



# Generative AI Demystified

Amanda Saunders | NVIDIA | GTC 2024

# Where do I begin?

ChatGPT 4 ▾



How can I help you today?

**Help me pick**  
an outfit that will look good on camera

**Plan an itinerary**  
to experience the wildlife in the Australian outback

**Recommend activities**  
for a team-building day with remote employees

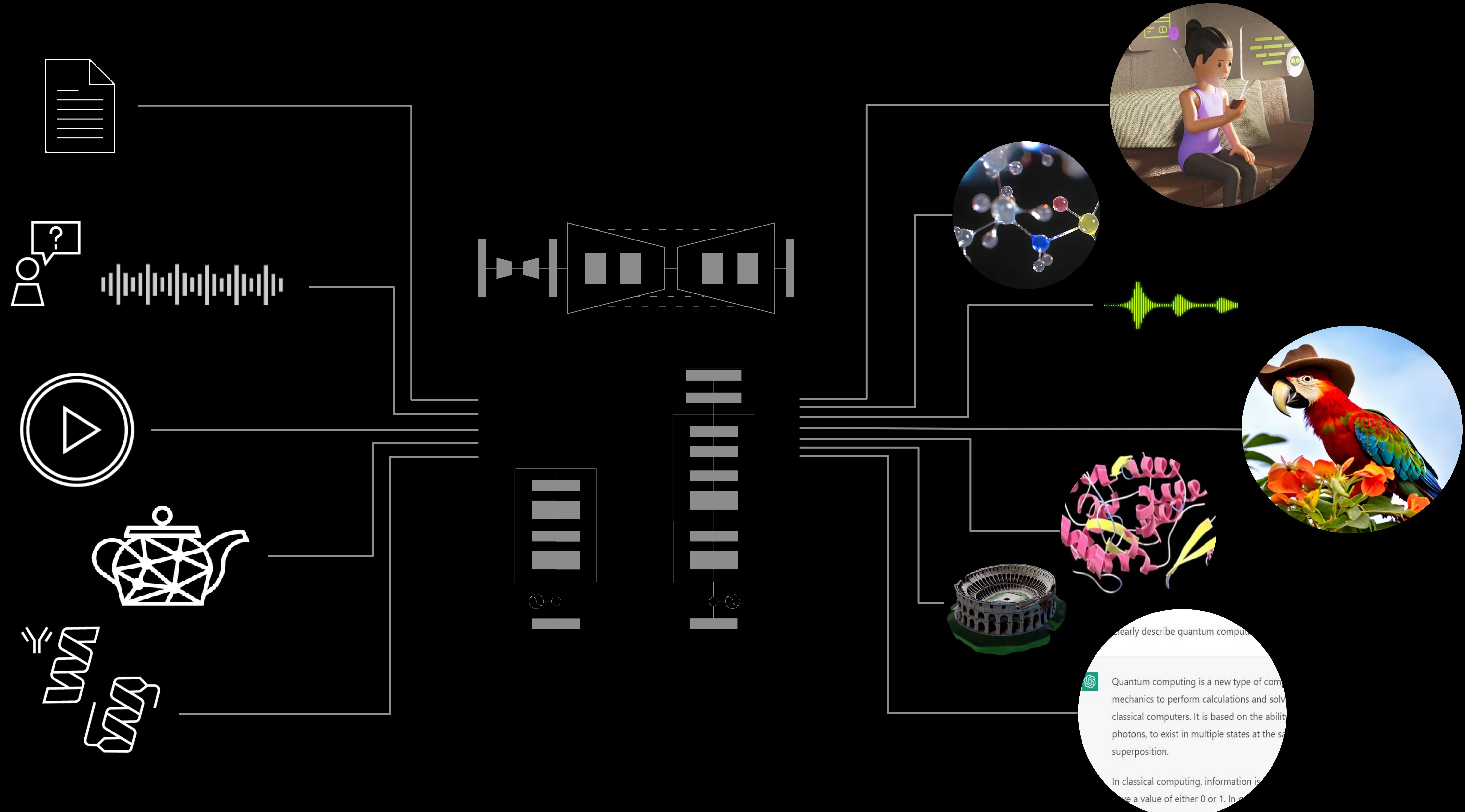
**Explain airplane turbulence**  
to someone who has never flown before

Message ChatGPT...

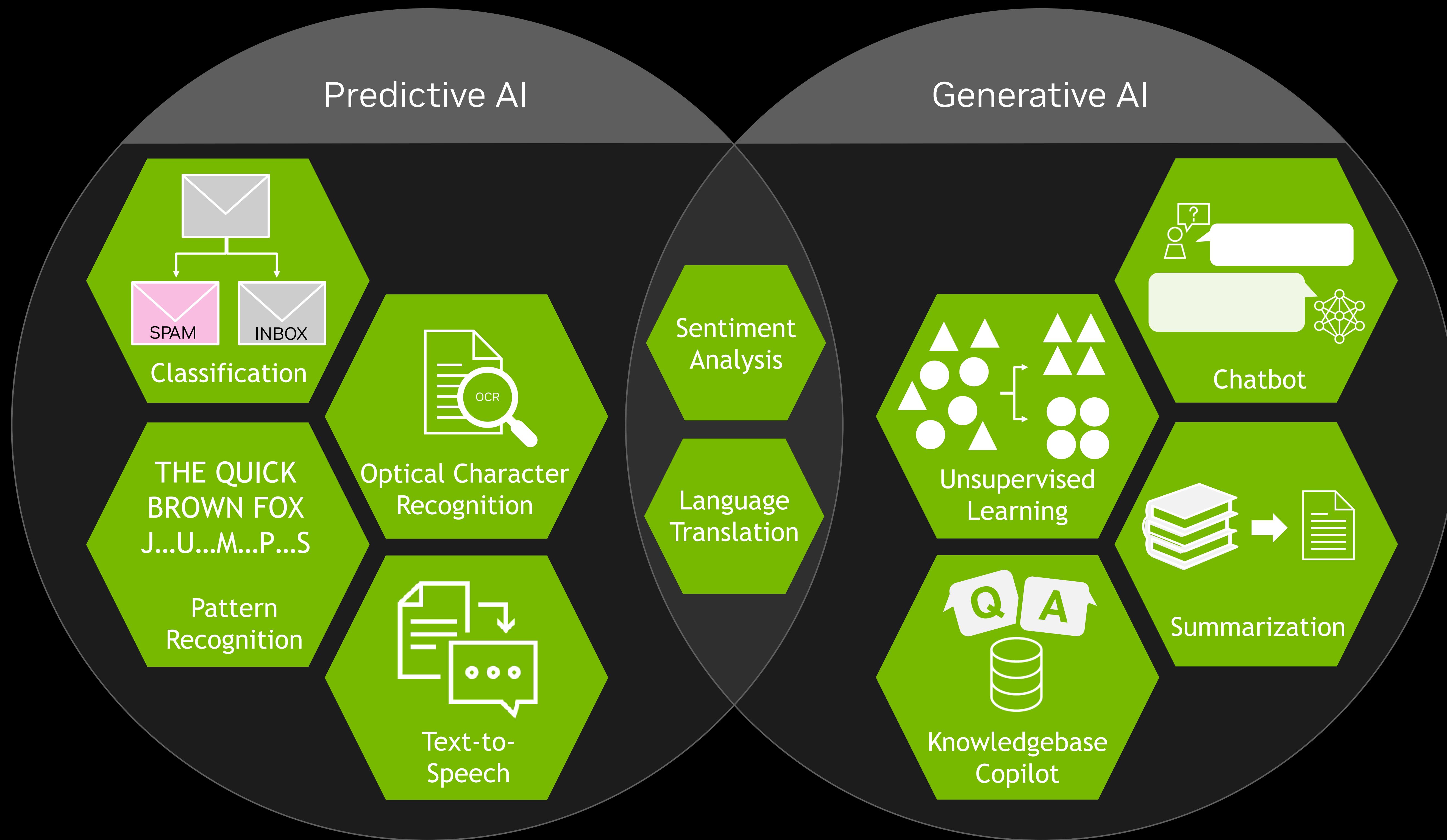
I

ChatGPT can make mistakes. Consider checking important information.

# What is Generative AI?



# When to Use Generative AI to Solve Challenges



Predictive AI focuses on understanding historical data and making accurate predictions

Generative AI creates new data based on patterns and trends learned from training data

# Generative AI Journey



## Explosion

ChatGPT gets announced late in 2022, gaining over 100 million users in just two months. Users of all levels can experience AI and feel the benefits firsthand.



## Experimentation

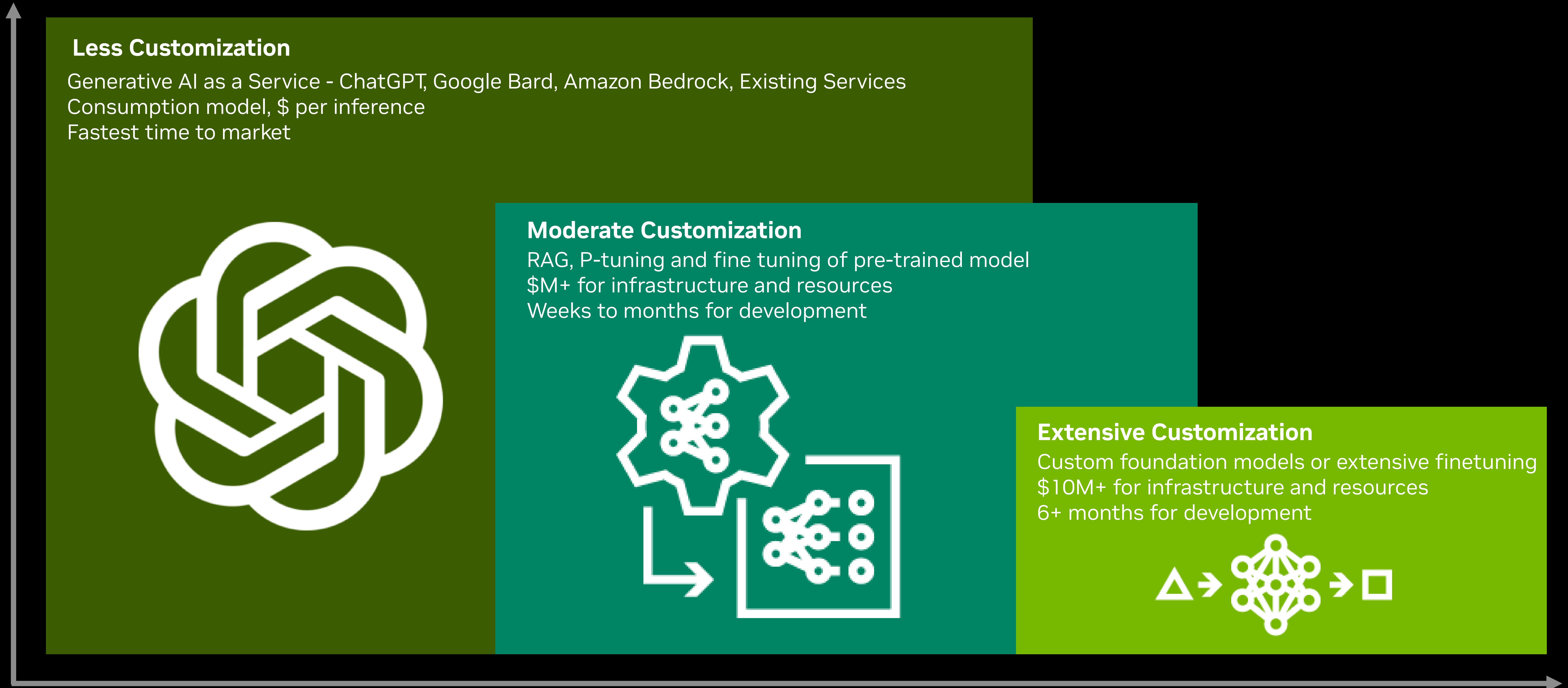
Enterprise application developers kick off POCs for generative AI applications with API services and open models including Llama 2, Mistral, NVIDIA, and others.



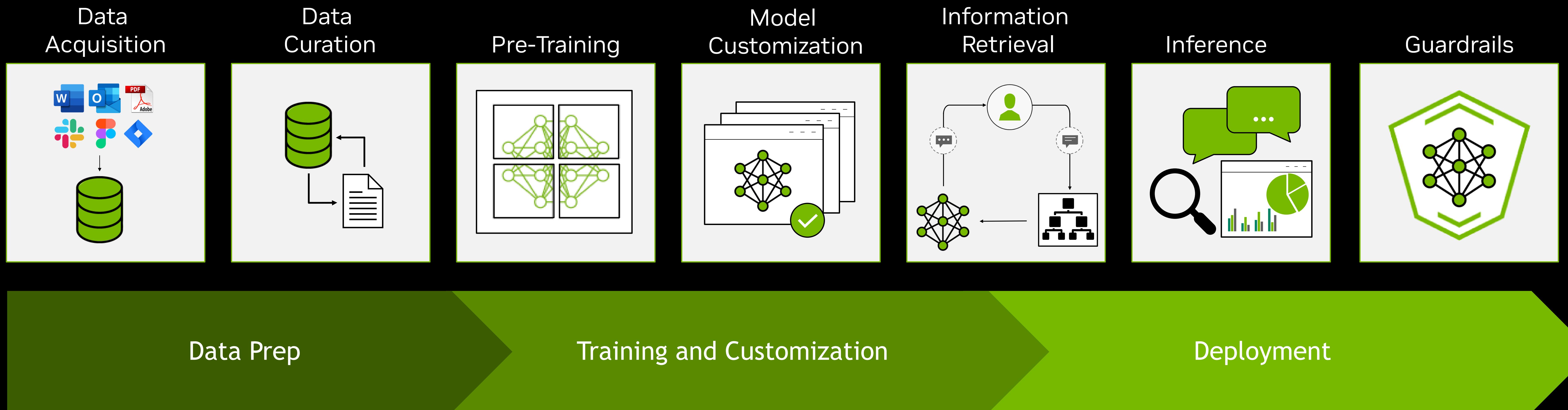
## Production

Organizations have set aside budget and are ramping up efforts to build accelerated infrastructure to support generative AI in production.

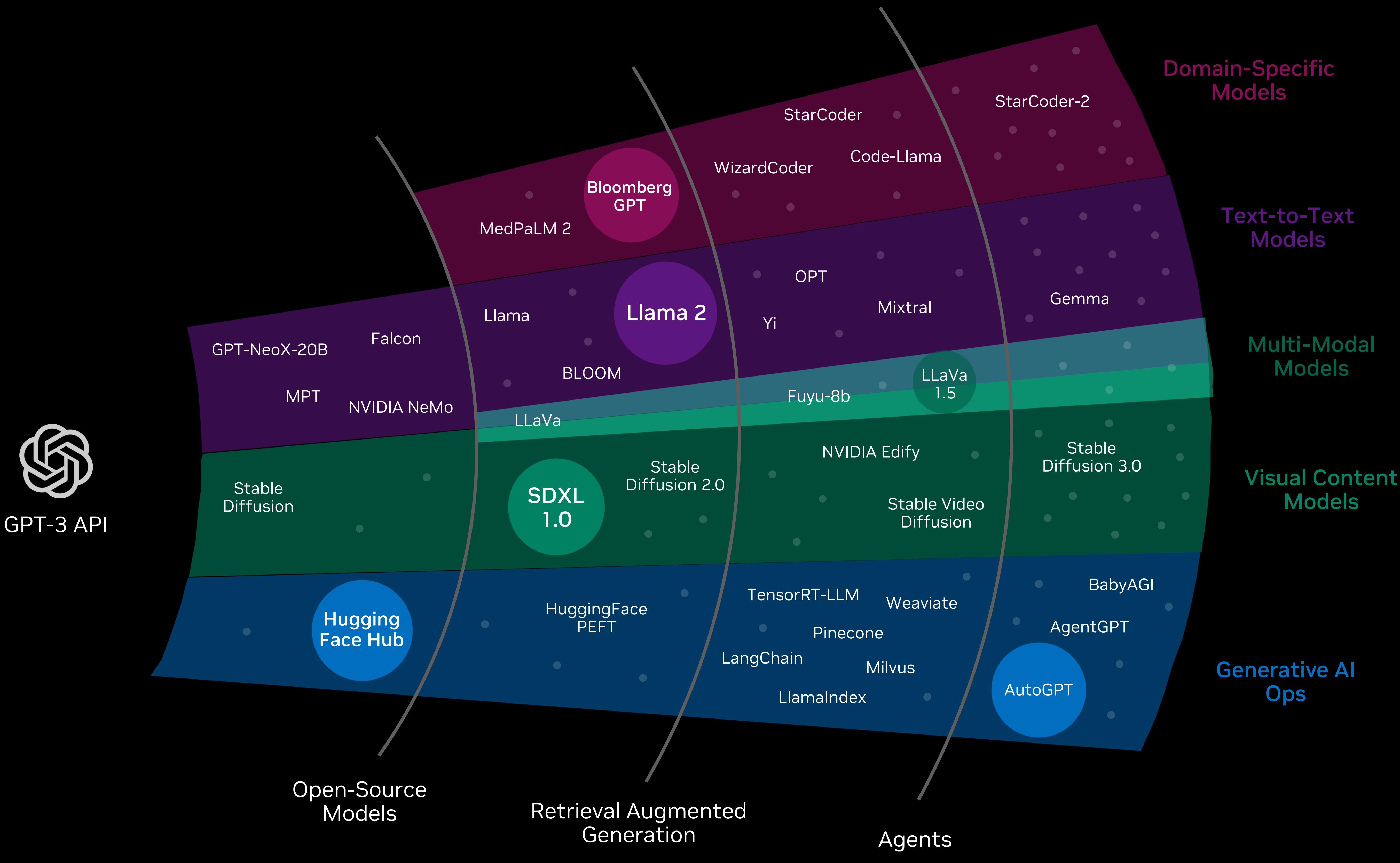
# Generative AI Is a Spectrum



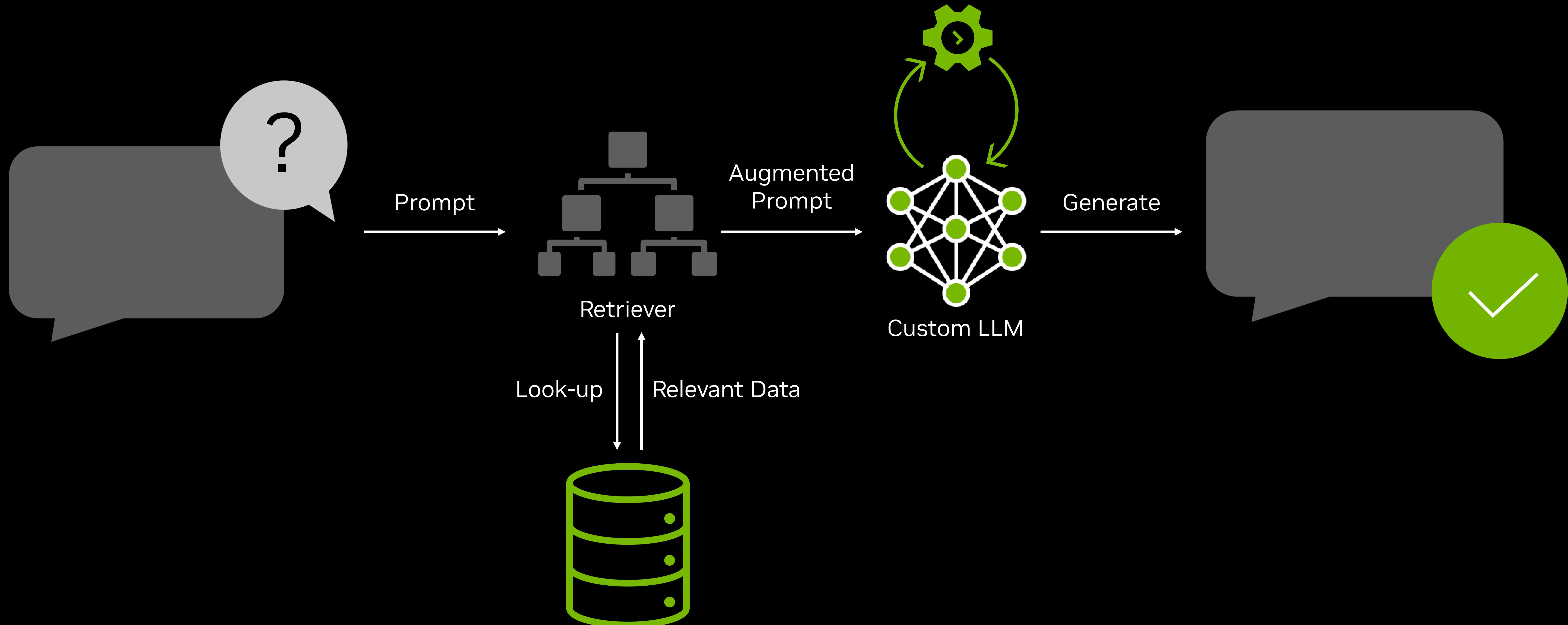
# Building Generative AI Applications



# Explosion of Generative AI Models



# Connecting LLMs to Enterprise Data with Retrieval Augmented Generation

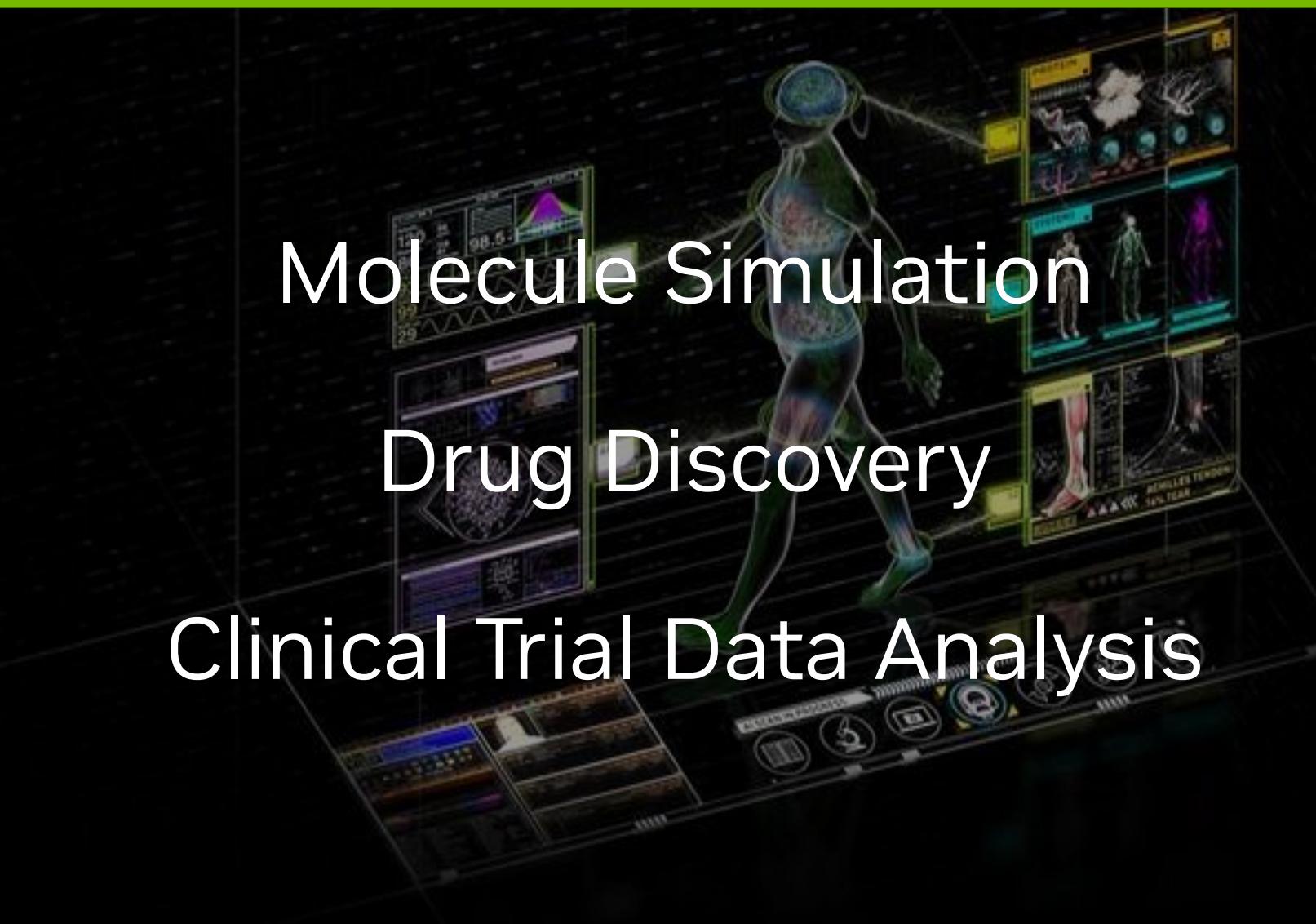


# Generative AI Adoption Across Industries

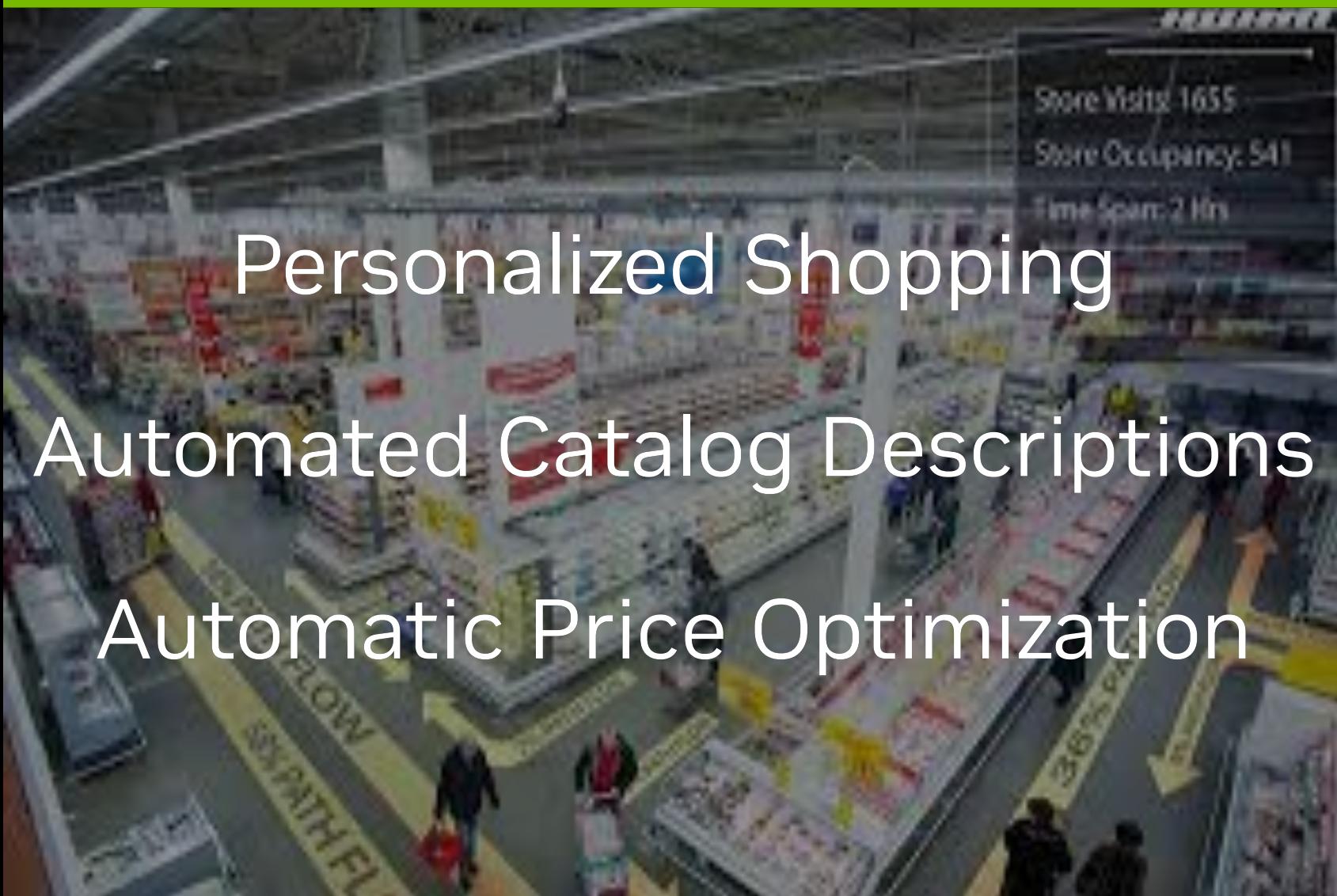
## Finance



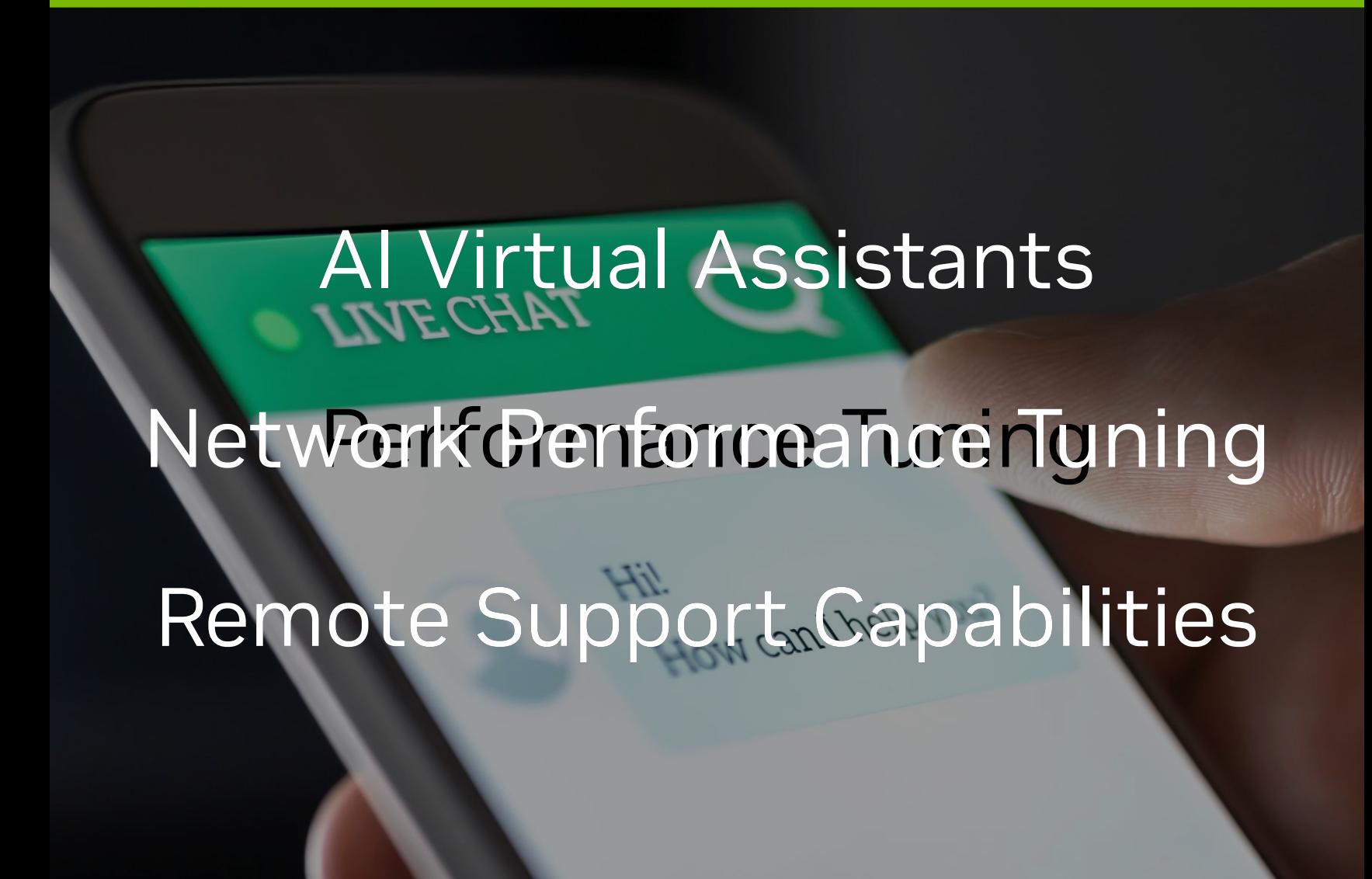
## Healthcare



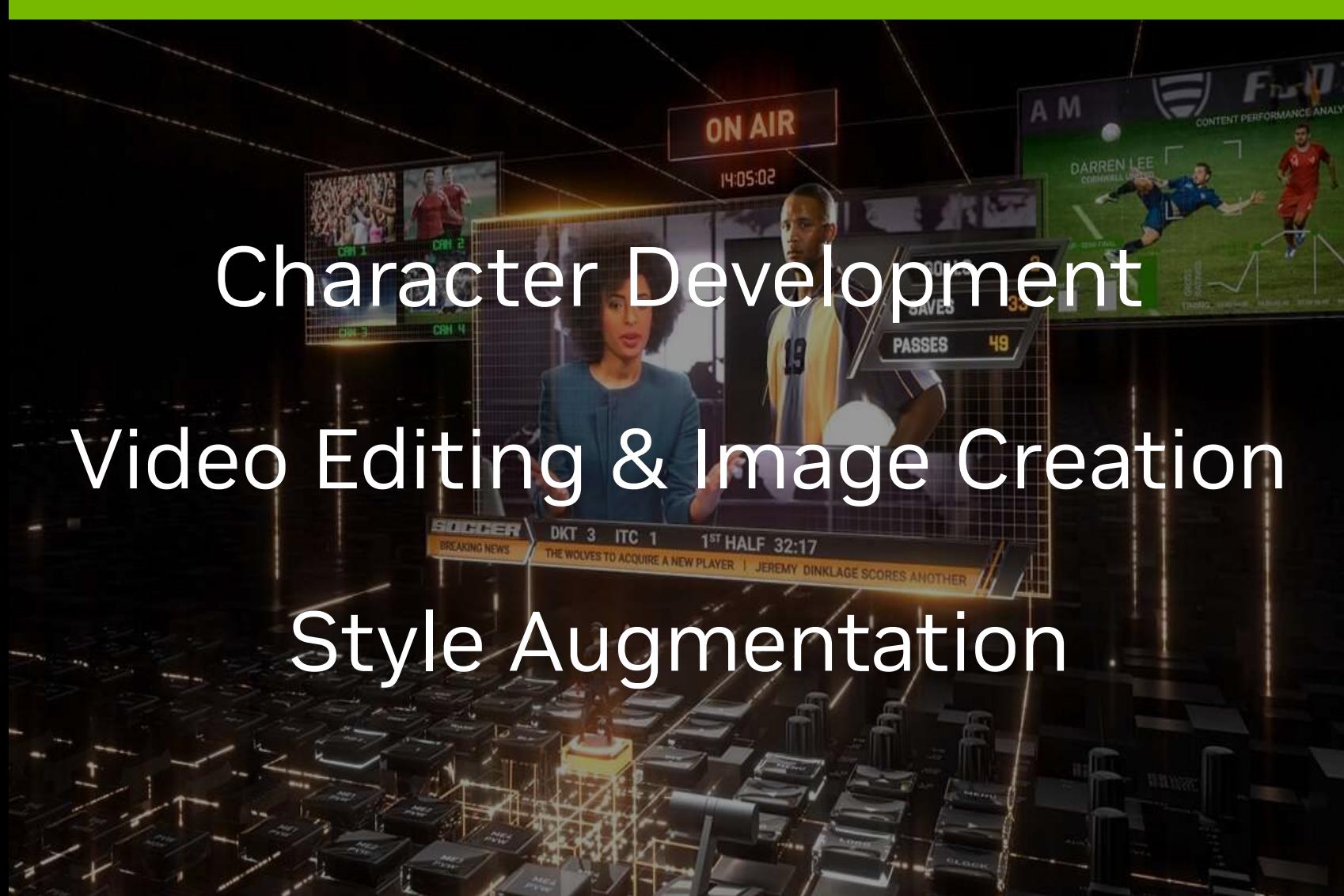
## Retail



## Telecommunications



## Media & Entertainment



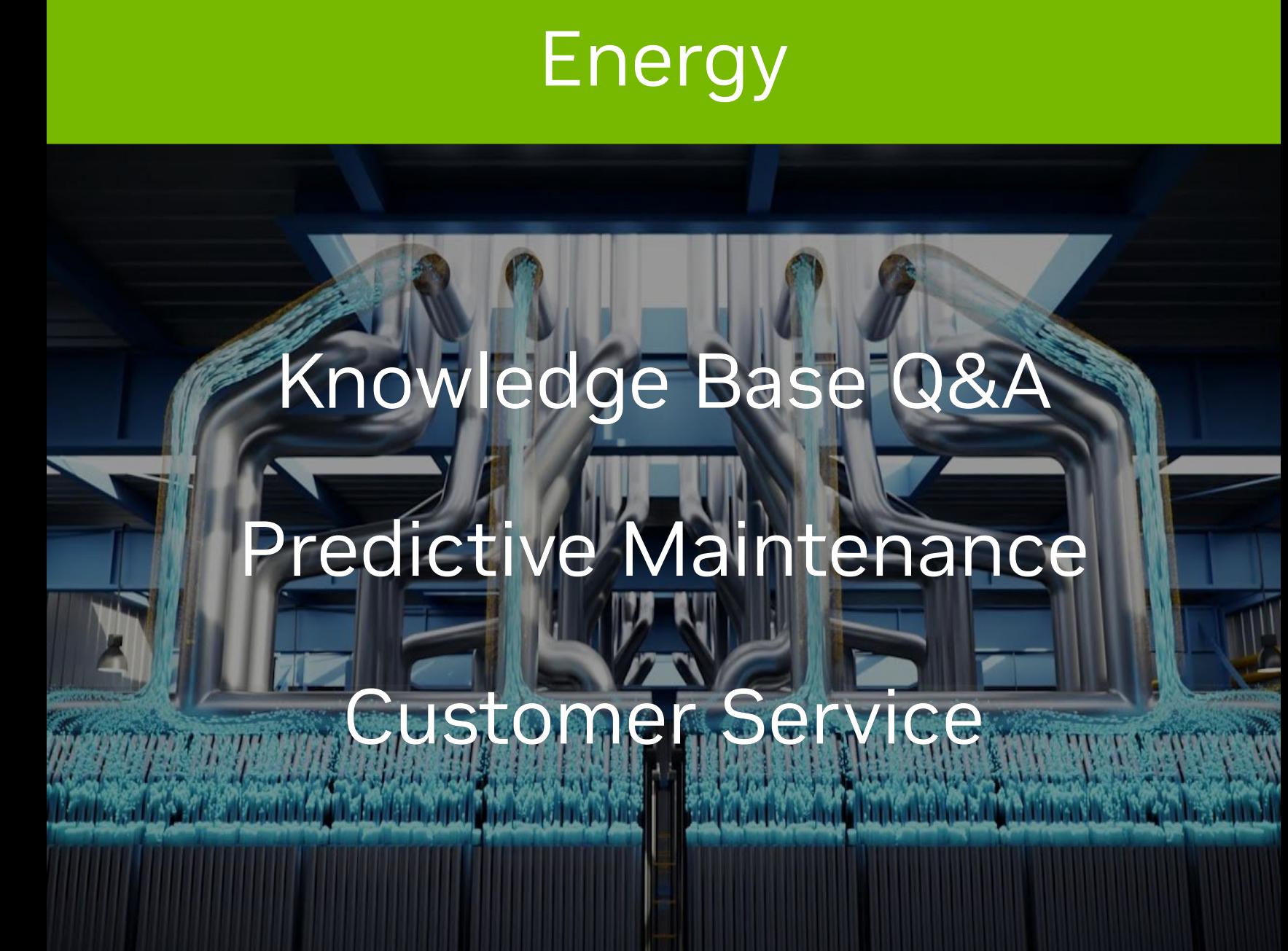
## Manufacturing



## Federal



## Energy



# Custom Generative AI for Enterprise IT

ServiceNow and NVIDIA have partnered to develop generative AI capabilities aimed at enhancing workflow automation across various business processes.

Leveraging NVIDIA's technology, ServiceNow is creating large language models trained on its specific data. This will enhance ServiceNow's existing AI functionality, enabling new applications of generative AI across the enterprise, including IT, customer service, and developers, to bolster workflow automation and boost productivity.

This innovative AI solution will provide higher accuracy and value in IT tasks, reshape customer service, and improve the employee experience.

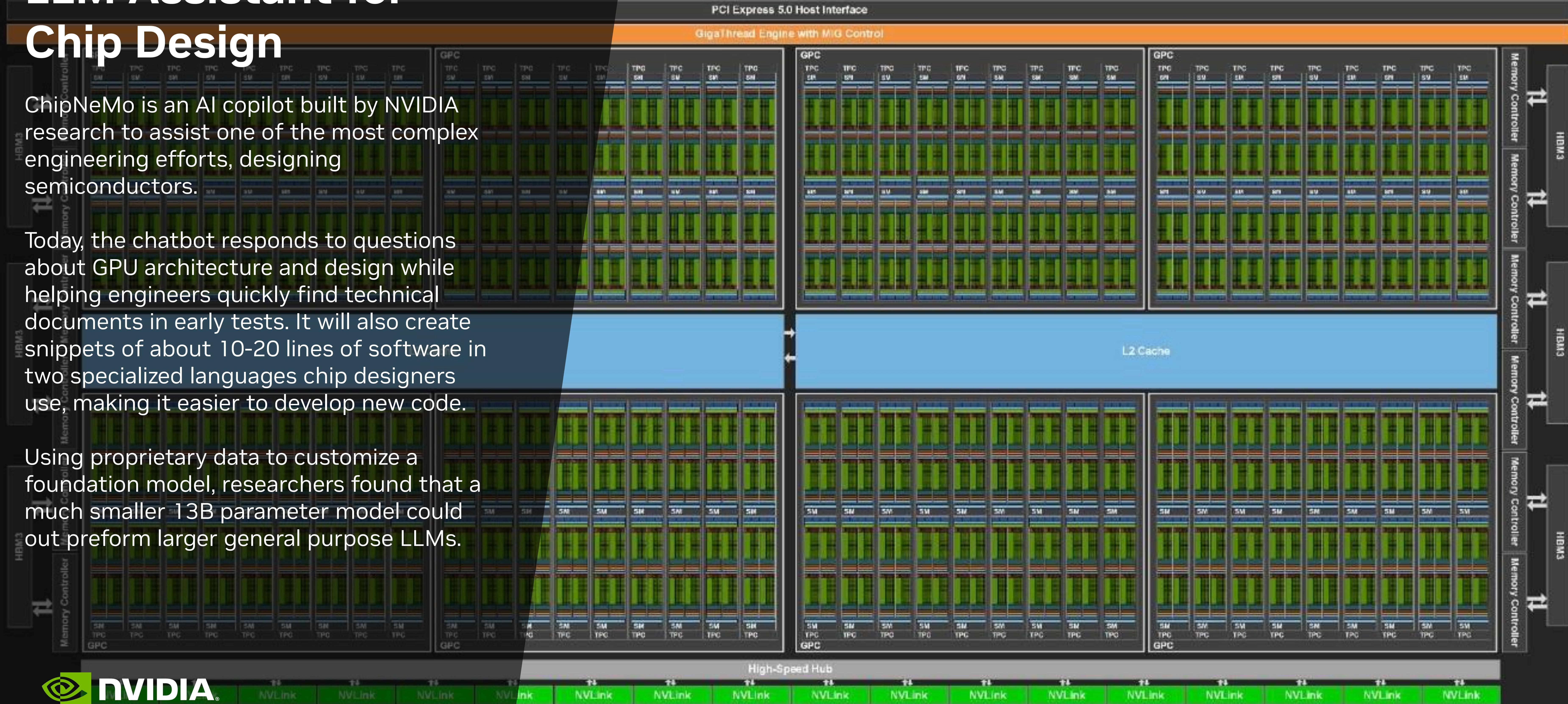


# LLM Assistant for Chip Design

ChipNeMo is an AI copilot built by NVIDIA research to assist one of the most complex engineering efforts, designing semiconductors.

Today, the chatbot responds to questions about GPU architecture and design while helping engineers quickly find technical documents in early tests. It will also create snippets of about 10-20 lines of software in two specialized languages chip designers use, making it easier to develop new code.

Using proprietary data to customize a foundation model, researchers found that a much smaller 13B parameter model could outperform larger general purpose LLMs.





Hippocratic AI  
— Do No Harm —



# Key Things to Remember About Generative AI

# Key Things to Remember About Generative AI

1

Generative AI is a tool, it's not  
the full solution

# Key Things to Remember About Generative AI

1

Generative AI is a tool, it's not  
the full solution

2

The space is moving quickly, build applications that  
have room for flexibility and adaptability

# Key Things to Remember About Generative AI

1

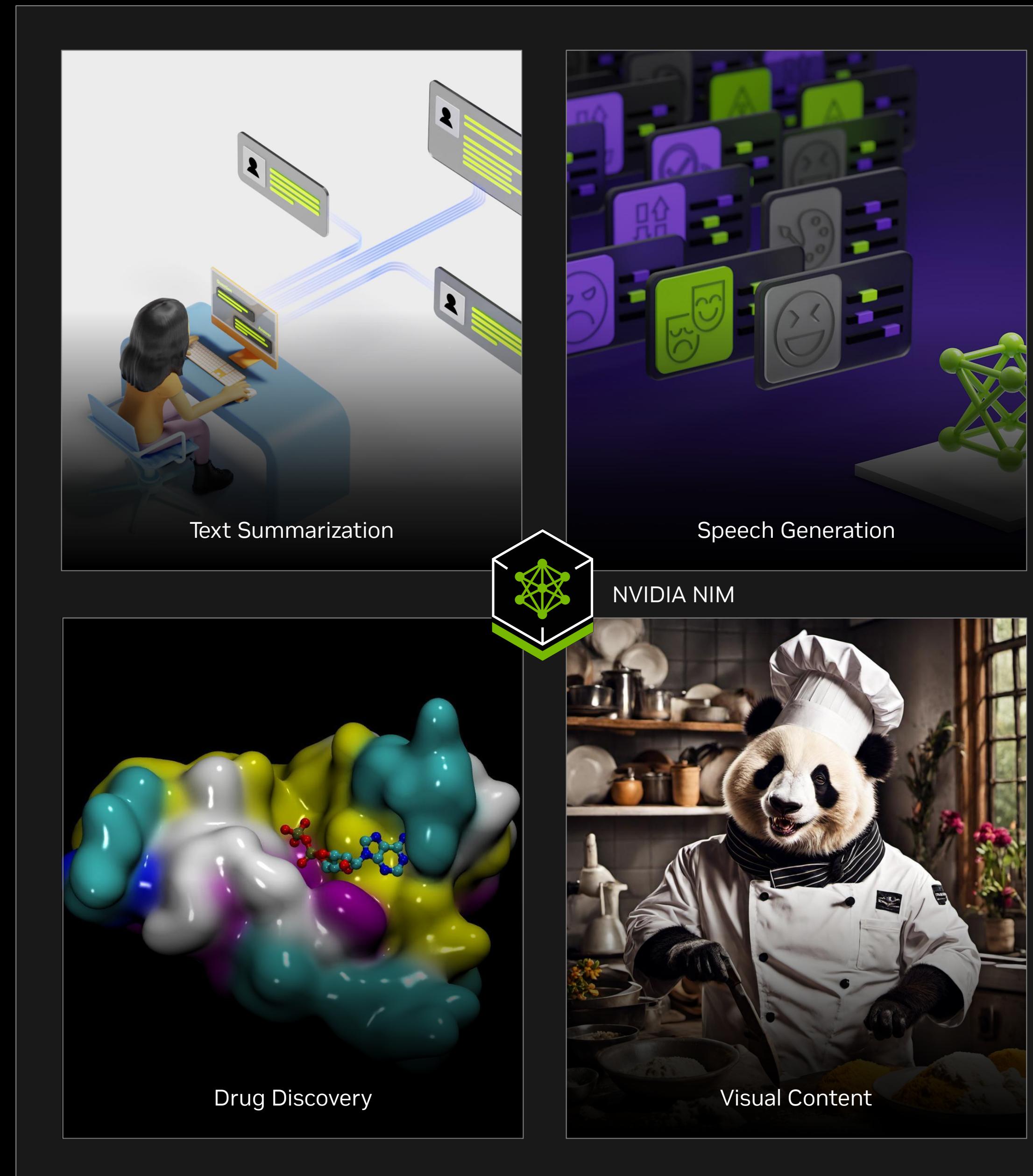
Generative AI is a tool, it's not  
the full solution

2

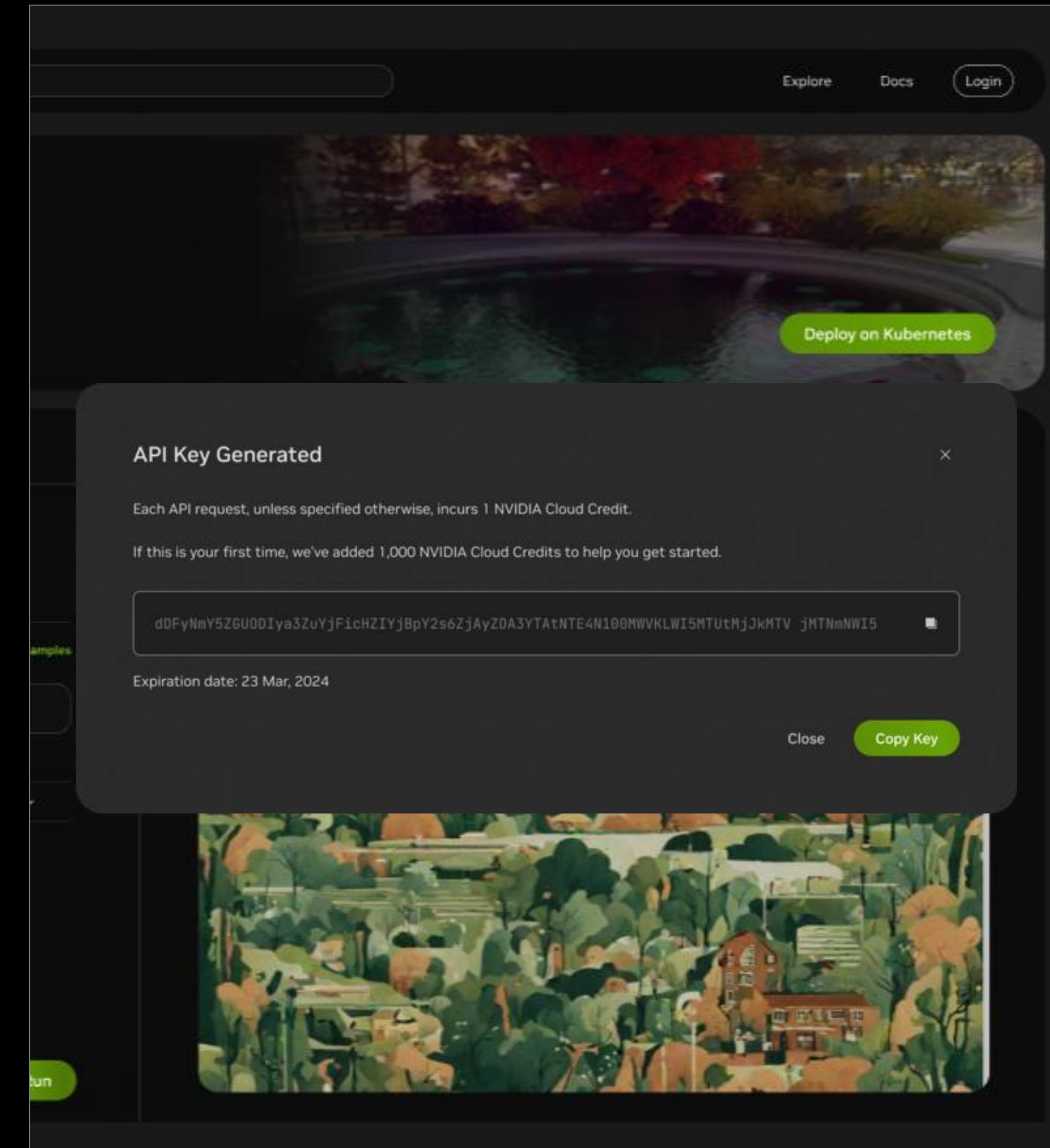
The space is moving quickly, build applications that  
have room for flexibility and adaptability

3

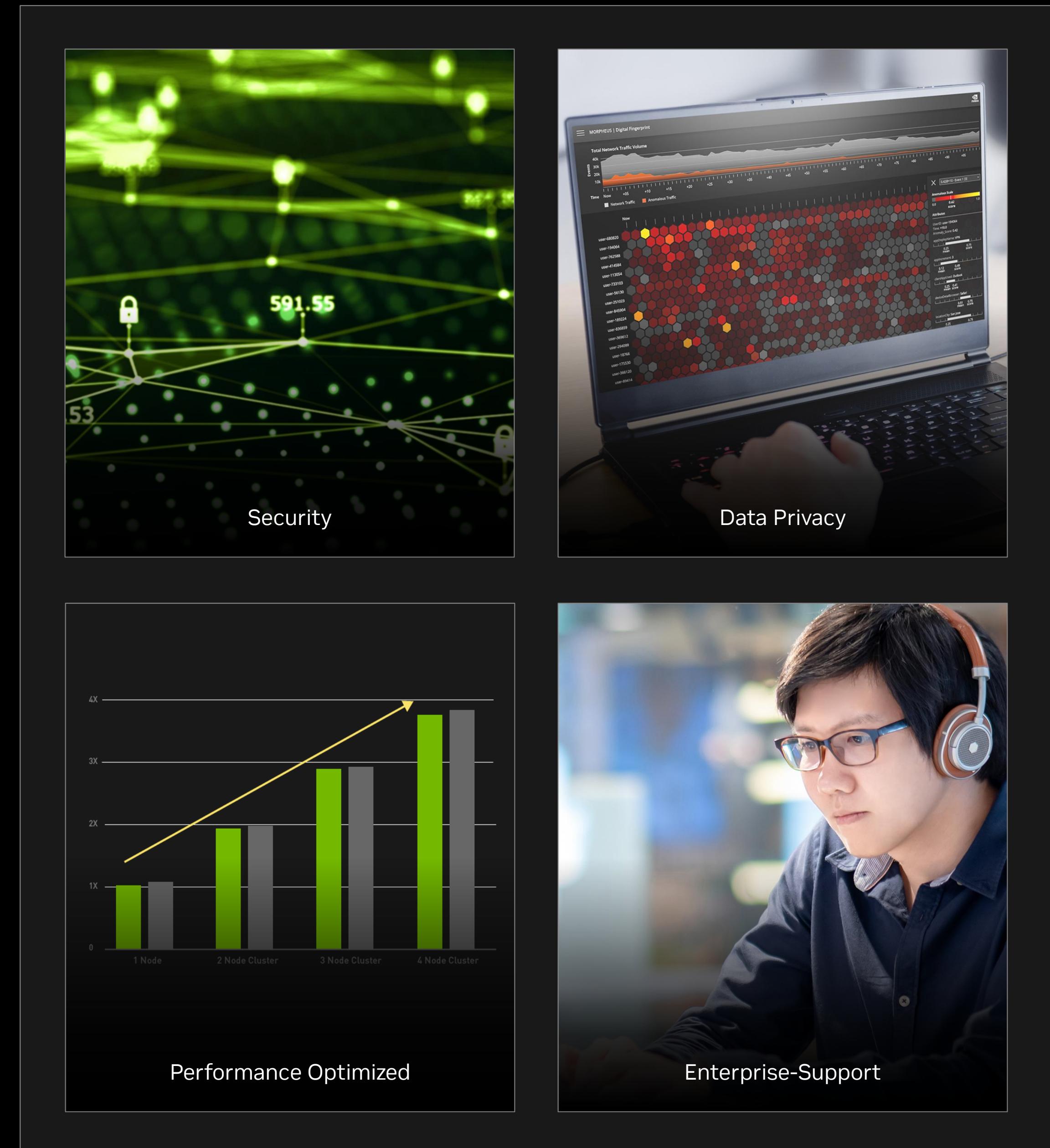
The future is bright, buy some shades



Experience Models



Prototype with APIs



Deploy with NIMs

