# Building LLM Applications with Prompt Engineering

Introduction

# Building LLM Applications with Prompt Engineering

Introduction

**Course Environment**

✏ EDIT    ⚙    ⧉    ⤓    🗑    ⠿

NVIDIA | DEEP LEARNING INSTITUTE

▶

START

# Large Language Models

Artificial Intelligence system designed to understand, generate, and process human language

## 2022

## 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.

## 2023

## Experimentation

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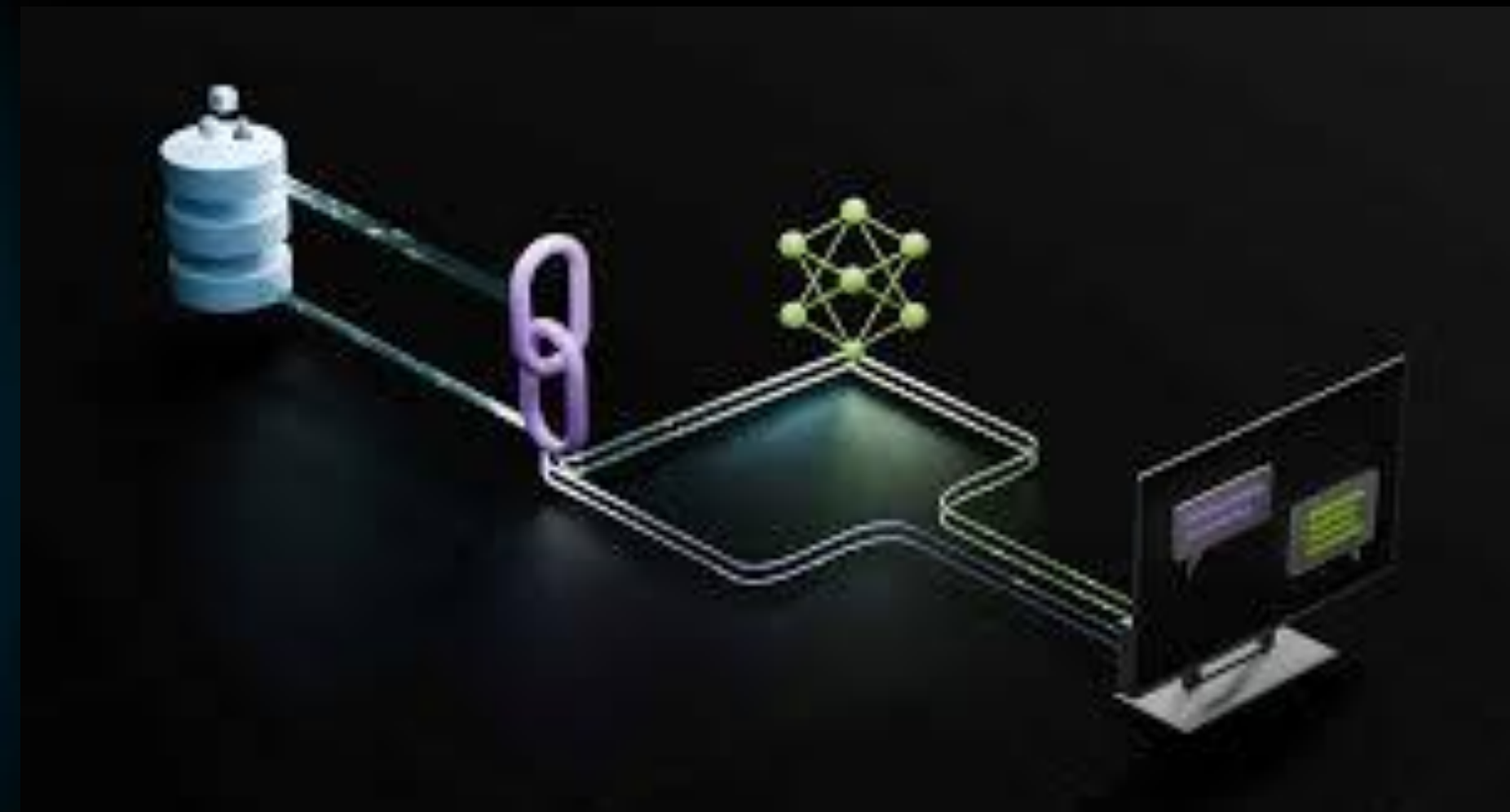
## TODAY

## Production

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

# 3 Methods for Interacting with LLMs

**Prompt Engineering**

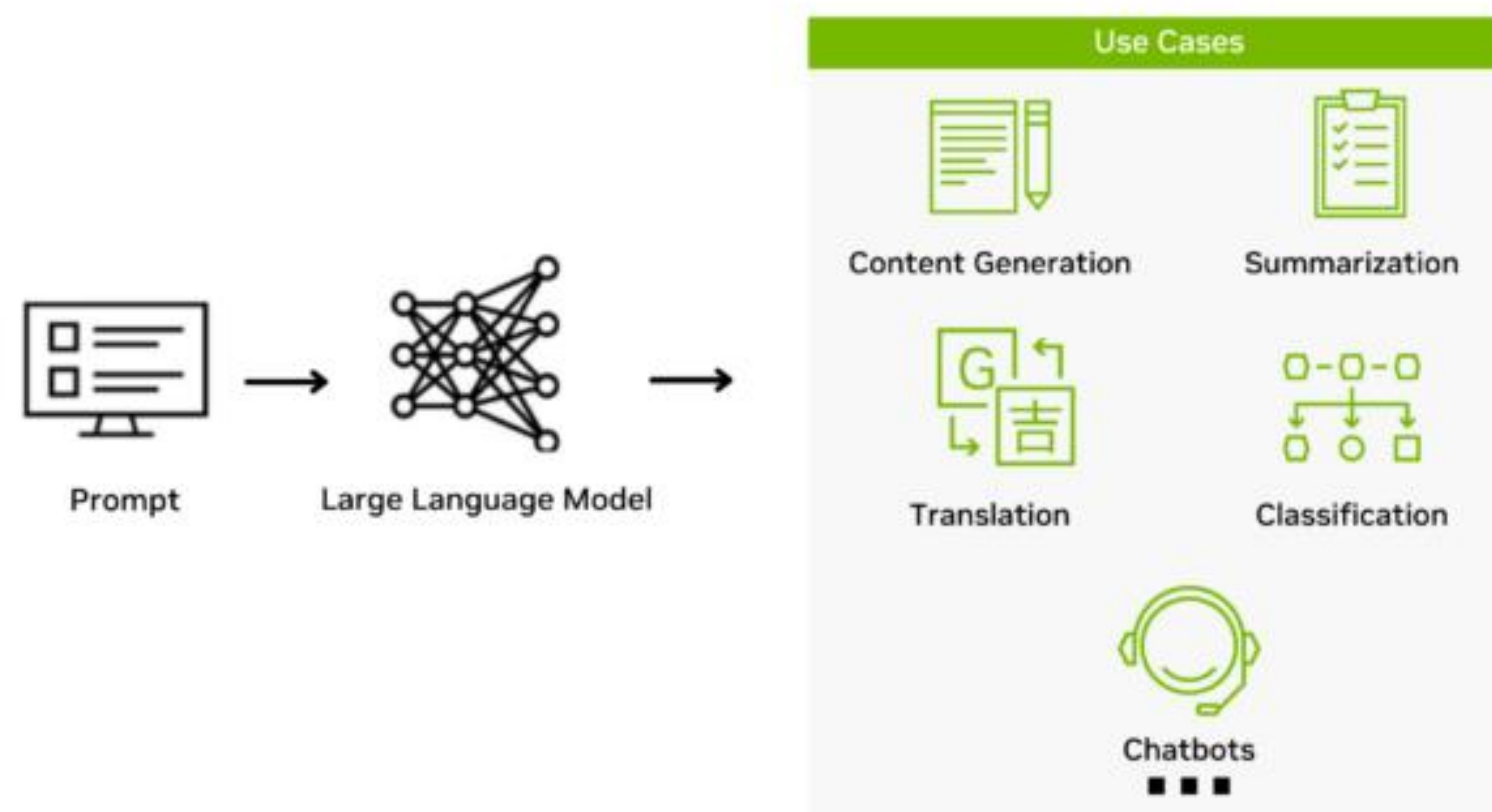**Retrieval Augmented Generation**

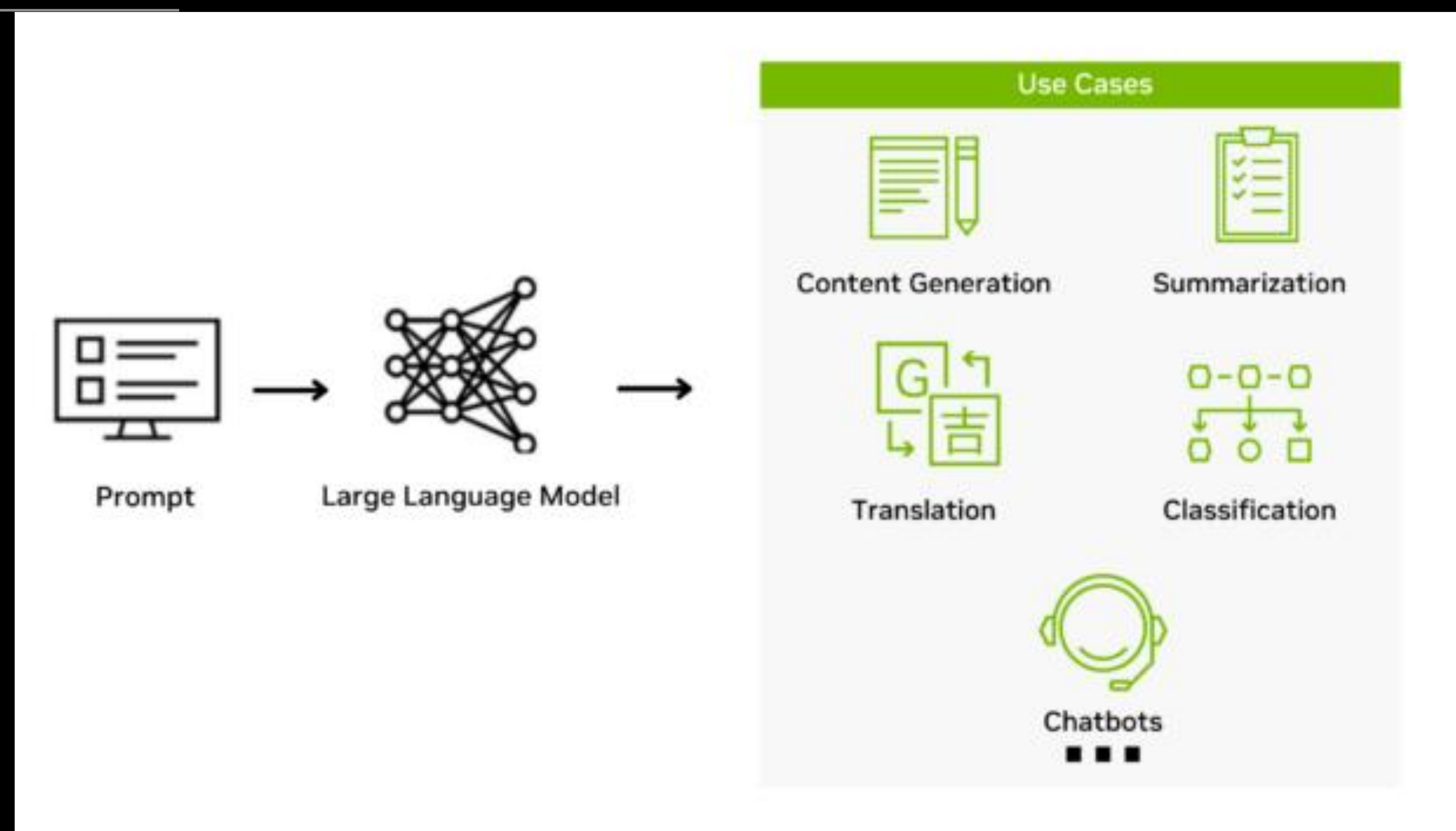**(Parameter Efficient) Fine-tuning**

# Goal of this Course

- Learn to interact programmatically with Chat-variant LLMs

- Be capable of using LLMs for a wide variety of application use cases

- Become fluent with fundamental LangChain techniques

- Develop good habits around iterative prompt engineering



NVIDIA.

# Goal of this Course

- Introduction to Large Language Models (slides)

- NVIDIA NIM (slides)

- Intro to Prompting (interactive)

- LCEL Chains (interactive)

- PE techniques w/ messages (interactive)

- Chatbots (interactive)

- Structured data and document tagging (interactive)

- Tools and Agents (interactive)

# Course Content

- Introduction to Large Language Models (slides)

- NVIDIA NIM (slides)

- Intro to Prompting (interactive)

- LCEL Chains (interactive)

- PE techniques w/ messages (interactive)

- Chatbots (interactive)

- Structured data and document tagging (interactive)

- Tools and Agents (interactive)

# Prerequisites
## Building LLM Applications with Prompt Engineering

- Basic Python Proficiency

- Prior LLM Exposure
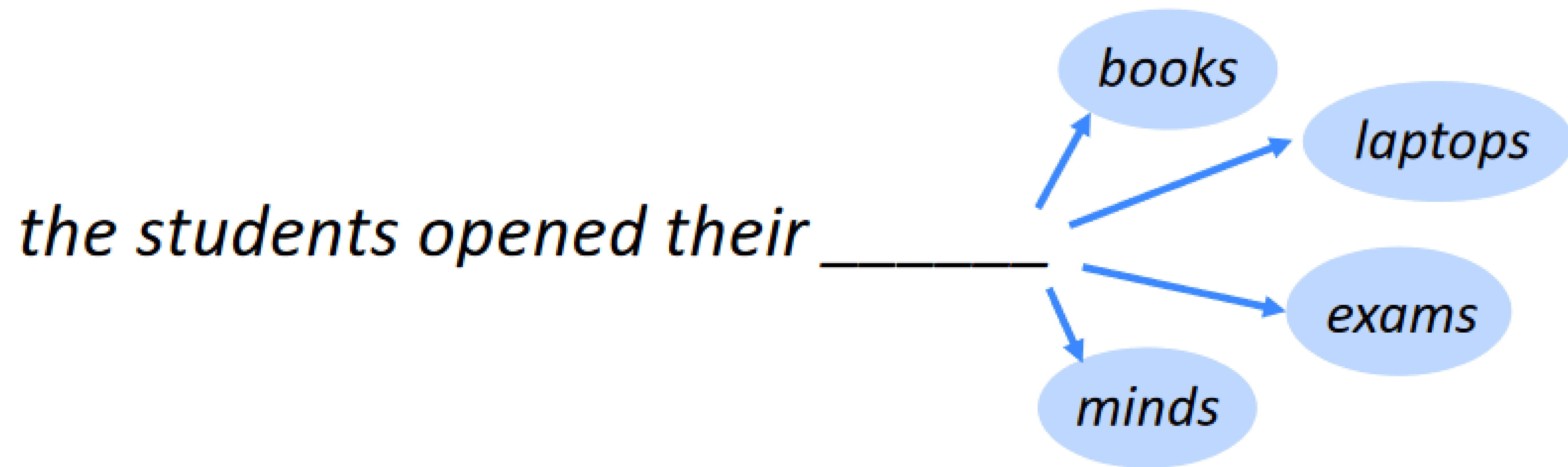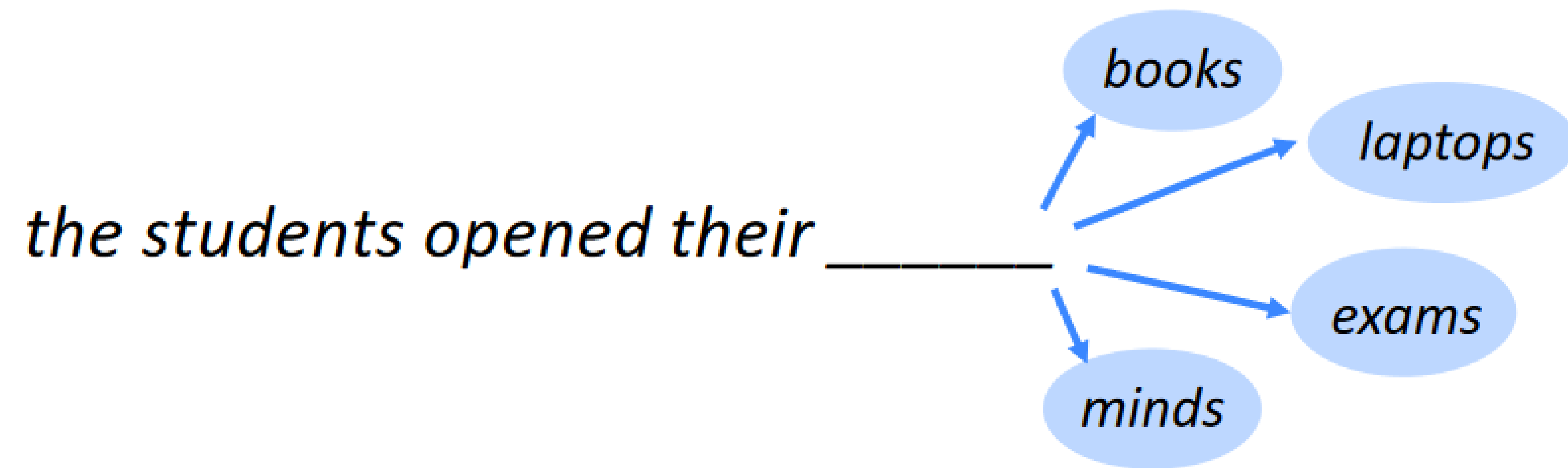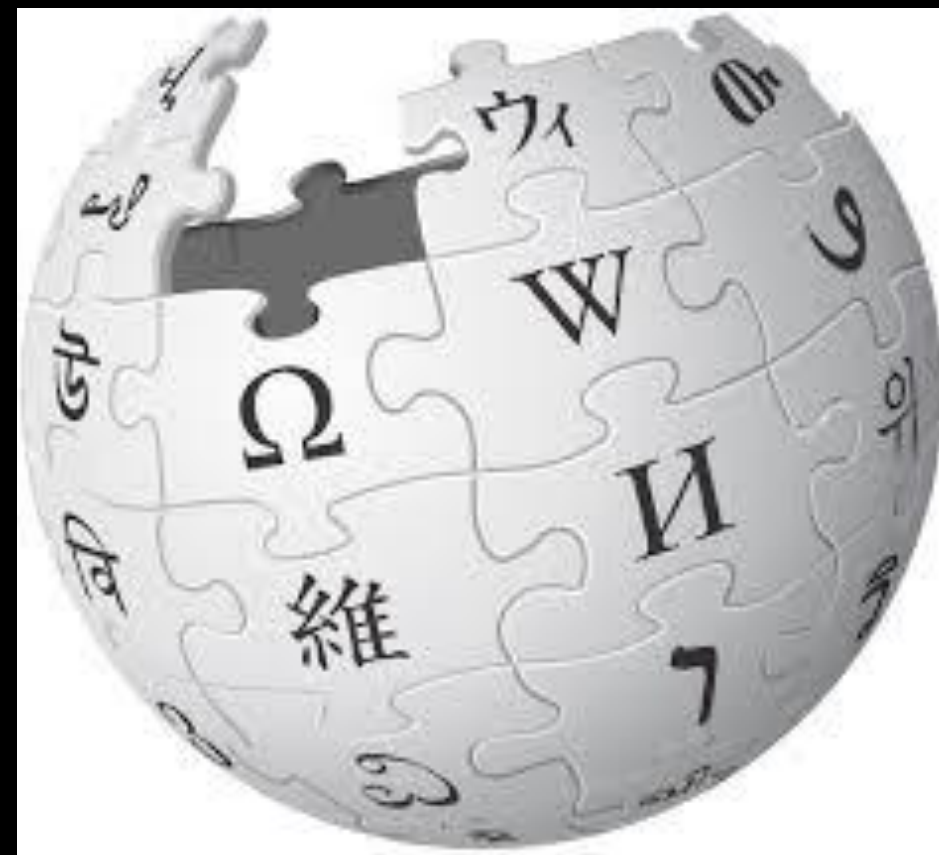
- Plus: Experience with LangChain or Knowledge in Natural Language Processing

# Building LLM Applications with Prompt Engineering
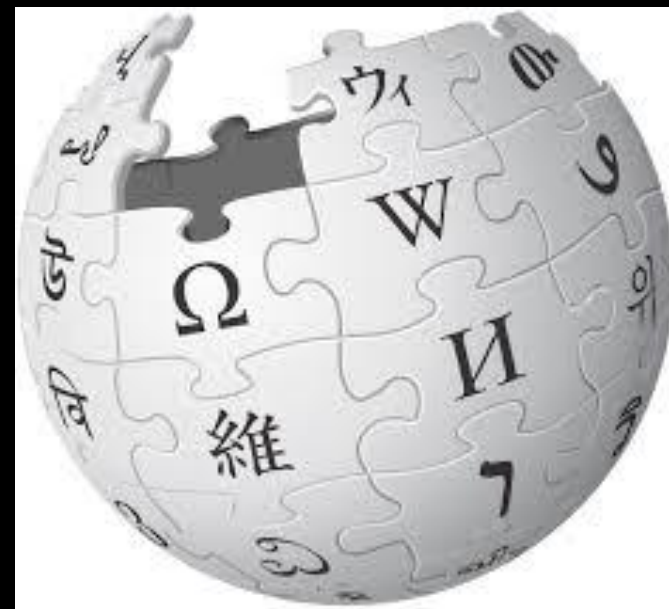
Introduction to Large Language Models

# What is Language Modeling?



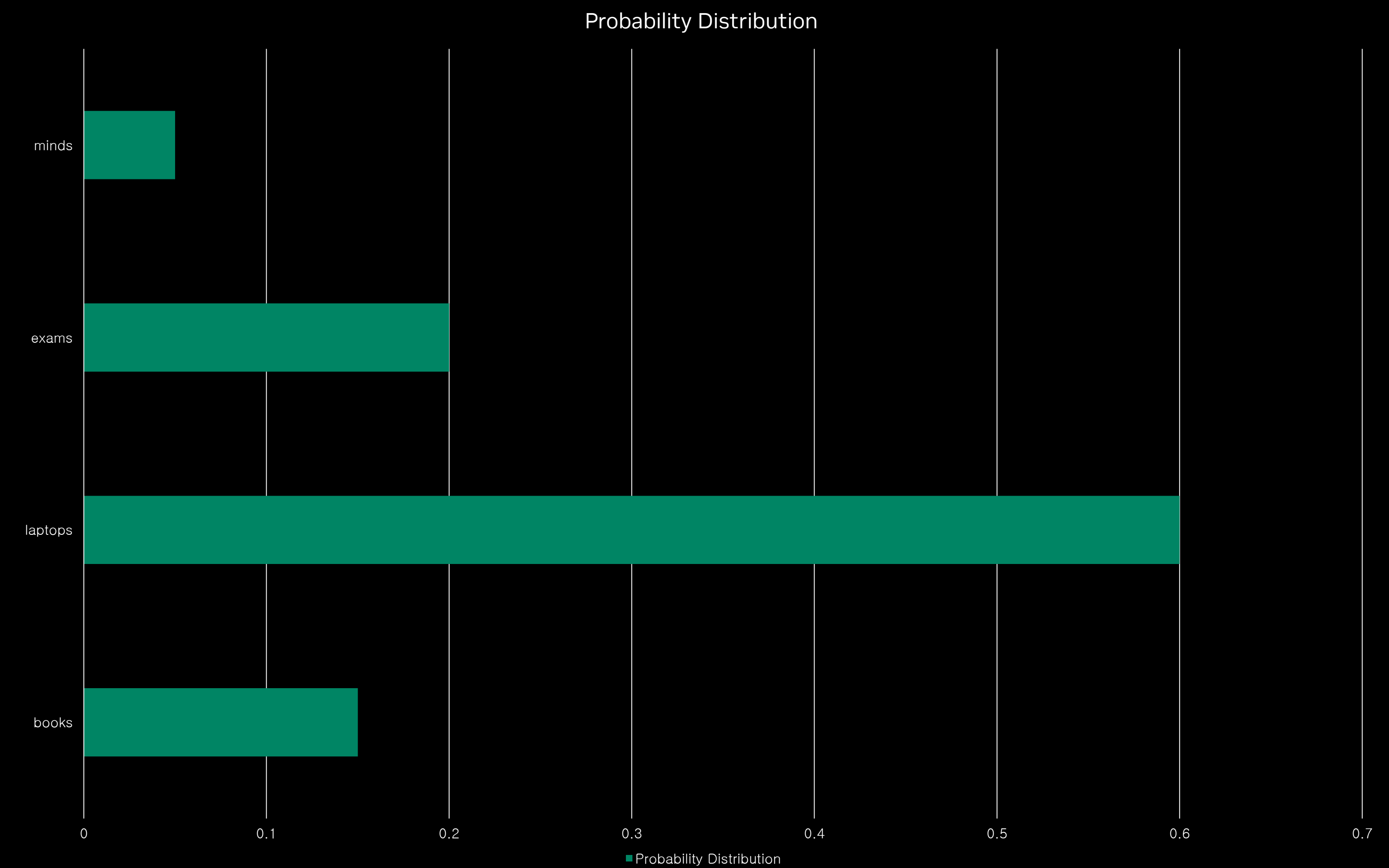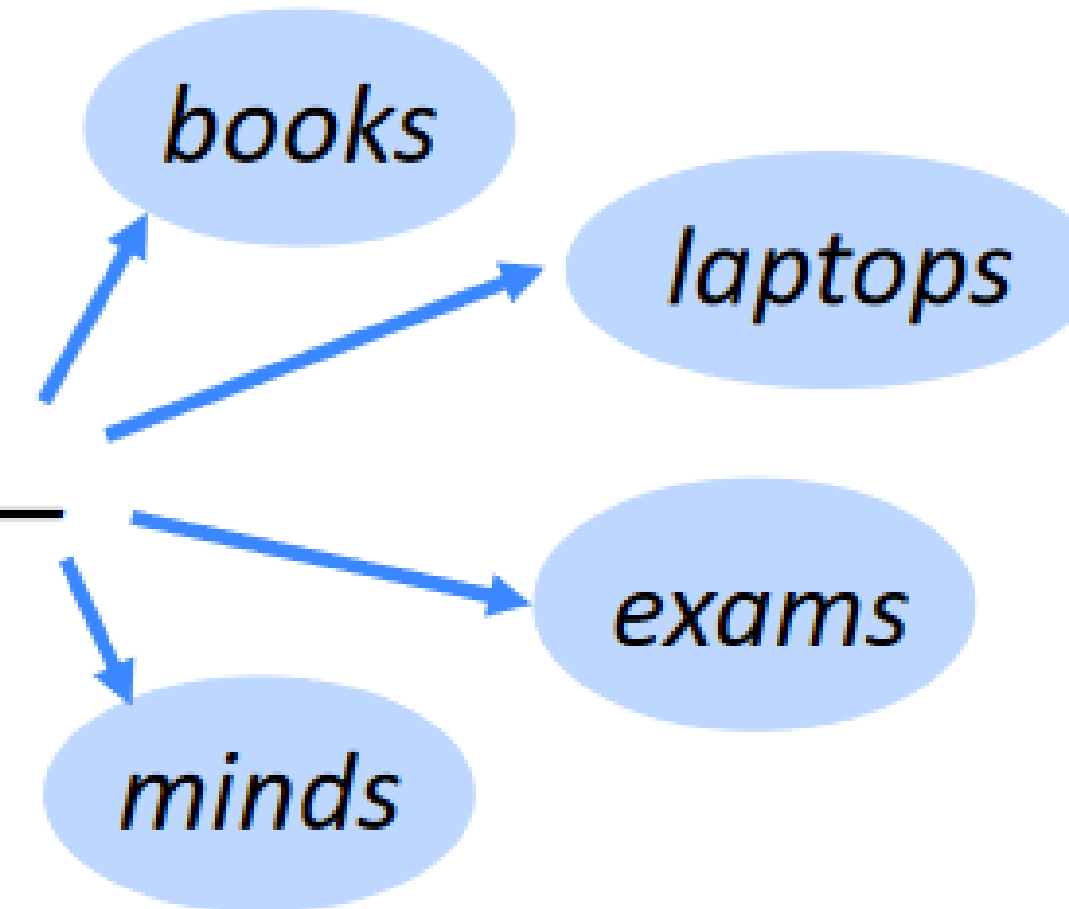the students opened their _____ → books, laptops, exams, minds

# Large Language Models Pipeline

# Large Language Models

# Training LLMs: Tokenization

# Training LLMs: Embedding

# Training LLMs: Embedding



$E(\cdots) + E(pos) + E(word)$

Per-Token Feed-Forward
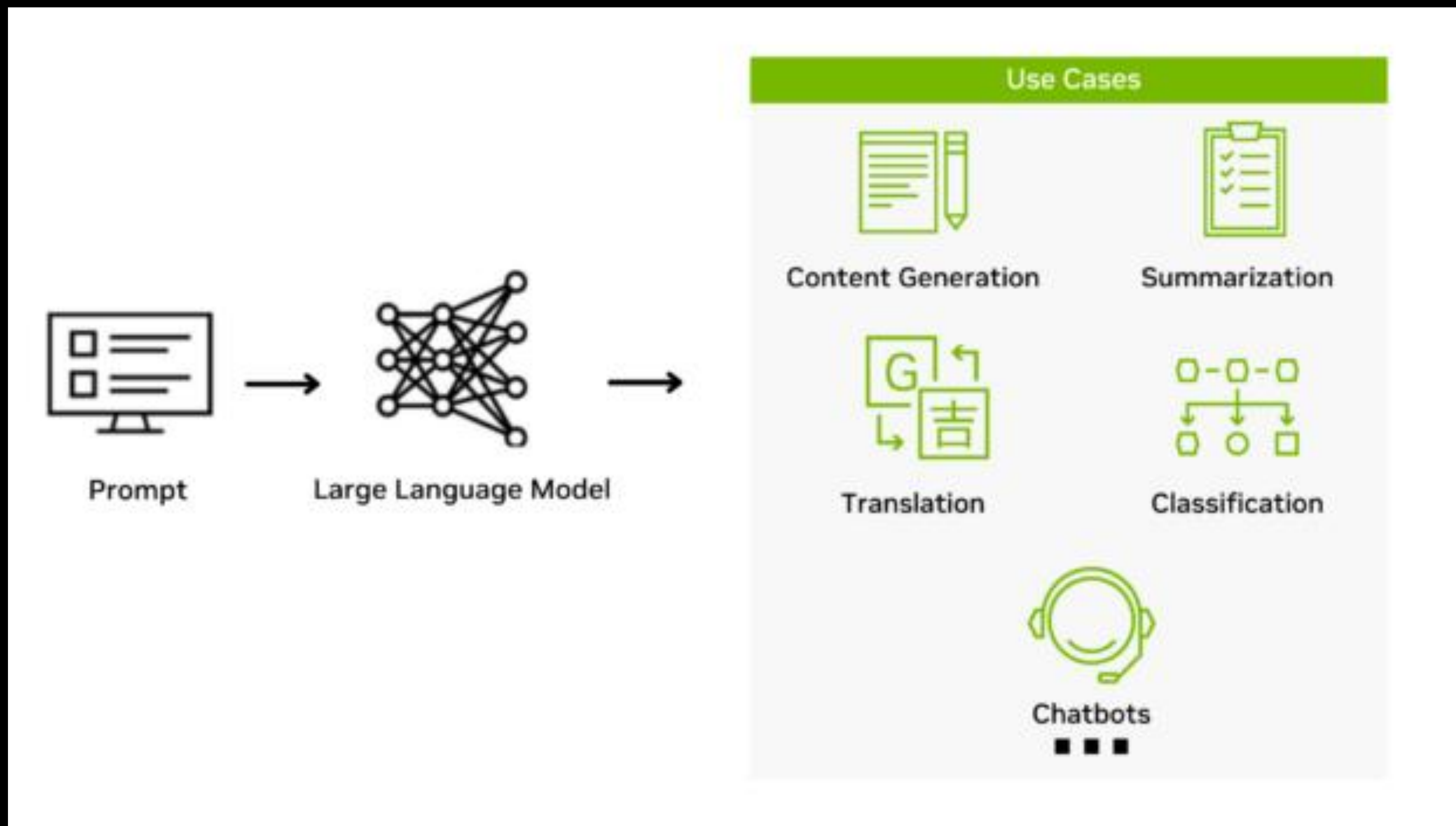
Per-Token Predictions

# Large Language Models

**Next-token prediction task**
•Massively Multi-task learner capable of doing all sorts of tasks

**Hard to control the text generation**

| Task | Example sentence in pre-training that would teach that task |
|------|-----------------------------------------------------------|
| Grammar | In my free time, I like to {code, apples} |
| Lexical Semantic | I went to the store to buy papaya, dragon fruit, and {pineapples, croissant} |
| World Knowledge | The capital of South Korea is {Seoul, Paris} |
| Sentiment Analysis | Movie review: I was engaged and on the edge of my seat the whole time. The movie was {bad, good} |
| Math Question | Arithmetic exam answer key: 3 + 8 + 4 = {15, 384} |
| Translation | The word for pretty in Spanish is {bonita, hola} |

# Prompt Engineering



Prompt Engineering is asking the right question to the LLM depending on your use case.
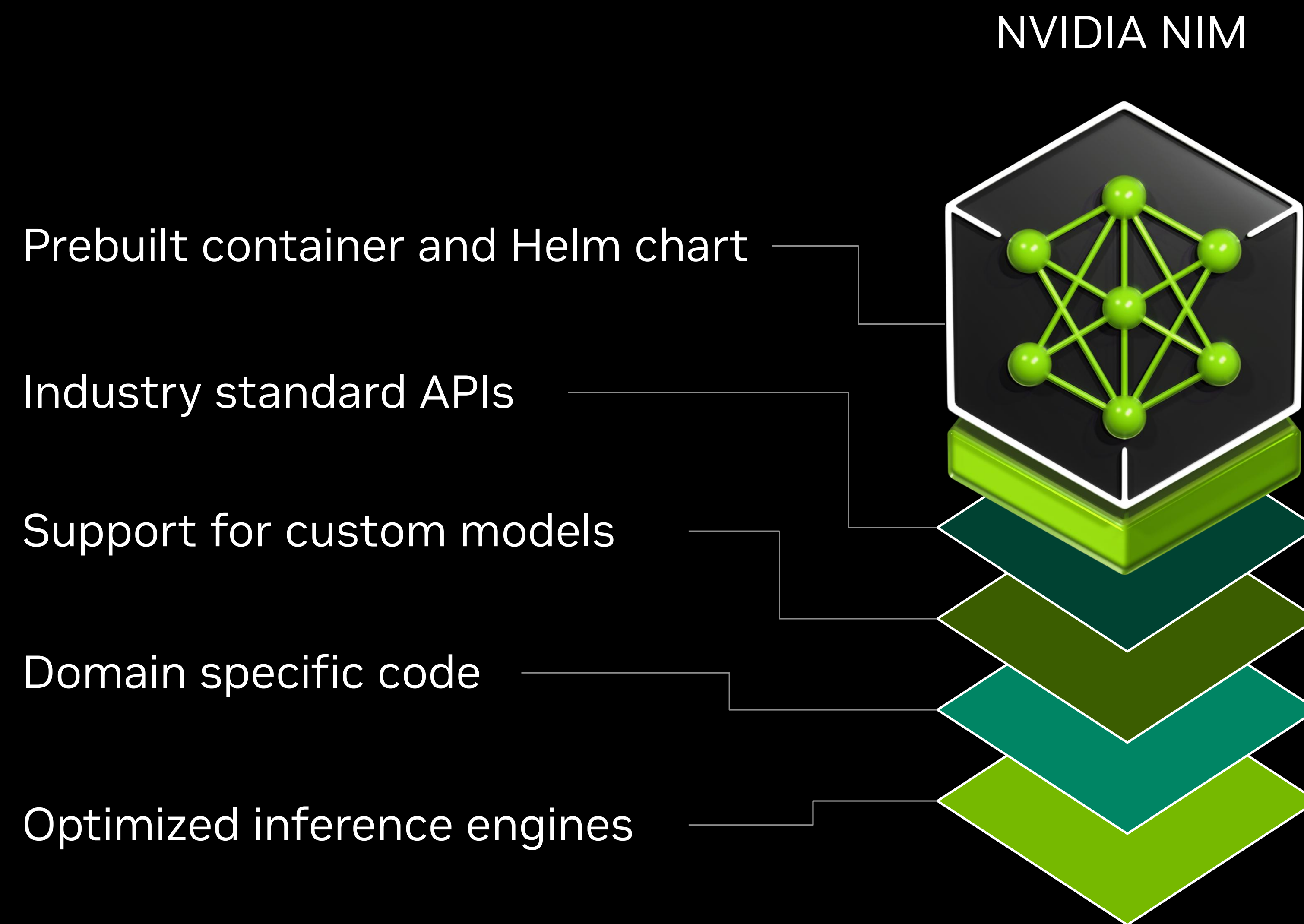
# NVIDIA NIM Optimized Inference Microservices

NVIDIA NIM is the fastest way to deploy AI models on accelerated infrastructure across cloud, data center, and PC

NVIDIA NIM



Prebuilt container and Helm chart
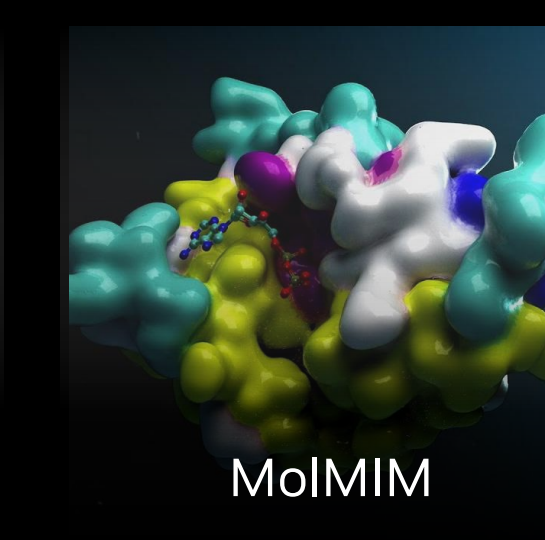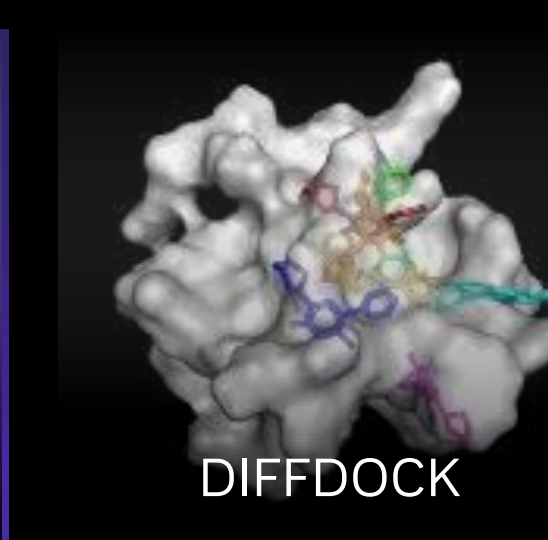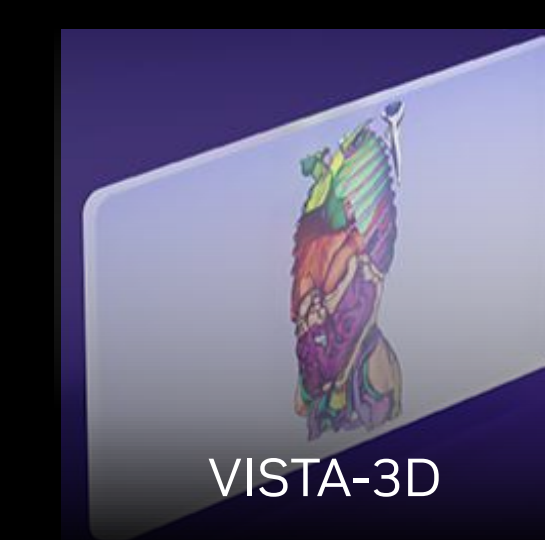
Industry standard APIs

Support for custom models

Domain specific code

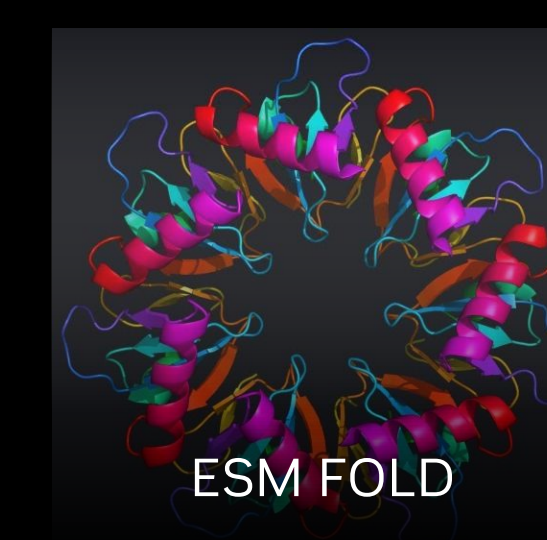Optimized inference engines

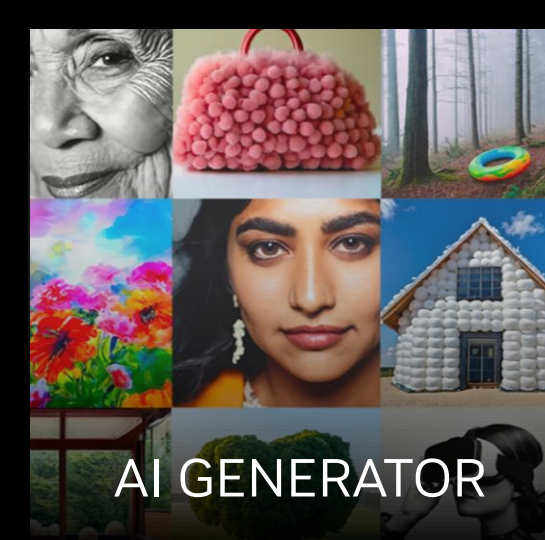**Deploy anywhere with security and control**

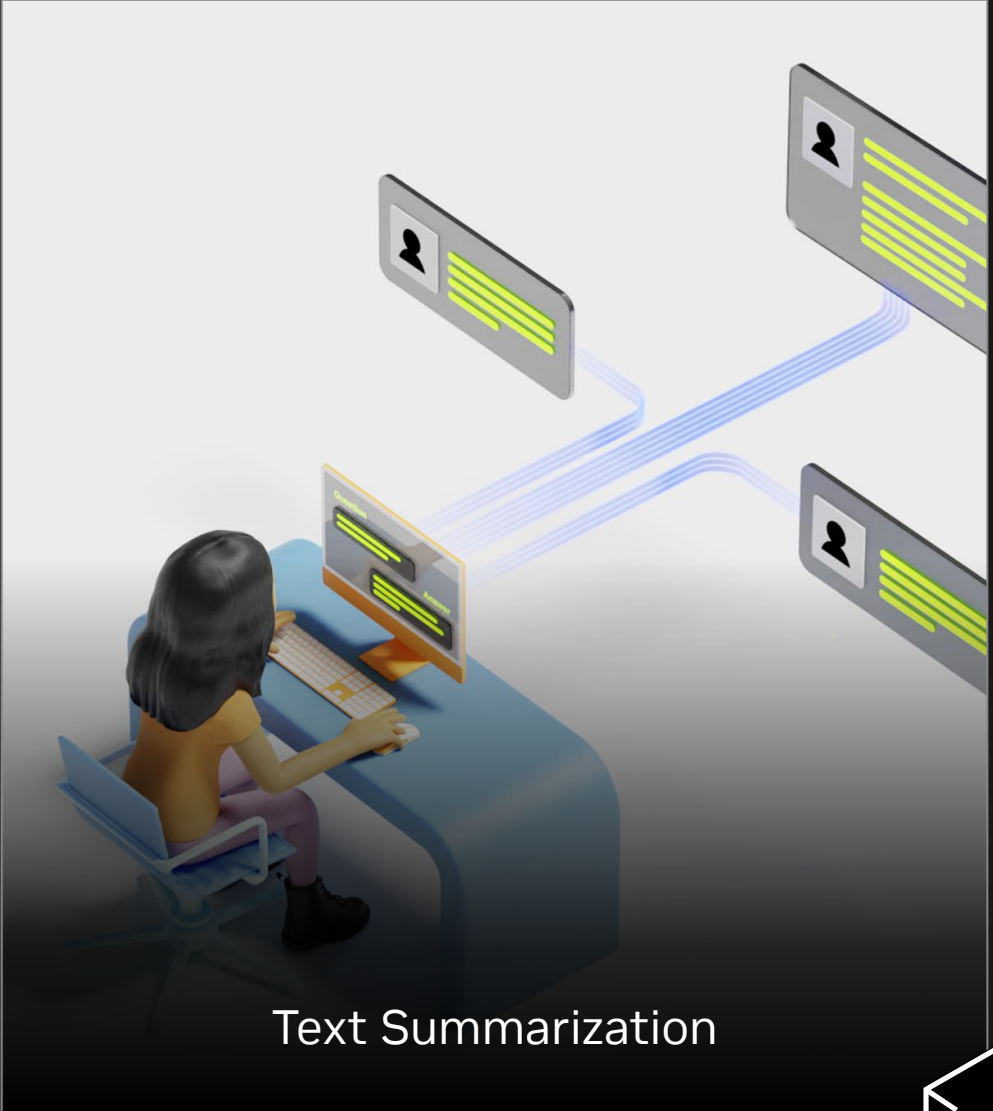**Speed time to market**

**Empower developers**

**Optimize throughput**

**Boost accuracy**

**Deploy in production**

LLAMA 3.1   MIXTRAL 8x7B   GEMMA 7B   FUYU   NEMO RETRIEVER   AI GENERATOR   3D GENERATOR   AUDIO2FACE   ESM FOLD   VISTA-3D   DIFFDOCK   MolMIM

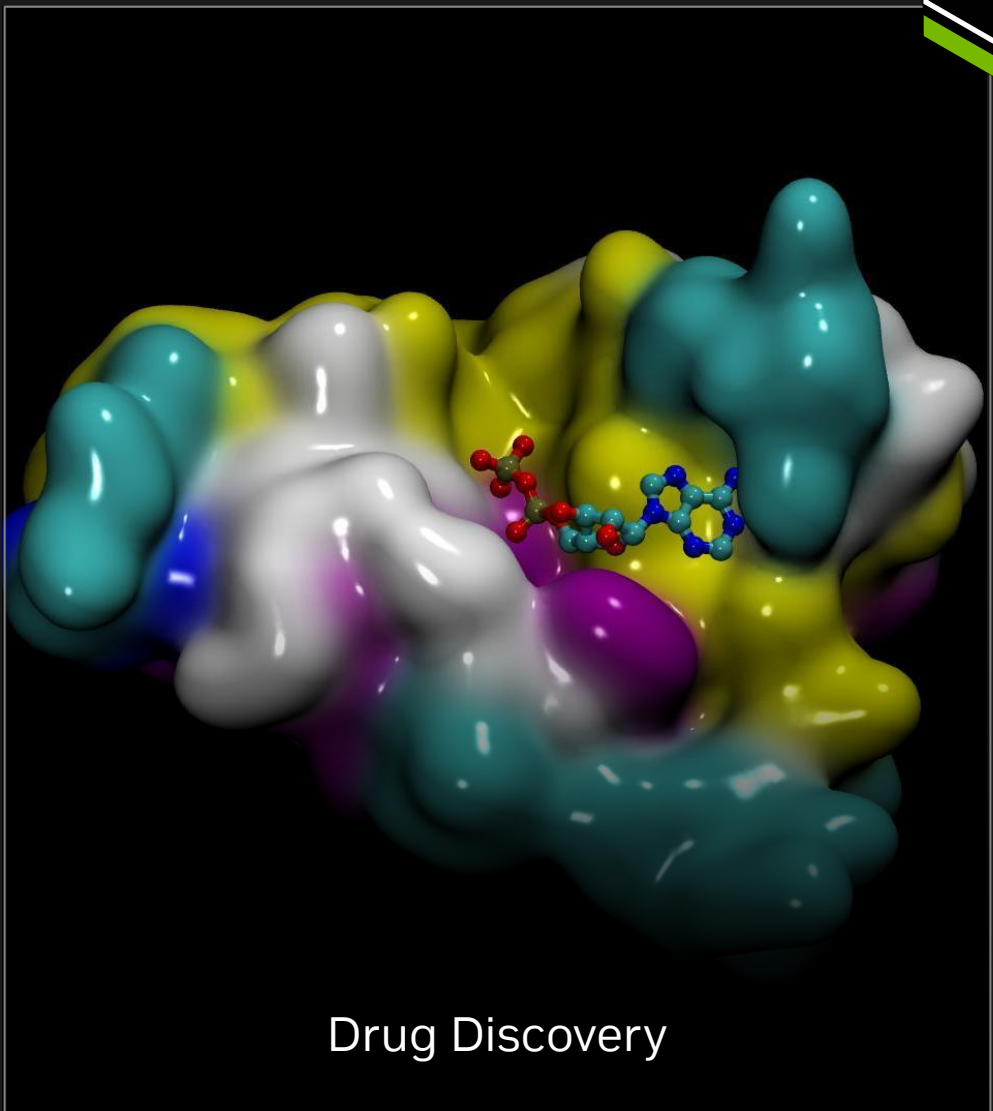# Experience and Run Enterprise Generative AI Models Anywhere

## Seamlessly integrate AI in business applications with NVIDIA AI APIs



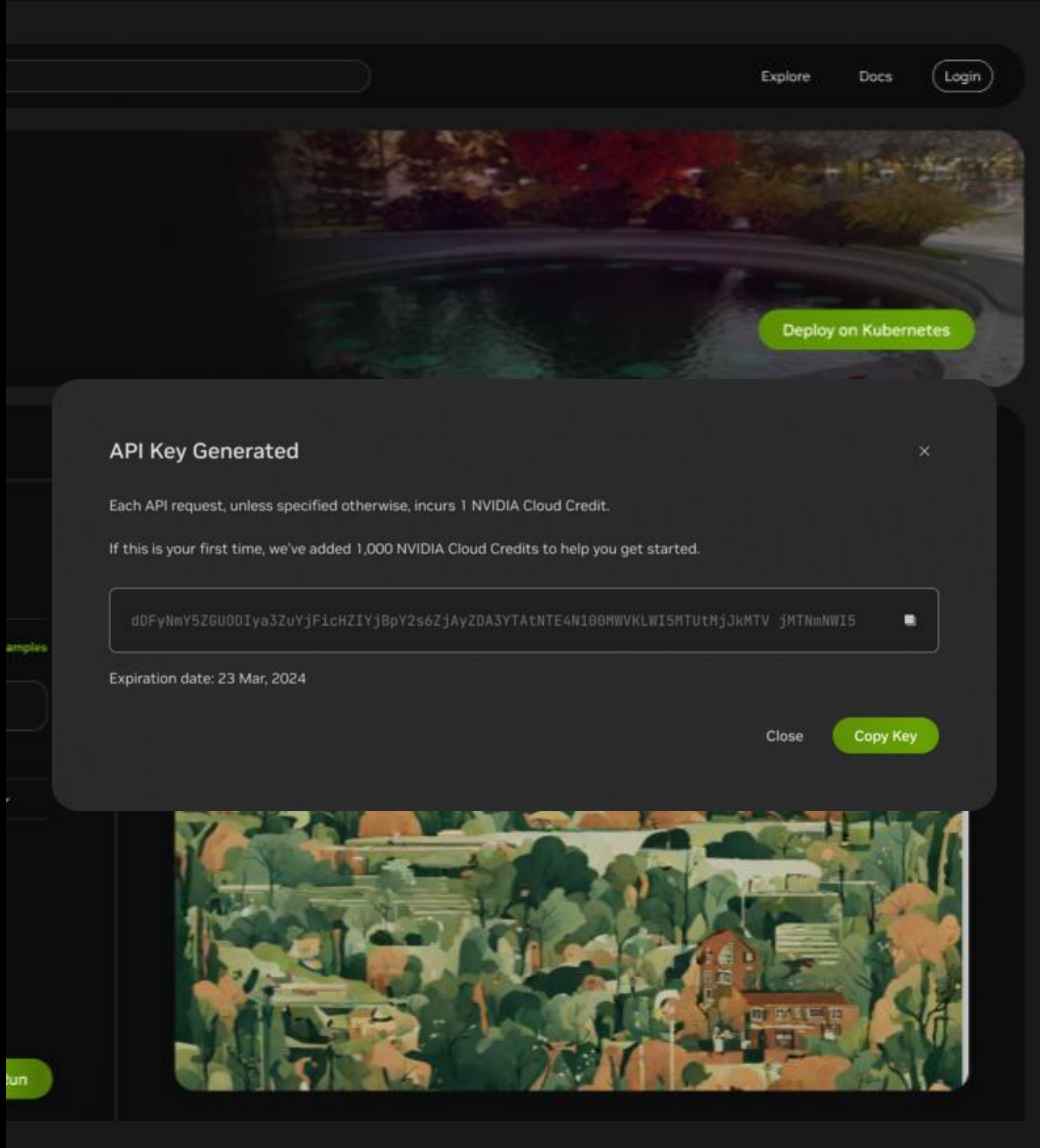Text Summarization

Speech Generation

NVIDIA NIM
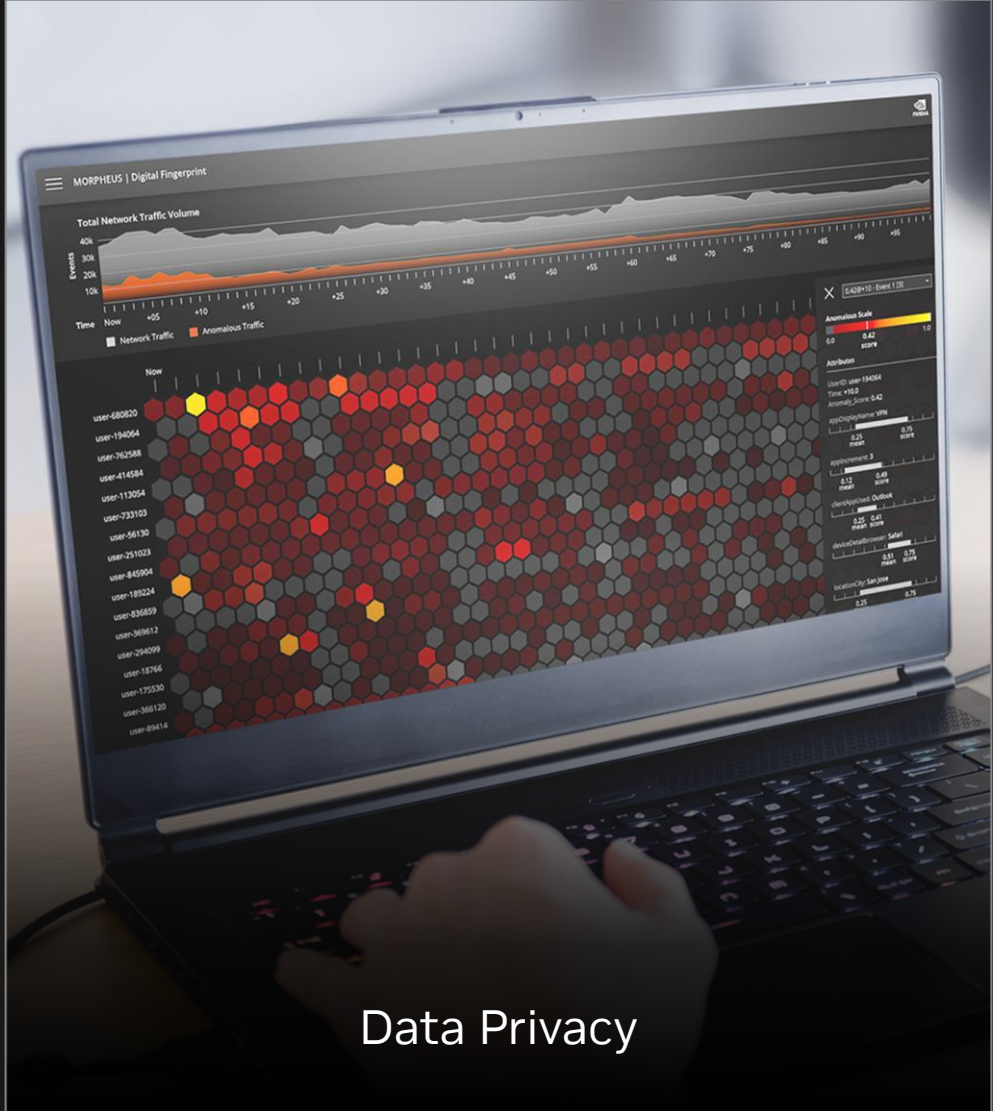
Drug Discovery

Visual Content

Explore    Docs    Login

Deploy on Kubernetes

**API Key Generated**                                                    ✕

Each API request, unless specified otherwise, incurs 1 NVIDIA Cloud Credit.

If this is your first time, we've added 1,000 NVIDIA Cloud Credits to help you get started.

dDFyNmY5ZGUODIya3ZuYjFicHZIYjBpY2s6ZjAyZDA3YTAtNTE4N10GMWVKLWI5MTUtMjJkMTV jMTNmNWI5

Expiration date: 23 Mar, 2024

Close    Copy Key

Security

Data Privacy

591.55

53
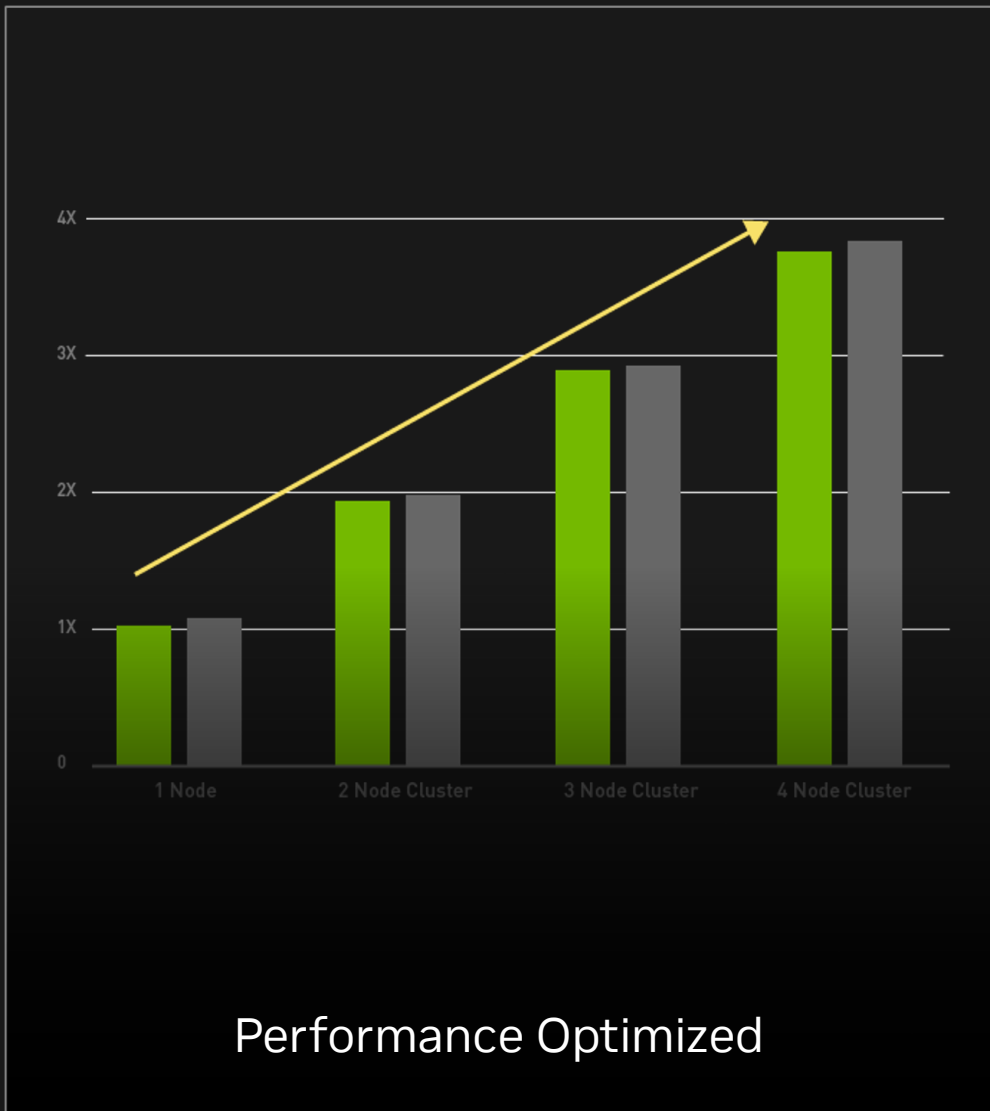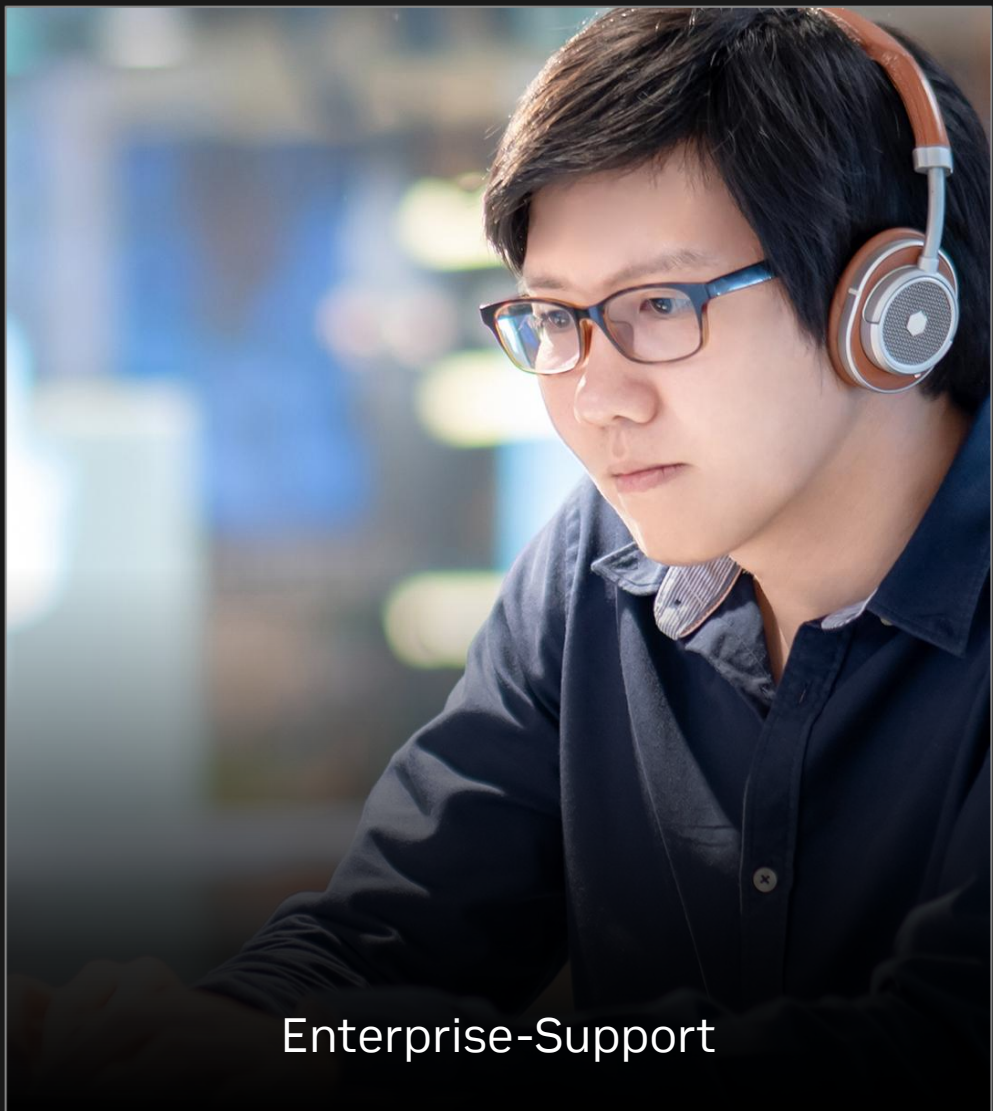
Performance Optimized

Enterprise-Support

- ## Experience Models
- ## Prototype with APIs
- ## Deploy with NIMs

# NIM Deployment Cycle



Locally-hosted LLMs reduce latency as opposed to models hosted on a public server

Scalable Deployment

Ownership of your intellectual property

# Building LLM Applications with Prompt Engineering