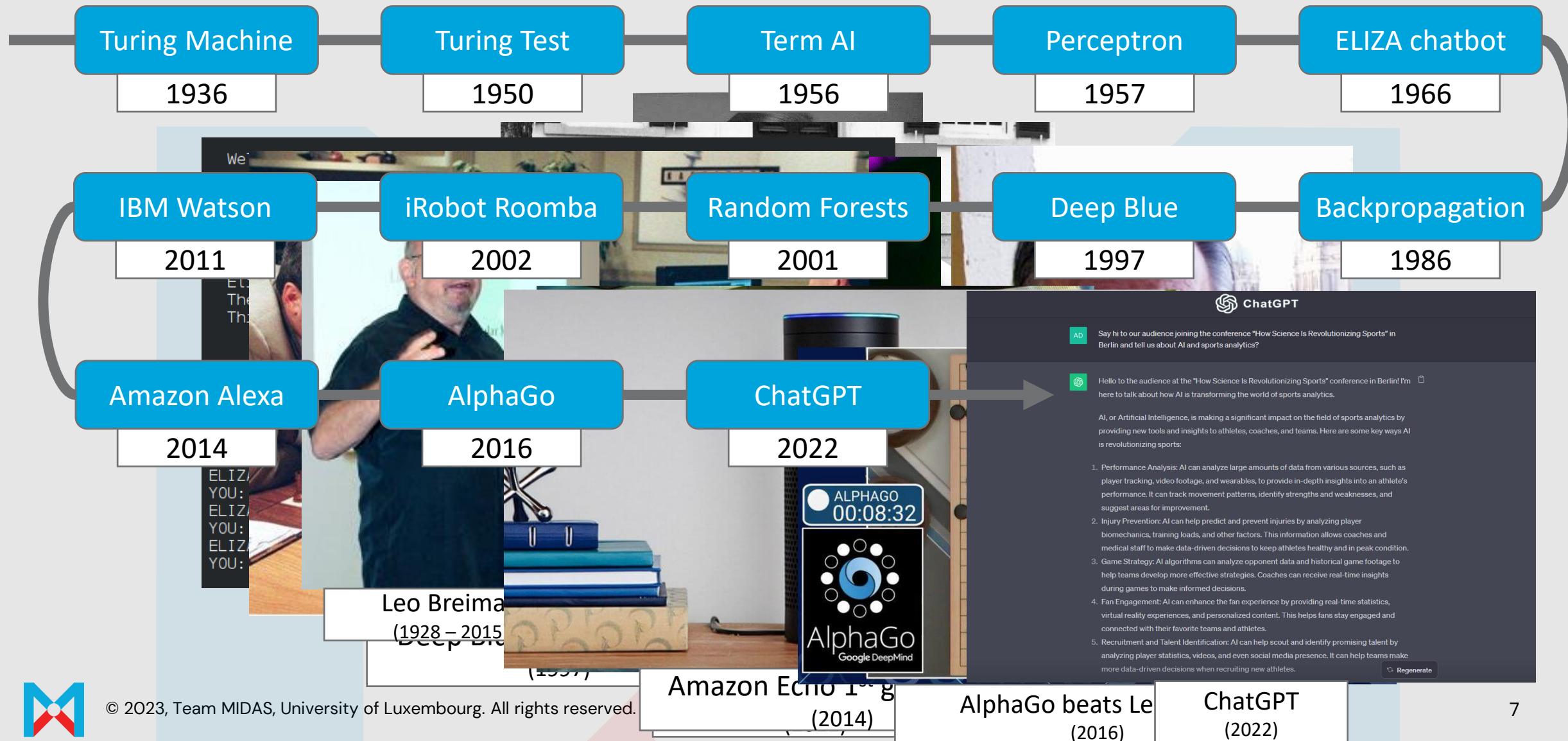
Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT press.



Felice, Florian, et al. "Statistically Enhanced Learning: a feature engineering framework to boost (any) learning algorithms." *arXiv preprint arXiv:2306.17006* (2023).



The history of AI



**Once upon a time, there was a
great fairy tale.**



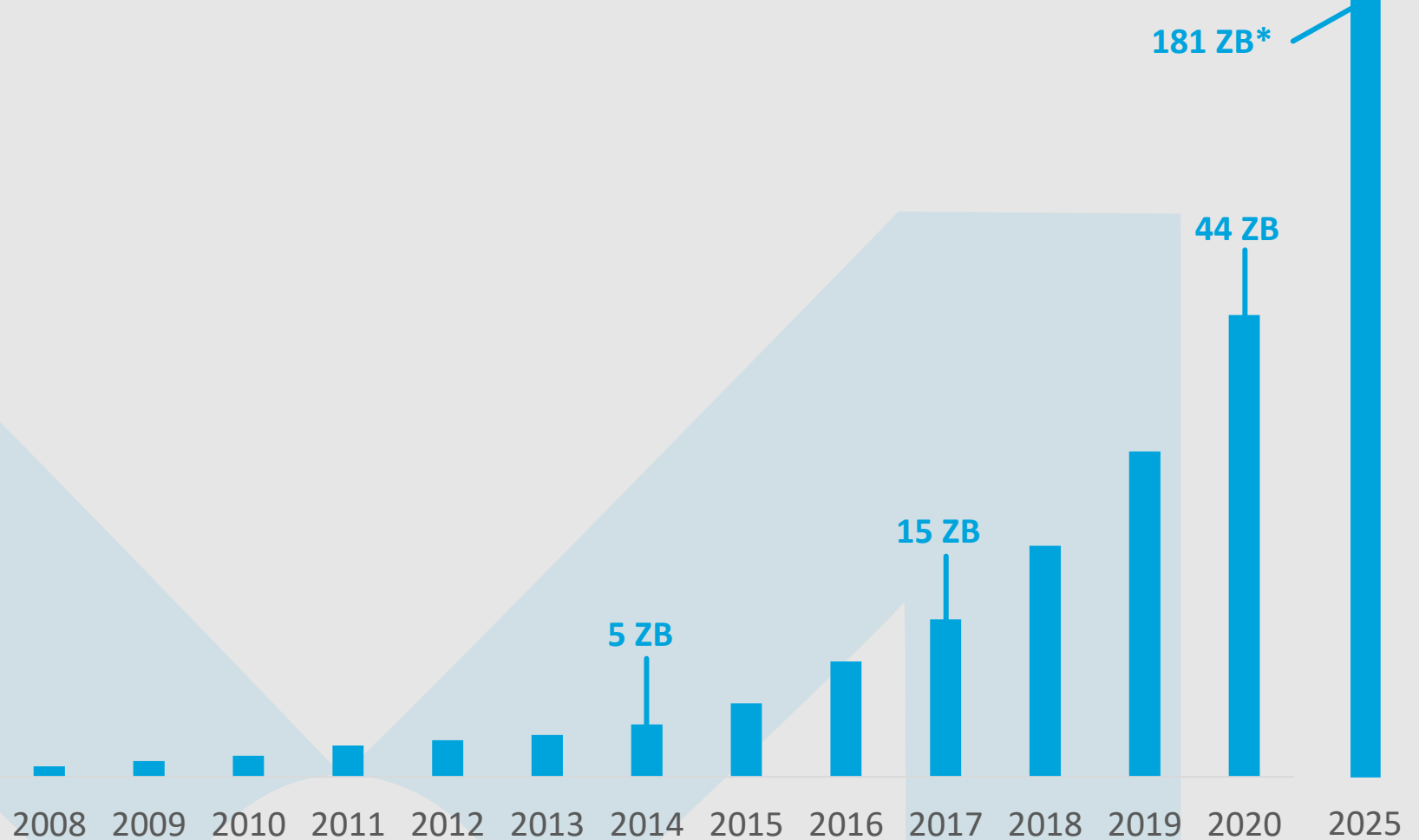
Evolution of stored data over time



1 ZB = 10^9 TB



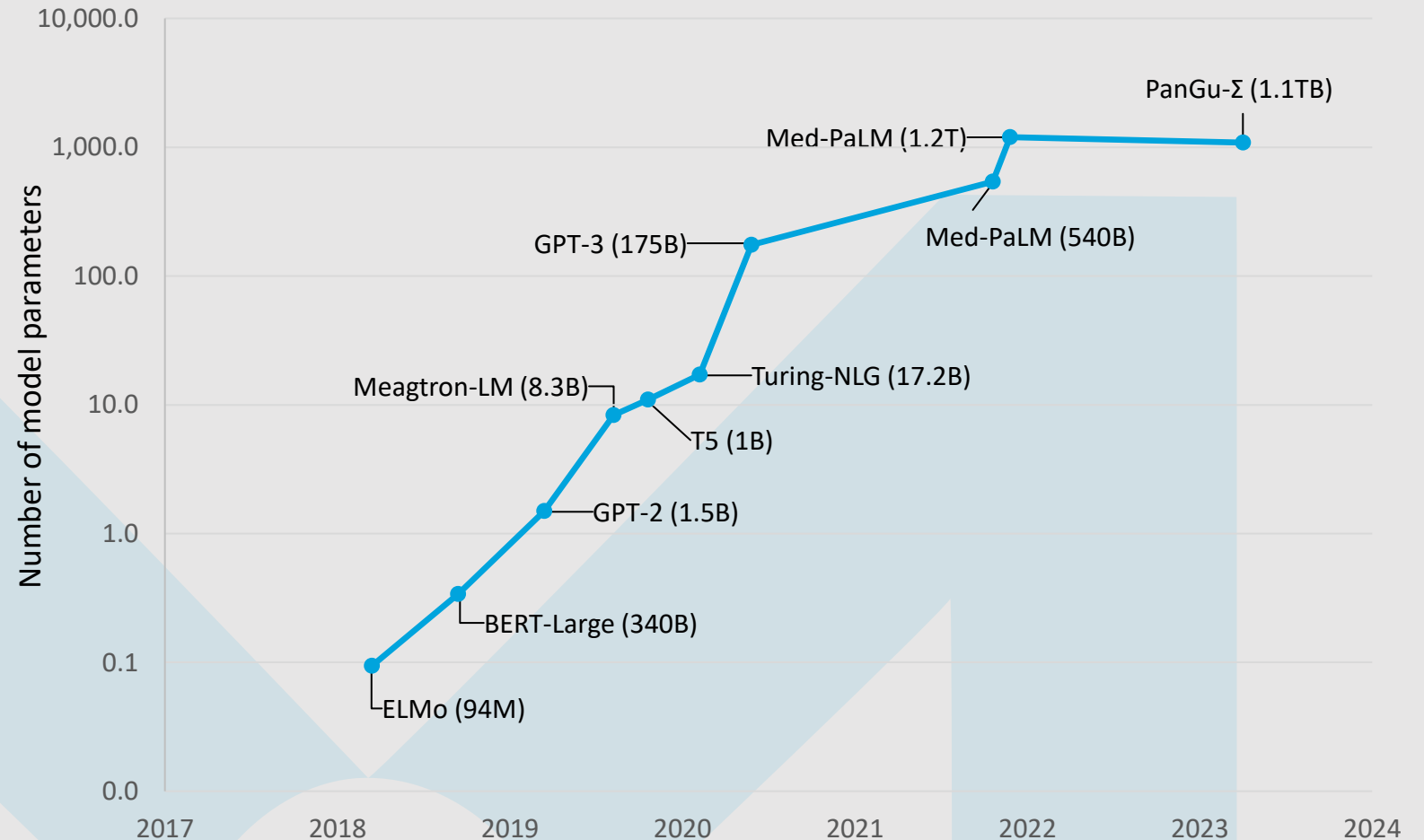
Hammad, Khalid Adam Ismail, et al. "Big data analysis and storage." *International conference on operations excellence and service engineering*. 2015.



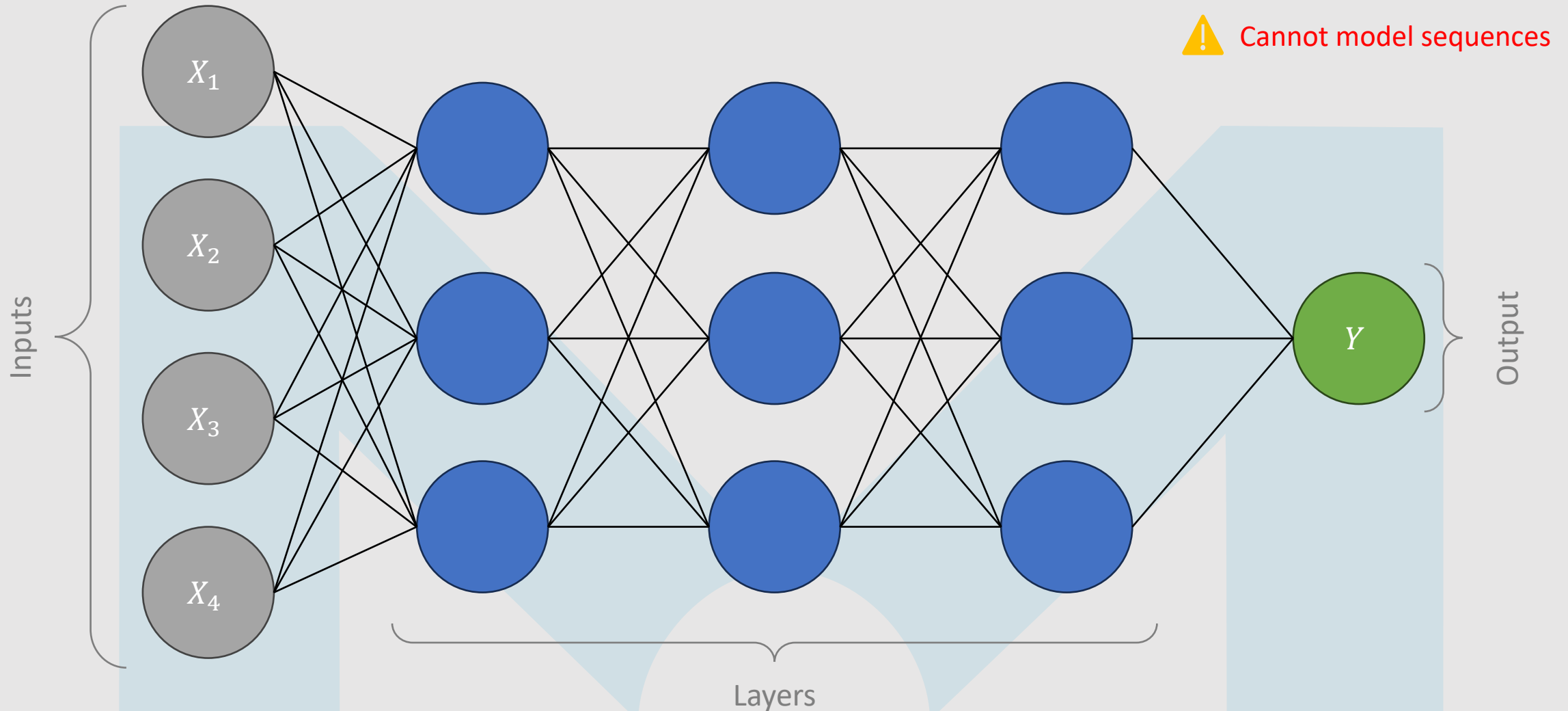
Evolution of model size



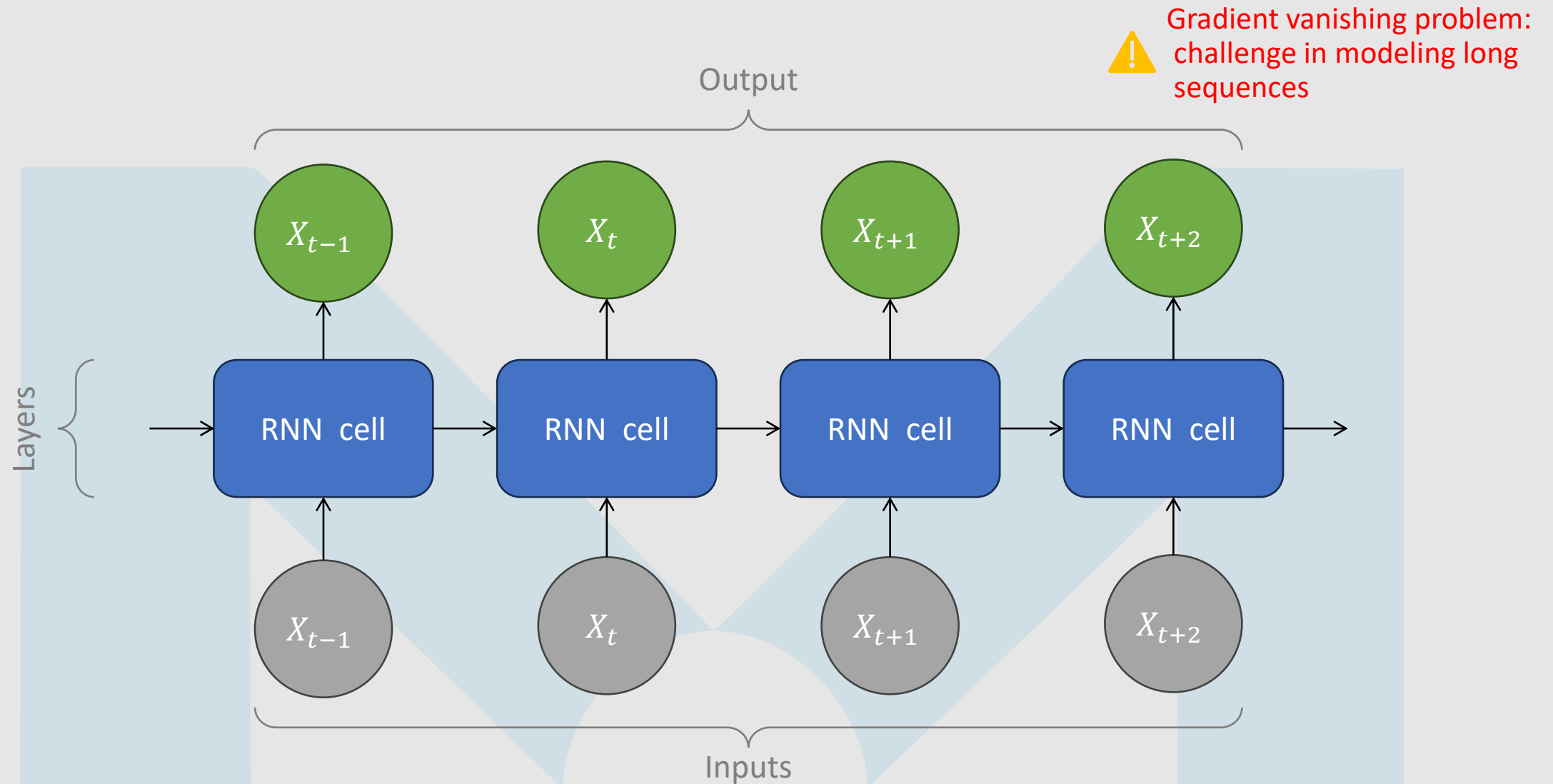
Simon, Julien. "Large language models: A new Moore's law."
Hugging Face blog (2021).



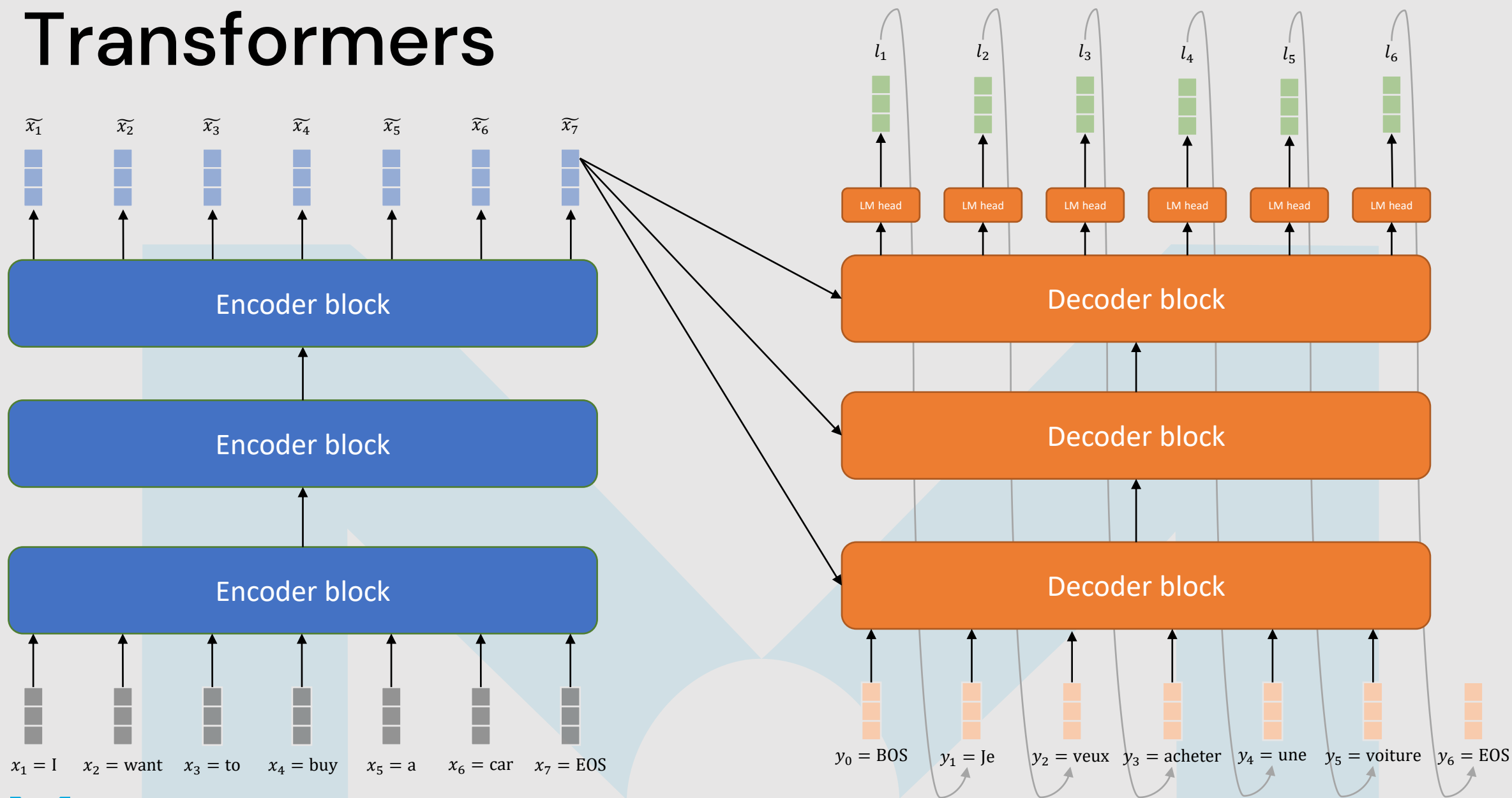
Multi-Layered Perceptron (MLP)



Recurrent Neural Nets (RNN)



Transformers



Generative AI capabilities



Generative AI for images

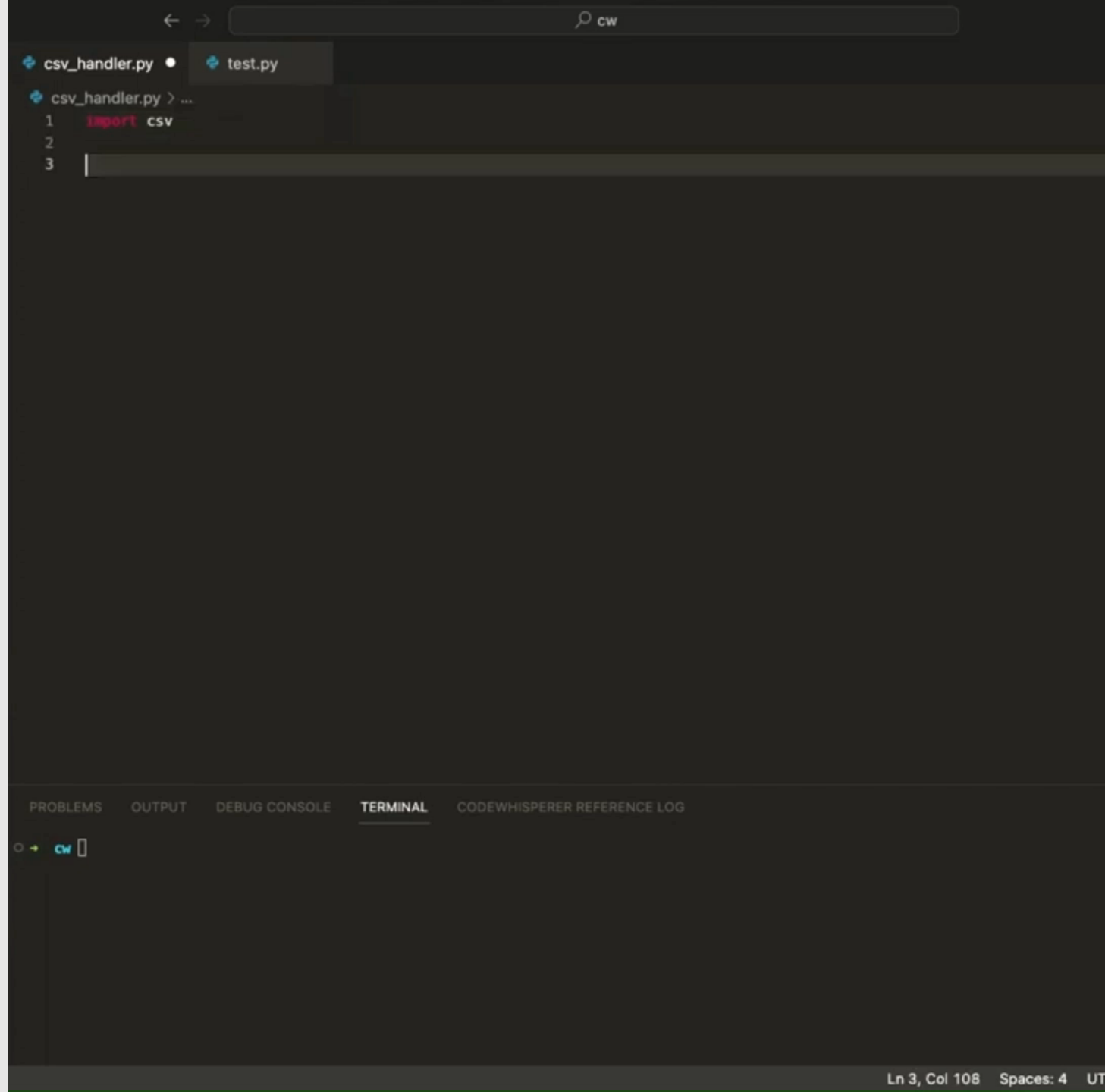
Prompt: *"Create a painting of a skateboarding cat in the style of Andy Warhol."*

Generated by Stable Diffusion XL via Amazon Bedrock



Generative AI for coding

- Support developers be more productive
- Help explain code
- Help debug



Generative AI for document summary

- Retrieval Augmented Generation

▼ Prompt

CLEAR PROMPT

Write a prompt and then click Submit

SUBMIT

▼ Response

☒ Markdown

The model will generate a response after you click Submit

Model may display inaccurate or offensive information that doesn't represent Google's view. Not all languages are supported. [Learn more](#).



Generative AI capabilities

Sentiment
analysis

Language
understanding

Text classification



Image
editing

Q&A

Text to
speech

Jobs destruction

- Speculations on jobs being lost because of GenAI technologies
- Can be a productive helper tool (e.g. AWS CodeWhisperer)
- Can expect a similar revolution as with online search engines
- Education on how to use such tools



Arising concerns



Plagiarism, cheating & Intellectual Property

- Models are trained on large amount of data
- Can reproduce the training data (potentially including confidential or copy-righted content)
- Challenge in education systems of students using such tools
- Ethical, legal and educational considerations



Toxicity


- Offensive/inappropriate content
- Discrimination against groups/individuals
- Difficult to draw boundary between restricting toxic content and censorship



Hallucinations

- Outcome that sounds/looks reasonable but verifiably incorrect
- Over-creativity cannot be linked to online and verified content
- Root cause: model is predict next word from some distributions



A close-up of Mr. Bean's face. He has a smug, knowing expression with a slight smirk and wide, staring eyes. He is wearing his signature brown tweed jacket over a white shirt.

**You have
NO CHOICE**

**NOT
CONVINCED??**

Regulations around AI

- Level 1: Social Principles (2018 –)
 - [OECD Principles on AI](#)
- Level 2: National Frameworks and Social Guidelines (2020 –)
 - [NIST AI RMF](#) (USA)
 - [GDPR, Regulatory Framework proposal on AI](#) (EU)
 - [National AI Strategy](#) (UK)
- Level 3: Technical Guidelines and Standardization (2021 –)
 - [ISO/IEC JTC 1/SC 42 AI Standards](#)
 - [IEEE AI Standards](#)

Responsible AI

An overview



“AI that is innovative and trustworthy and respects human rights and democratic values”

Source: [OECD](#)

Some Responsible AI dimensions

Fairness

How a system impacts different subpopulations of users (e.g. by gender, ethnicity)

Explainability

Mechanisms to understand and evaluate the outputs of an AI system

Robustness

Mechanisms to ensure an AI system operates reliably

Privacy and Security

Data used in accordance with privacy considerations, and protected from theft and exposure

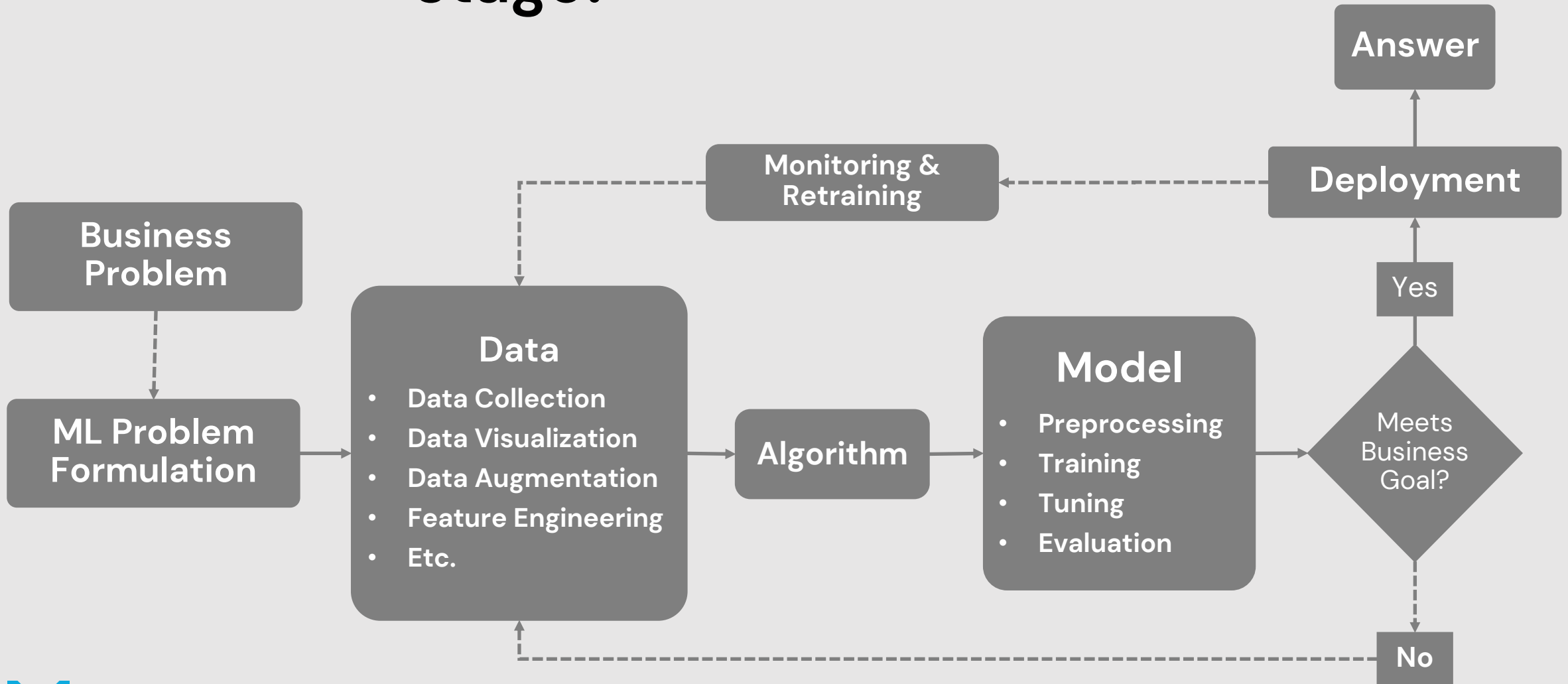
Governance

Processes to define, implement and enforce responsible AI practices within an organization

Transparency

Communicate information about an AI system so stakeholders can make informed choices about their use of the system

ML lifecycle... must have RAI components at any stage!



SO...

WHAT NOW???

makeameme.org



A responsible use of GenAI



Fairness in Generative AI


Once upon a time, a doctor called John (0.7)
Jane (0.2)

Dr. Hanson studied the patient's chart carefully, and he... (0.8)
she... (0.1)

Fairness in Generative AI

Once upon a time, a doctor called John (0.4)
Jane (0.35)

Dr. Hanson studied the patient's chart carefully, and he... (0.4)
she... (0.6)

 Not scalable!

Fairness in Generative AI


- Careful preparation of training data
 - Models will reflect the tone of the training data
- Human annotated data
 - Training methods such as RLHF or model alignment
- Human tested model
 - Check how model behaves for specific prompts


Fairness in Generative AI

- Limit the toxicity of a model by:
 - Working on training step
 - E.g. Reinforcement Learning with Human Feedback (RLHF)
- Giving instructions to the model
 - Prompt engineering
- Training guardrail models
 - Identify and filter out undesired content
- Understanding the behavior of your model
 - Explainability

Explainability in Generative AI

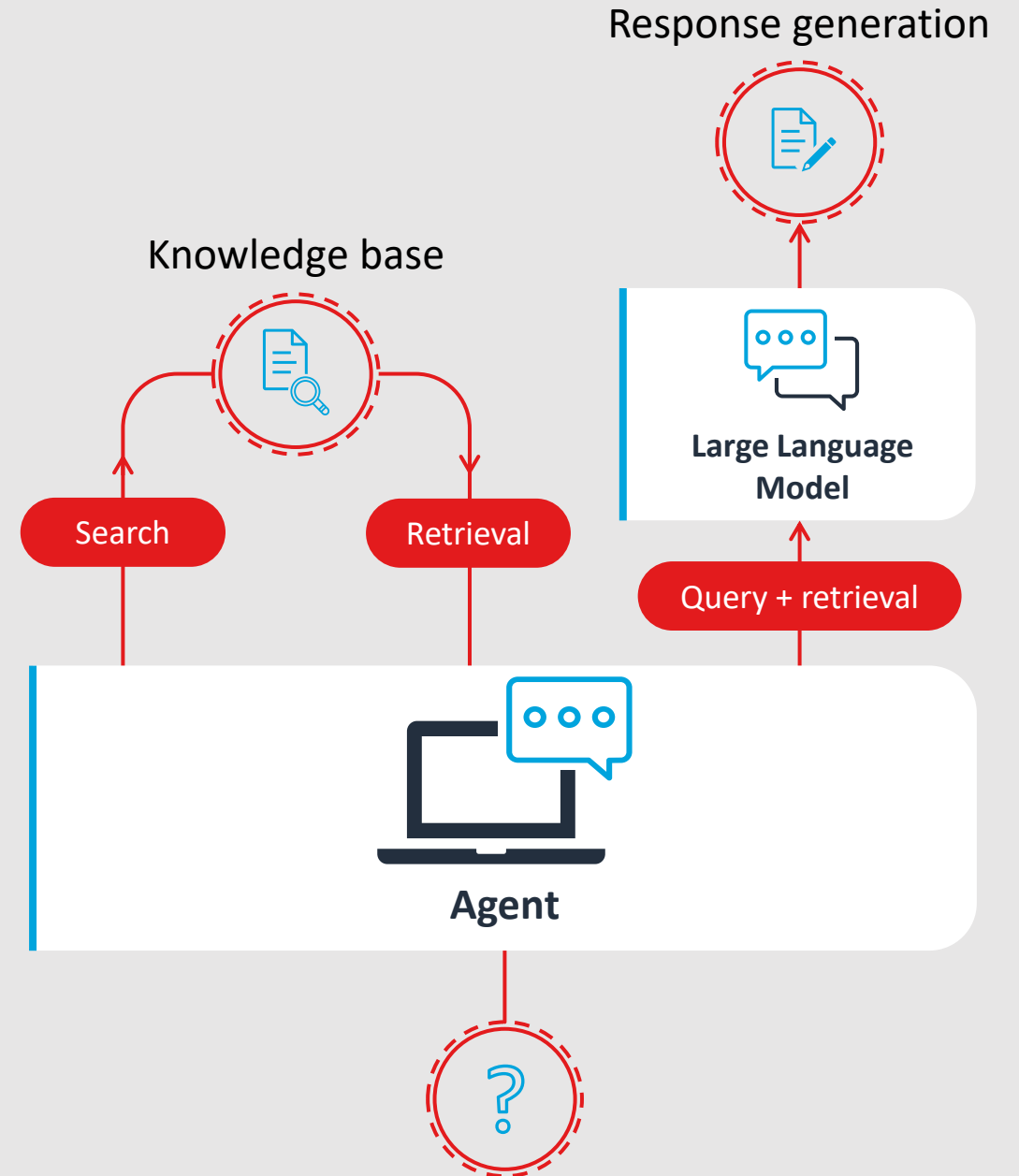
- Active field of research
- Requires definition of explanations for text generation
- Some literature available for classification models
- Chain-of-thoughts prompting
- Available AWS service: [SageMaker Clarify](#)

 Zafar, Muhammad Bilal, et al. "On the lack of robust interpretability of neural text classifiers." *arXiv preprint arXiv:2106.04631* (2021).

 Kokalj, Enja, et al. "BERT meets shapley: Extending SHAP explanations to transformer-based classifiers." *Proceedings of the EACL Hackashop on News Media Content Analysis and Automated Report Generation*. 2021.

Robustness in Generative AI

- **Human-in-the-Loop**
Expert feedback during development
Collection of user feedback in production
- **Retrieval Augmented Generation**



Privacy and security for Generative AI

- Data privacy
 - Careful preparation of input data
 - Remove personal/confidential information
- Data security
 - E.g. data encryption on S3
- Service security
 - Model endpoint on VPC
 - E.g. restrict access to internet, isolated subnet

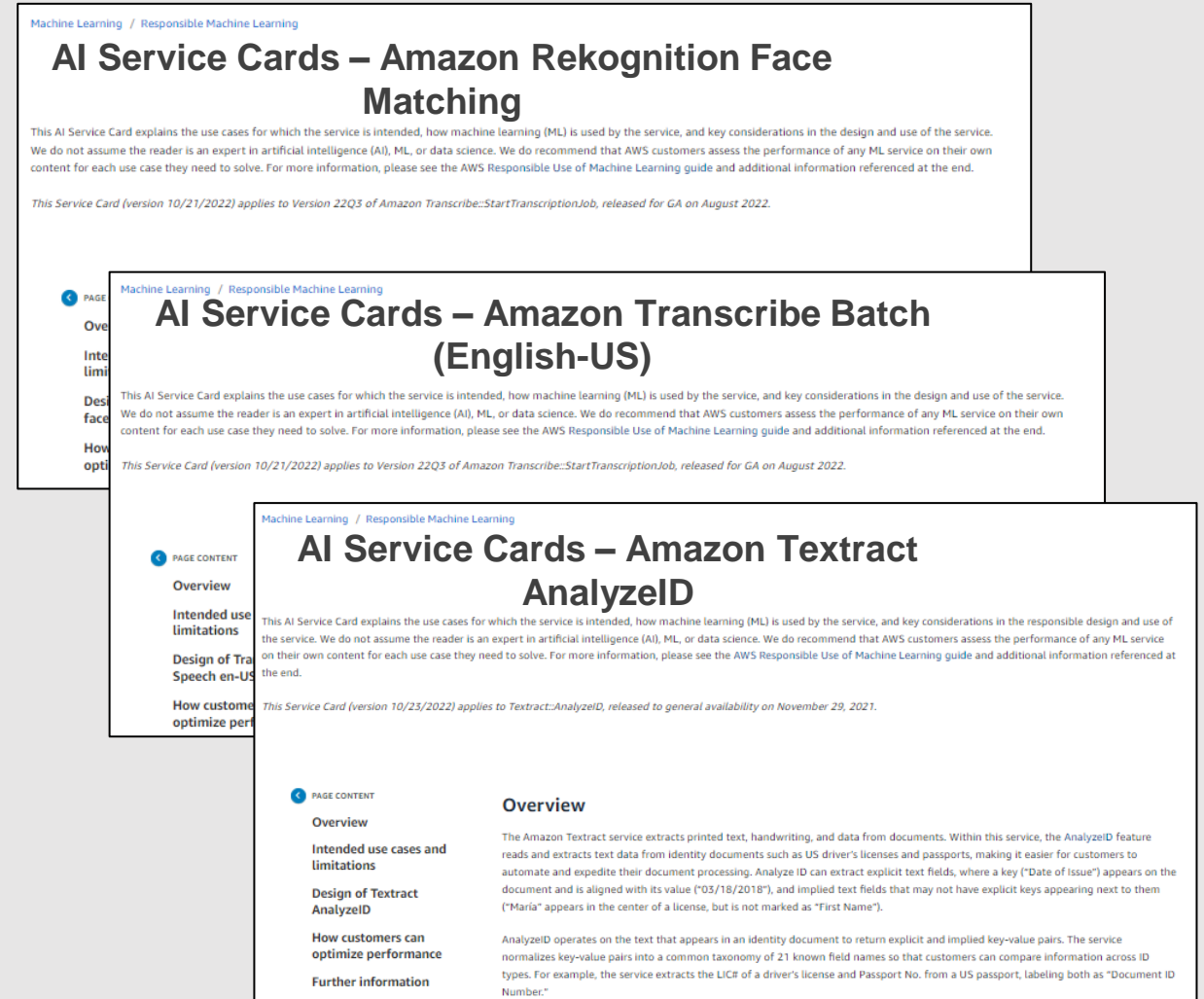
Governance for Generative AI

- Accountability and management of the service/tool
- Ensure RAI practices are carried out among stakeholders
- Regulatory frameworks (e.g. NIST – US) focus on governance
 - They insist on documentation
 - Accountability and permissions at each step of the model lifecycle

Transparency for Generative AI

New transparency resource to advance responsible AI

- Documents the intended use cases and fairness considerations of our AWS AI services
- Reflects our comprehensive development process
- Three new AI Service Cards published



Conclusion



Conclusion

- Generative AI is not all you need
 - It will not replace traditional ML: it will complement it
- Prepare your data carefully
- Responsible AI is everywhere
- Stay tuned! Research is in progress
- Read more on [Amazon Science](#)
 - [Responsible AI in the generative era](#) by [Michael Kearns](#)

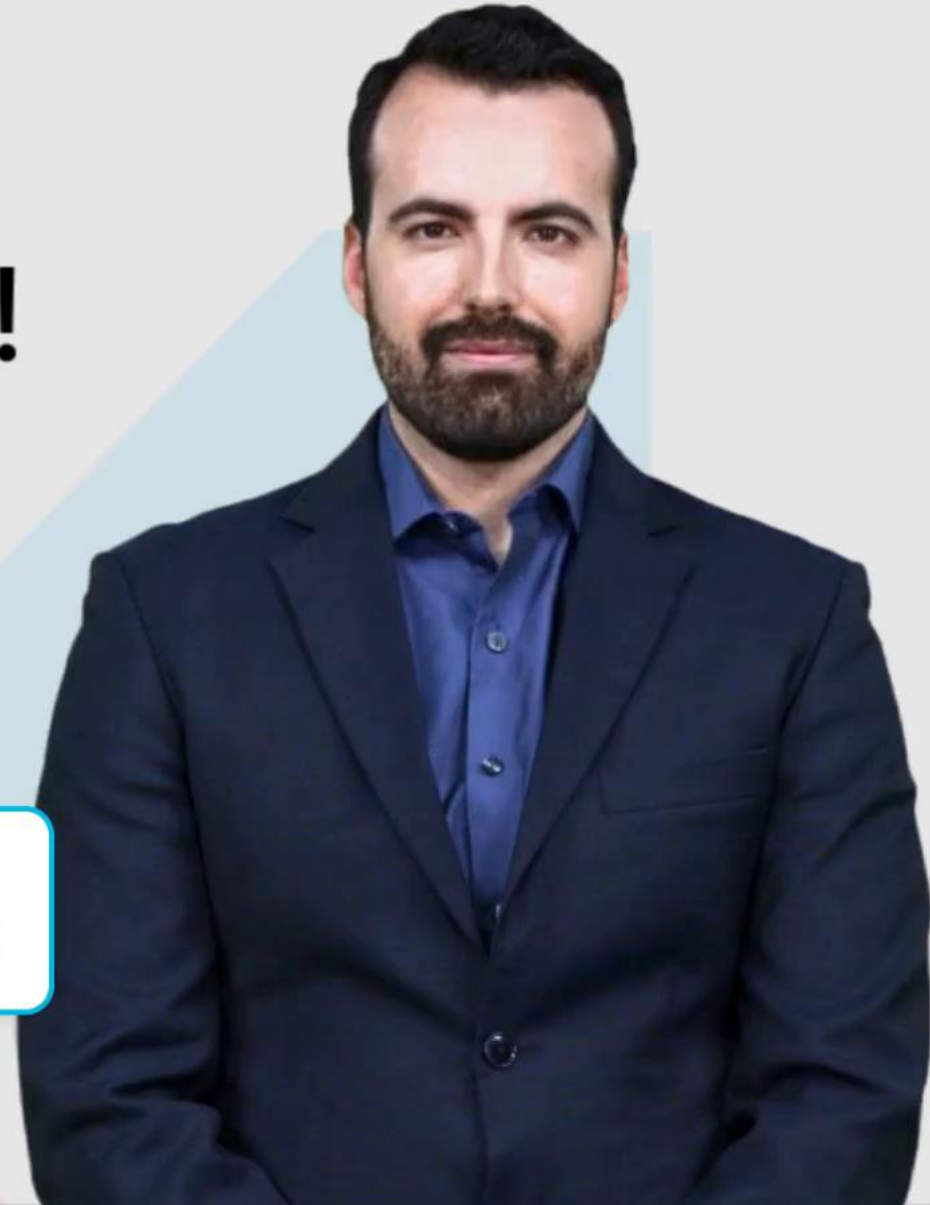


Thank you!



TRY FOR FREE : AI STUDIOS

www.deepbrain.io



Thank you!

