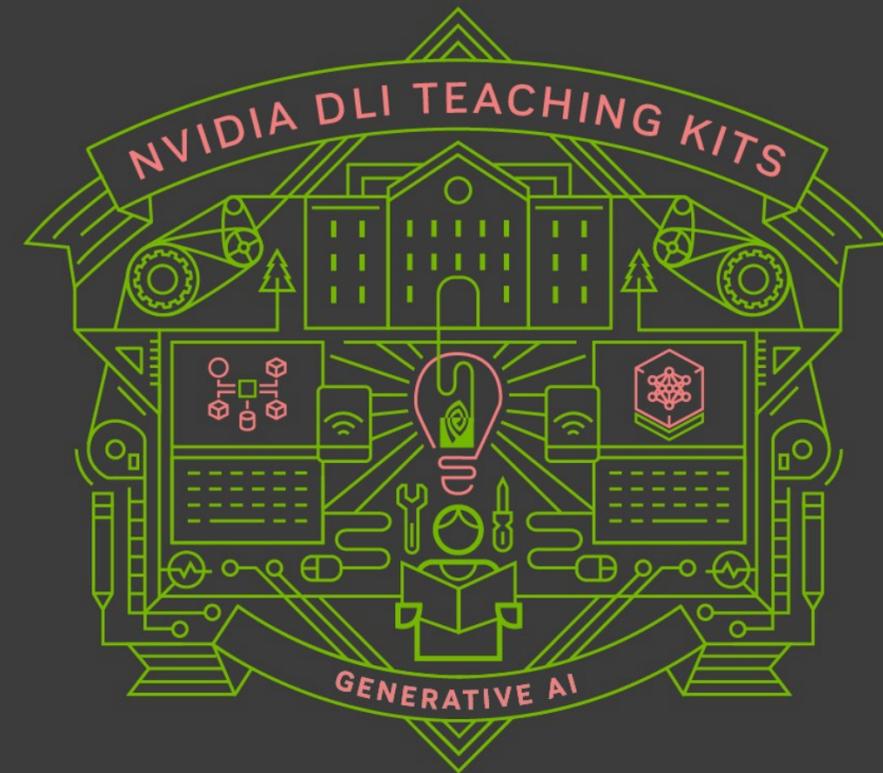# Lecture 1.1 – Introduction to the Generative AI Course

Generative AI Teaching Kit

The NVIDIA Deep Learning Institute Generative AI Teaching Kit is licensed by NVIDIA and Dartmouth College under the
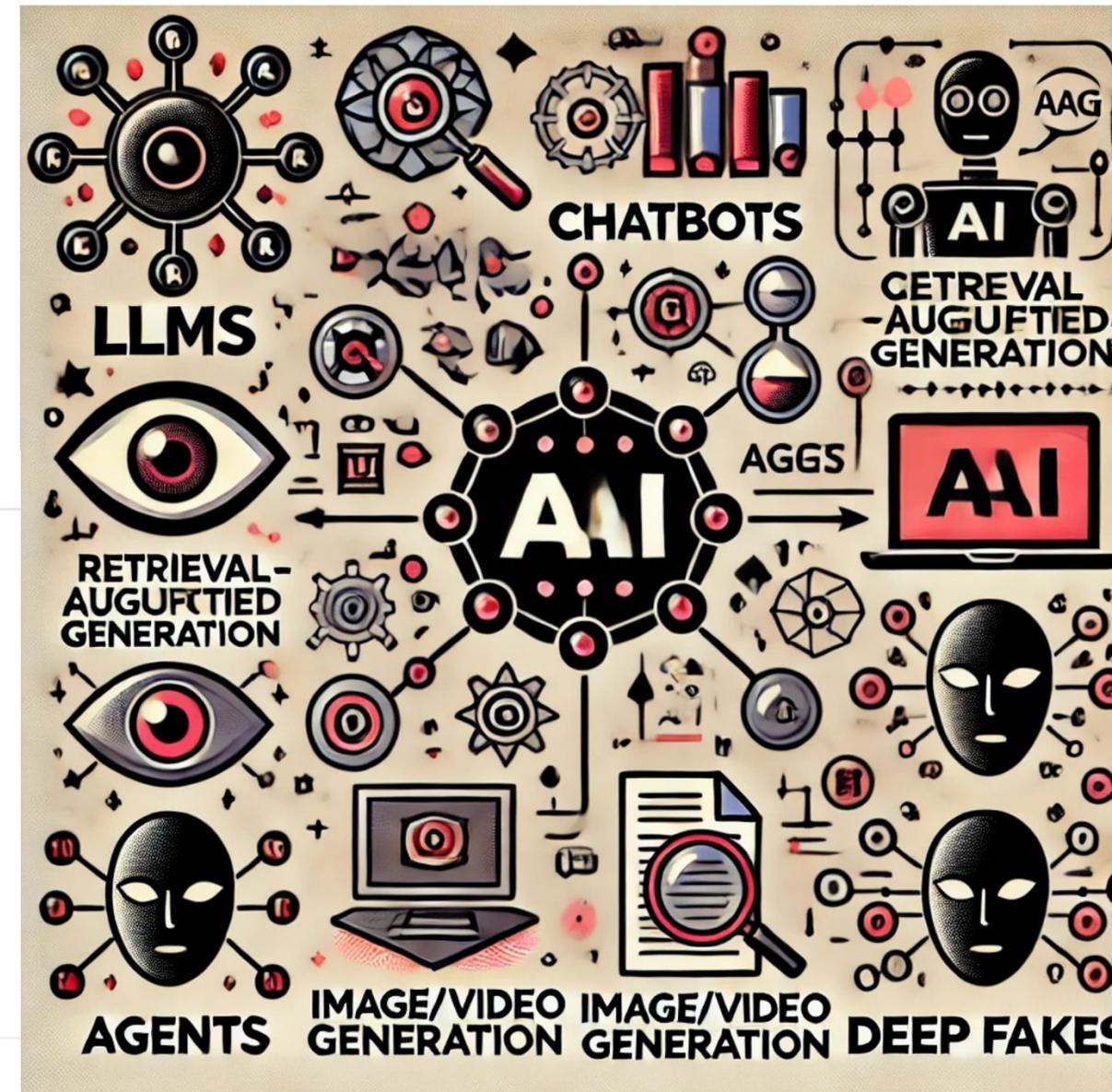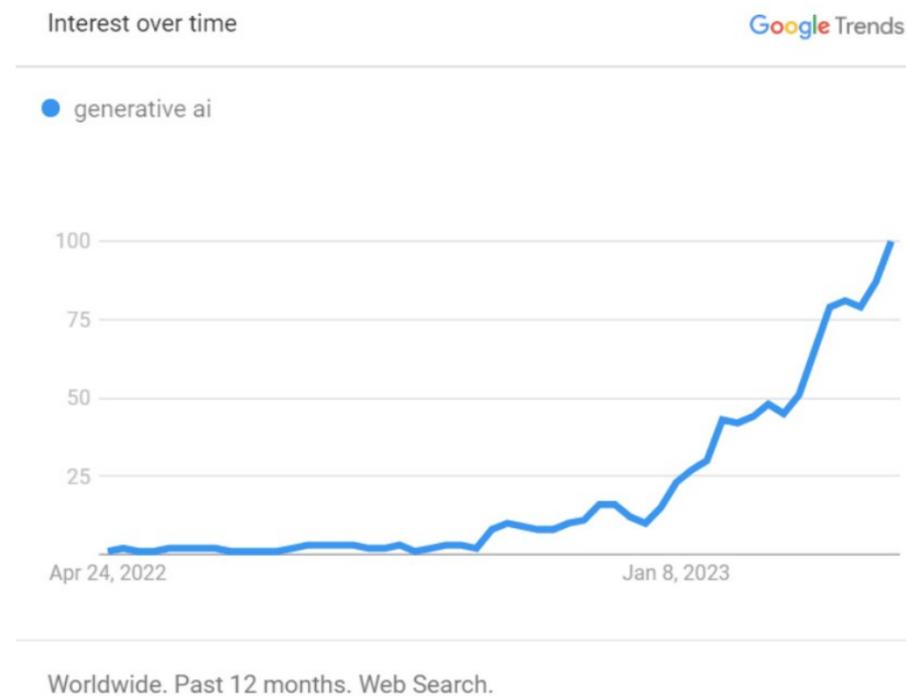
# Course Structure

# Motivation for this course

Over the last two years, you've probably heard all of these GenAI buzzwords:

- LLMs
- Chatbots
- RAG
- Agents
- Image/Video generation
- Deep Fakes

# Motivation for this course

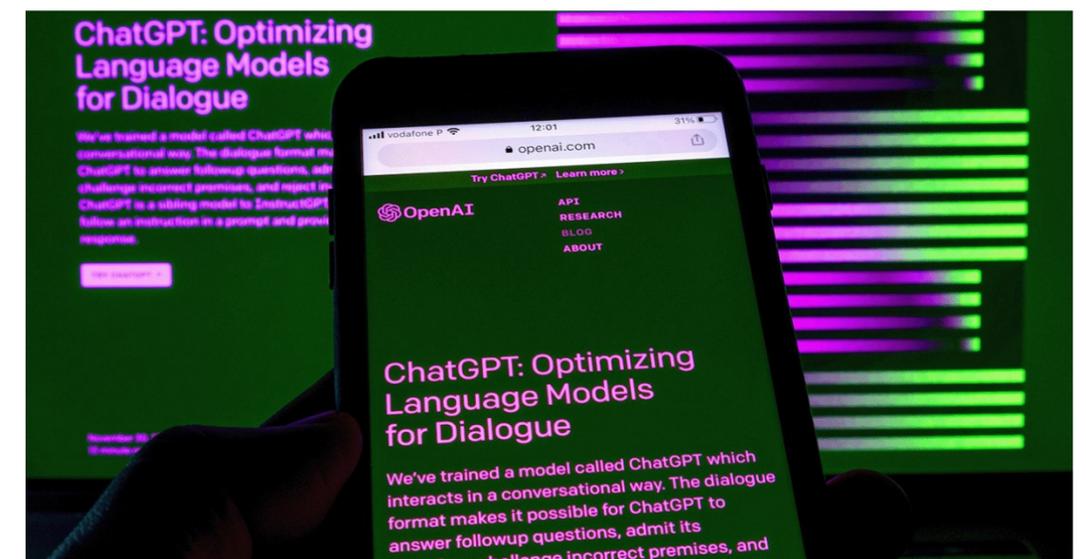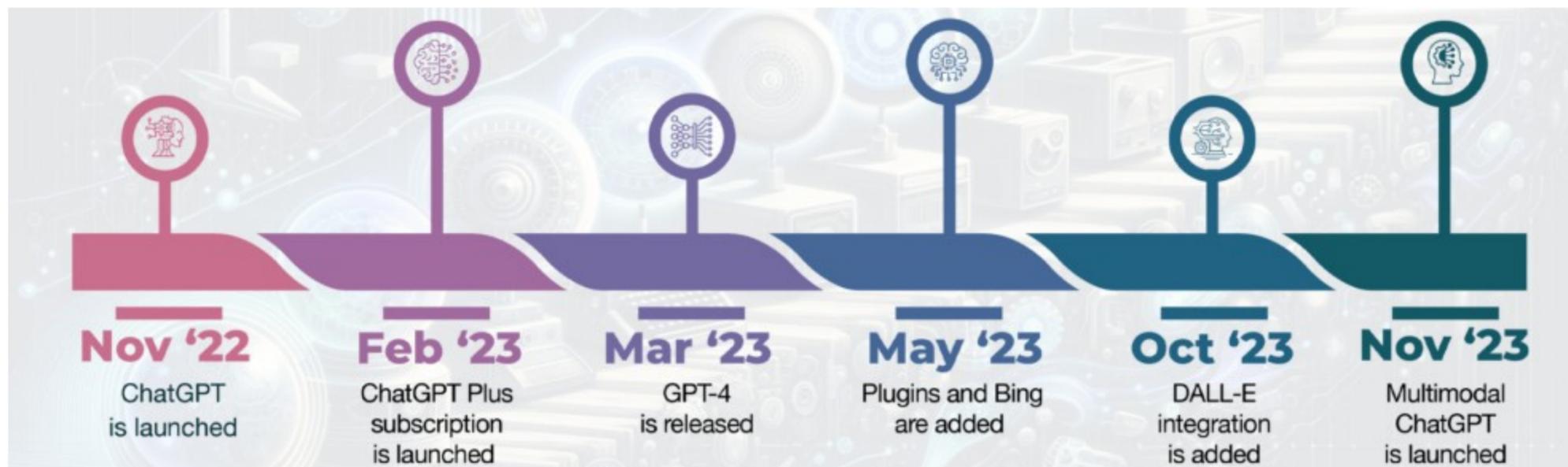ChatGPT Changed the world of knowledge forever

- Released Nov 2022 Took the world by storm
- How does it know so much?
- Is it "thinking" ?
- What is it made of?



ChatGPT Sprints to One Million Users
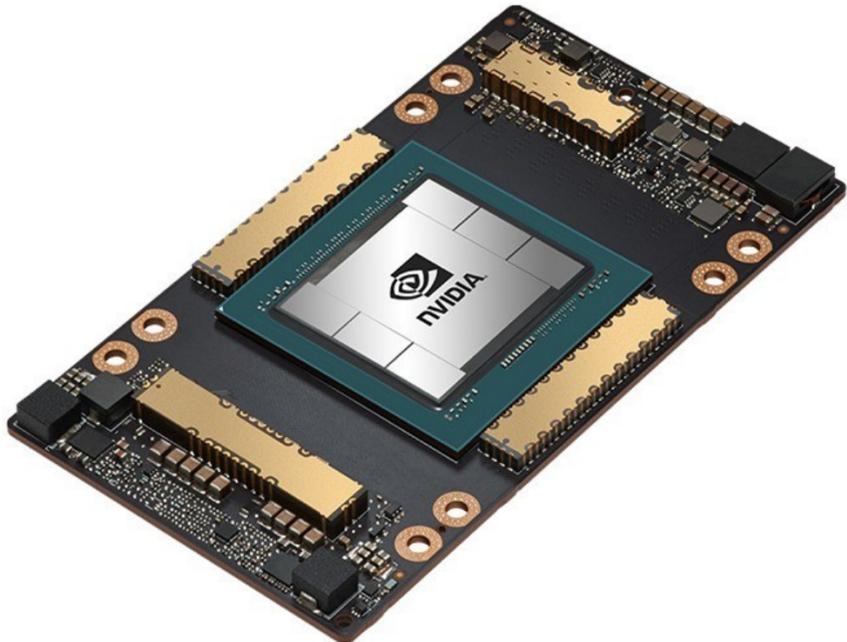Time it took for selected online services to reach one million users

| Service | Launched | Time |
|---|---|---|
| Netflix | 1999 | 3.5 years |
| Kickstarter* | 2009 | 2.5 years |
| Airbnb** | 2008 | 2.5 years |
| Twitter | 2006 | 2 years |
| Foursquare*** | 2009 | 13 months |
| Facebook | 2004 | 10 months |
| Dropbox | 2008 | 7 months |
| Spotify | 2008 | 5 months |
| Instagram*** | 2010 | 2.5 months |
| ChatGPT | 2022 | 5 days |

* one million backers   ** one million nights booked   *** one million downloads
Source: Company announcements via Business Insider/Linkedin

statista



**Nov '22** ChatGPT is launched

**Feb '23** ChatGPT Plus subscription is launched

**Mar '23** GPT-4 is released

**May '23** Plugins and Bing are added

**Oct '23** DALL-E integration is added

**Nov '23** Multimodal ChatGPT is launched



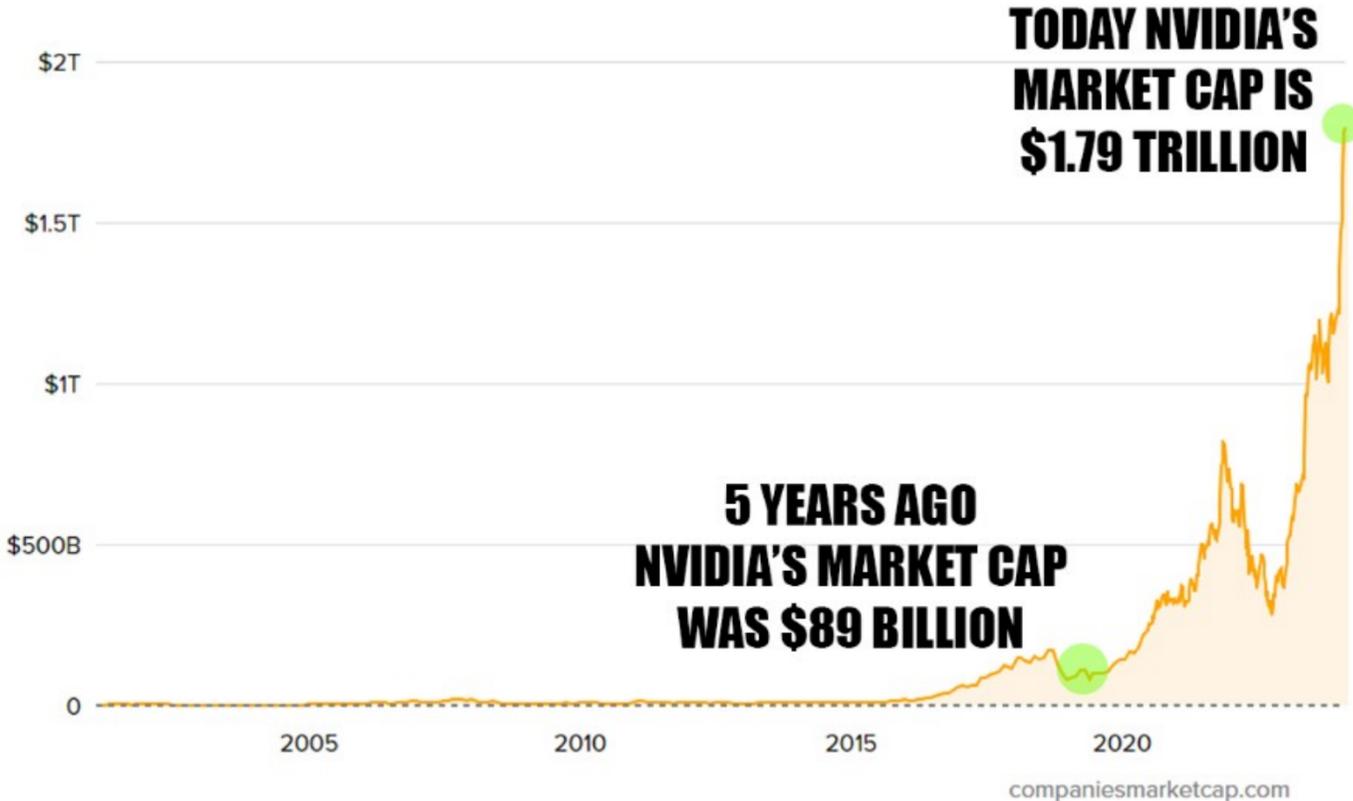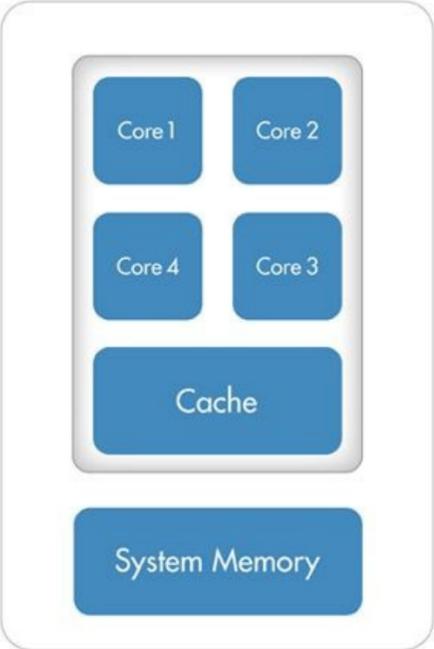ChatGPT: Optimizing Language Models for Dialogue

# Motivation for this course



In 2024, NVIDIA's market cap reached 3T! Why?
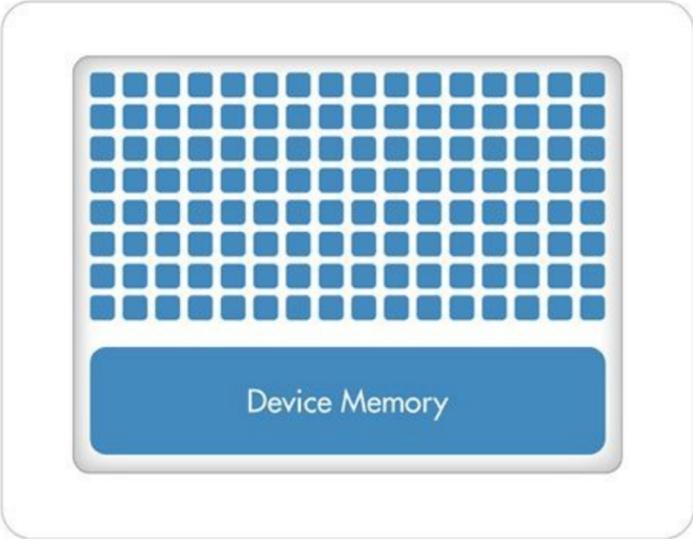
Why are GPUs now the hottest commodity for GenAI?



Market cap history of NVIDIA from 2001 to 2024

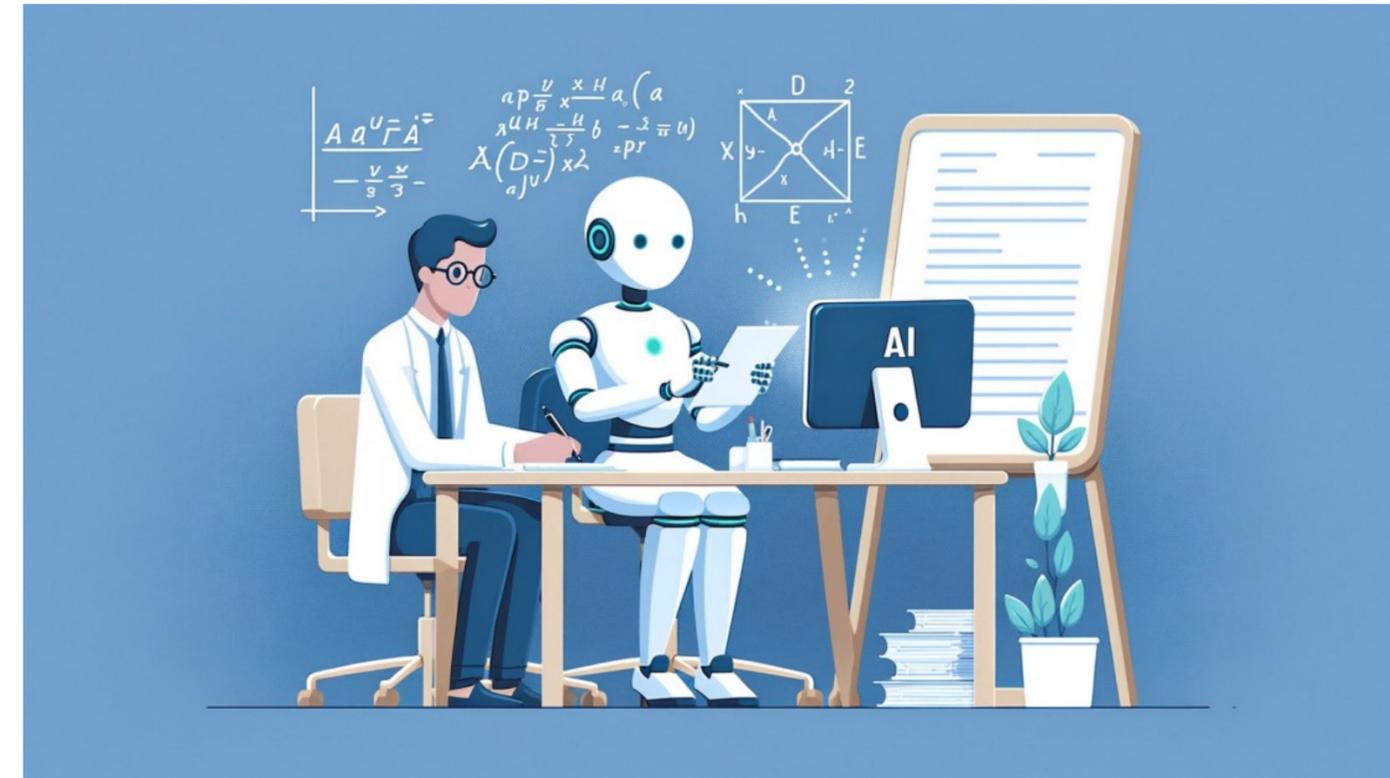TODAY NVIDIA'S MARKET CAP IS $1.79 TRILLION

5 YEARS AGO NVIDIA'S MARKET CAP WAS $89 BILLION

companiesmarketcap.com



**CPU (Multiple Cores)**

Core 1   Core 2

Core 4   Core 3

Cache

System Memory

**GPU (Hundreds of Cores)**

Device Memory

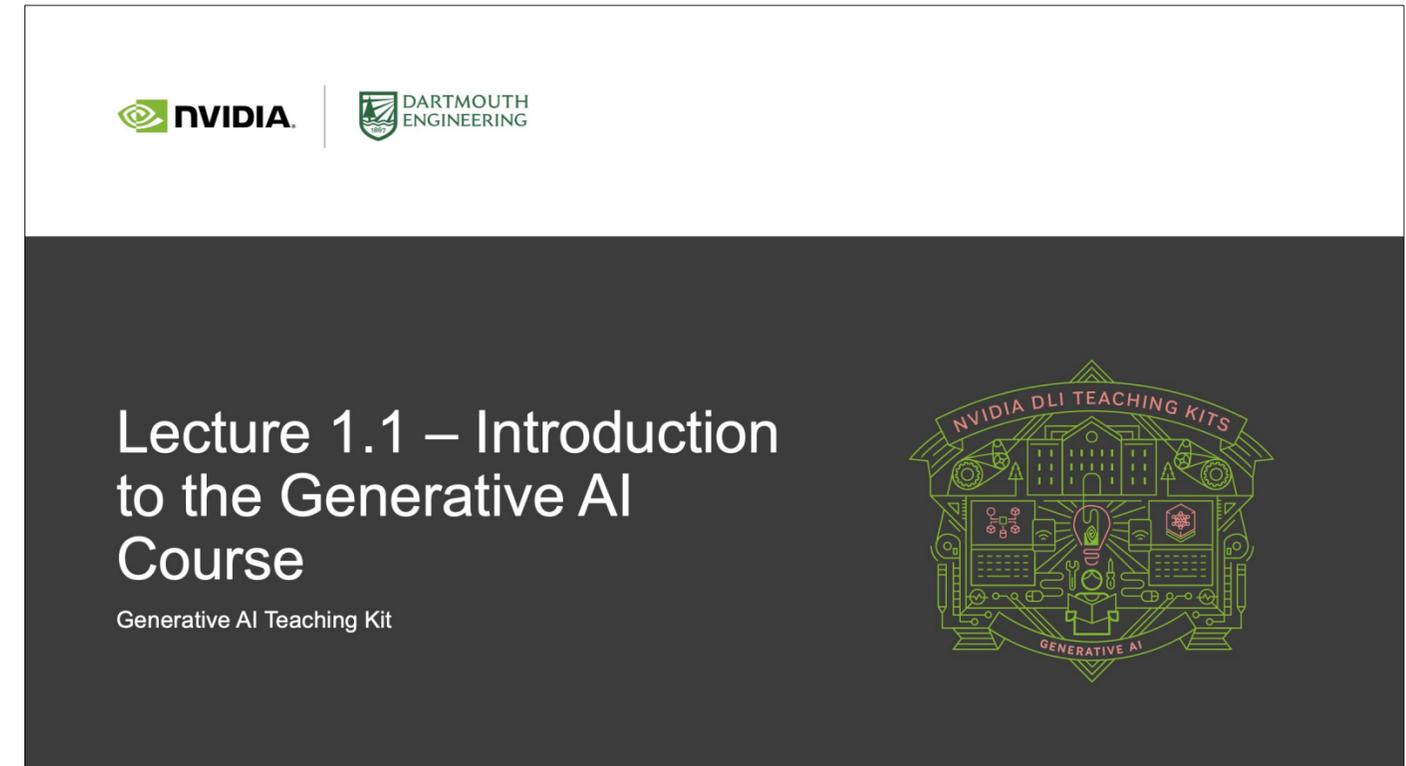DARTMOUTH ENGINEERING | nVIDIA

# Motivation for this course

- GenAI poses a challenge for Academia

- Pace of development has far outpaced adoption in schooling

- How do we do teaching given these extremely knowledgeable, highly accessible tools?

- Familiarity with the abilities and limits of these tools is vital to maintain competitiveness for both teachers and students

# Structure of the course

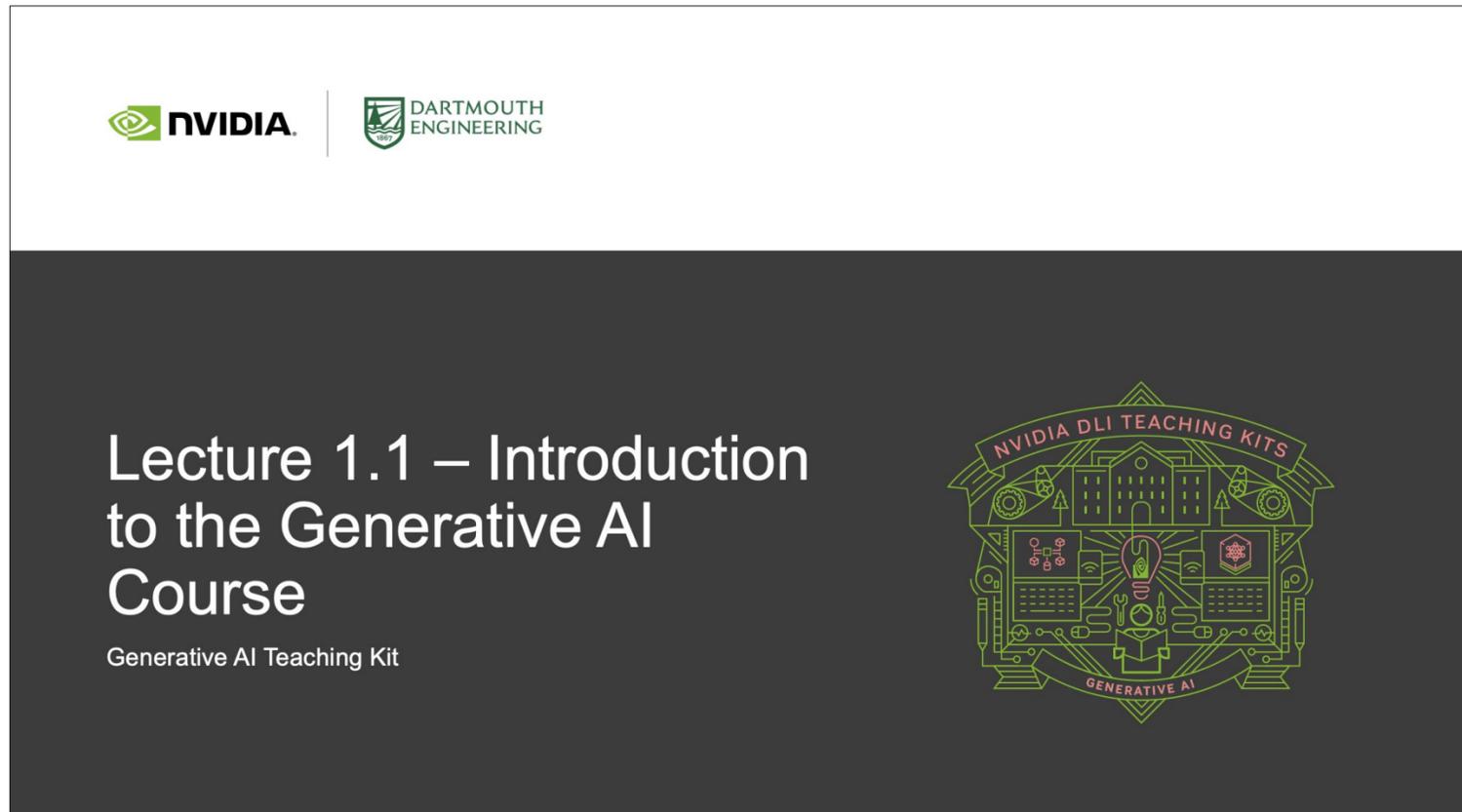How each module of the course is designed

- **[Online Syllabus](#)**
  - Lays out entire course including planned future content

- Lecture Slides
  - Introduce content

- Demo Notebooks
  - Show the tools in practice

- Knowledge Checks
  - Keeping track of the content

- Lab Notebooks
  - Your turn to build with the tools

- DLI Online Courses
  - Online, self-paced courses offering certificate



Lecture 1.1 – Introduction to the Generative AI Course

Generative AI Teaching Kit

# Structure of the course

Lecture Slides

- 2 or three lectures per module
- Designed to introduce the necessary information needed to understand what/how things work
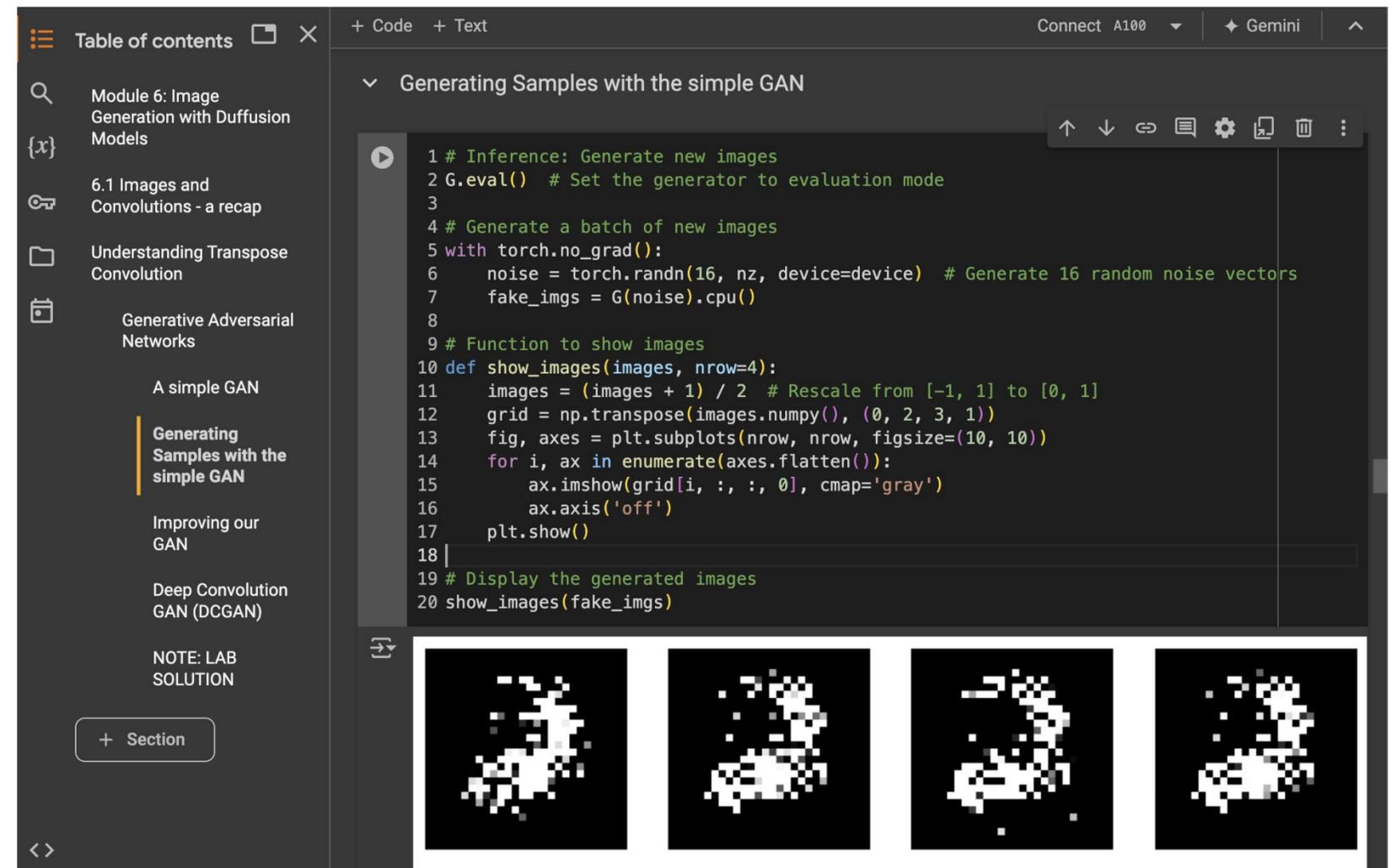- Not deep mathematical proofs, more applied science/engineering focus

# Structure of the course

Demo Notebooks

- Designed to be used during class time for the instructor to show how these tools/topics can work
- Built with an interactive Python notebook, namely Google Colab

# Structure of the course

## Knowledge Checks
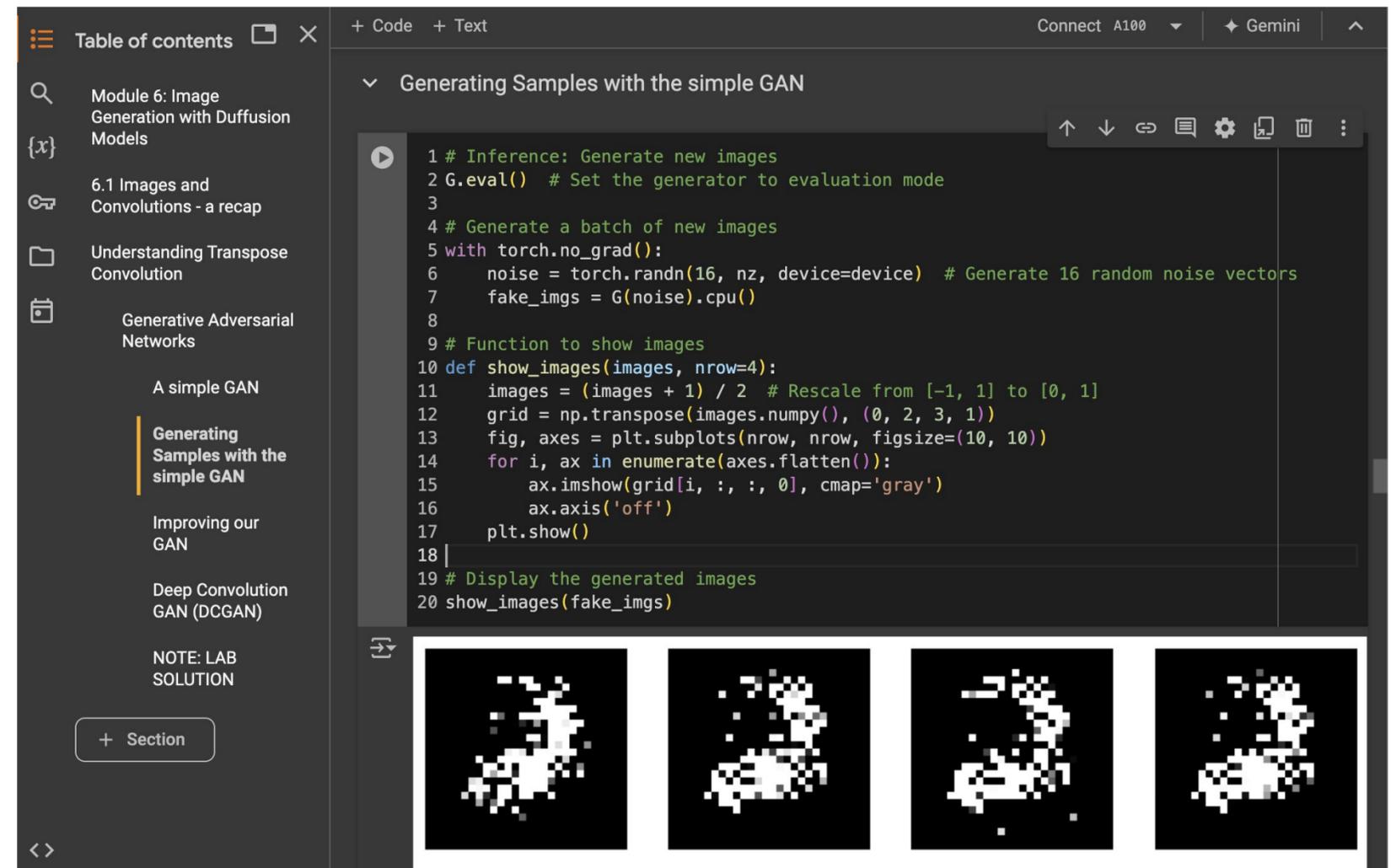
- Confirmation that you're developing a working knowledge of GenAI topics

| Question | Option A | Option B | Option C | Option D |
|---|---|---|---|---|
| What is the main advantage of Few-Shot Learning in LLMs? | Requires a large labeled dataset | Reduces the need for labeled data | Improves the model's interpretability | Requires extensive fine-tuning |
| How does self-supervised learning work? | It uses labeled data for predictions | It masks parts of the data as input/output pairs | It uses external labels for classification | It combines labeled and unlabeled data |
| What problem does self-supervised learning address? | It reduces the need for labeled data | It addresses data scarcity in supervised learning | It simplifies the model training process | It ensures consistency in model predic |
| What is the role of autoregressive learning in LLMs? | Predicts the next token based on prior tokens | Enhances the model's output accuracy | Identifies relationships between labeled data | Increases the model's parameter coun |
| What is Few-Shot Learning in the context of GPT-3? | Training the model with many labeled examples | Providing a few examples before the main input | Building a single model for multiple tasks | Training the model on multiple dataset |
| How does in-context learning differ from fine-tuning? | In-context learning is task-specific | In-context learning uses the model's parameters flexibly | In-context learning is only used during training | In-context learning replaces fine-tunin |
| What is the purpose of prompt templating in LLMs? | To create unique responses | To format inputs consistently | To standardize outputs | To introduce randomness in outputs |
| Why is temperature an important factor in LLM outputs? | It controls the sequence length | It affects the randomness of outputs | It sets the learning rate | It determines the output length |

# Structure of the course

Lab Notebooks

- Based on the tutorial notebooks

- These will give you the chance to implement the tools and explore the topics covered in each module

# Structure of the course



DLI Online Courses

- Related self-contained, self-paced DLI courses

- Free for students using promo codes provided

- Runs on DLI Platform instead of Colab

- Offers student certificates of competency based on built-in assessments

# Tools used in GenAI

# GenAI Tools in this Course

**Python** is the most widely used language in the AI/ML community, making it the go-to language for developing Generative AI models.
**Python's** simple syntax and readability allow developers and researchers to quickly prototype and implement complex algorithms.



- TensorFlow & PyTorch: Essential for building and training deep learning models, which are at the core of GenAI.
- Hugging Face Transformers: A popular library for working with large language models (LLMs), crucial for tasks like text generation and RAG.
- OpenCV & PIL: Important for image and video manipulation in GenAI applications like deep fakes and image generation
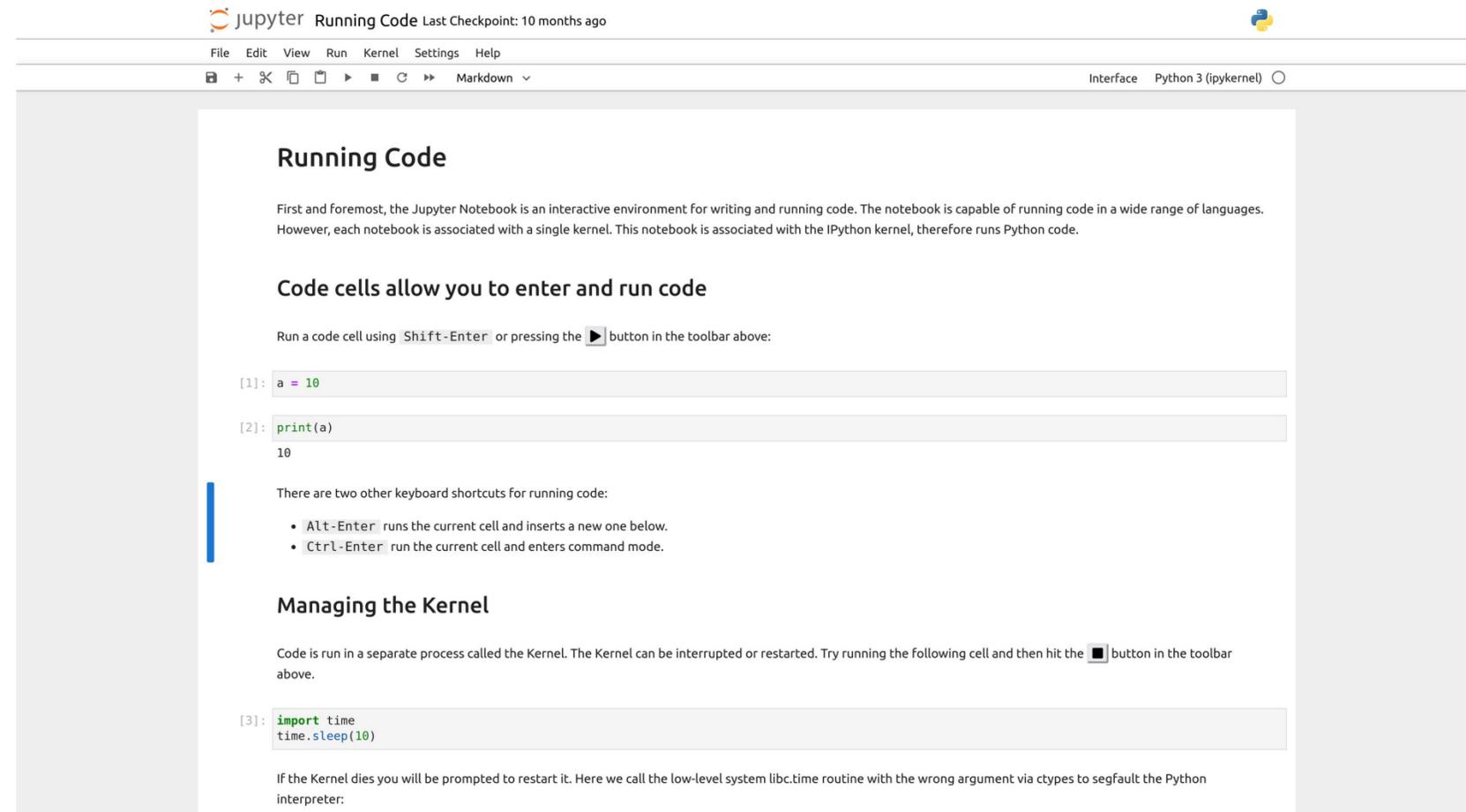
# GenAI Tools in this Course

Coding Notebooks

The primary platform to work with these AI tools and code libraries will be through the Google Colab notebooks

These notebooks allow for easy interaction with the data and models in a visual manner

Platforms like Google Colab and other cloud notebook providers also supply access to cloud-based GPUs and compute infrastructure.
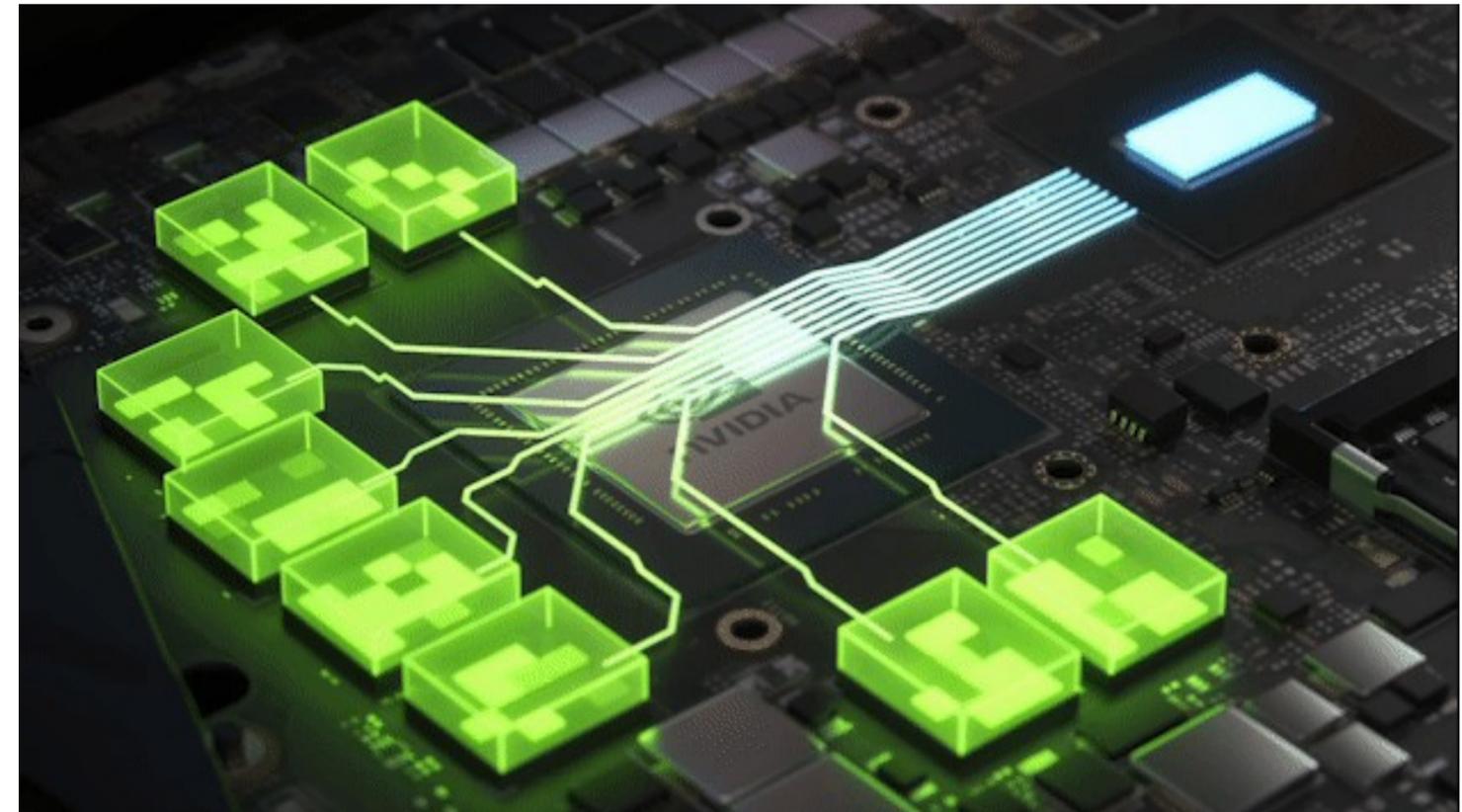
# GenAI Tools in this Course

## The Importance of GPUs in GenAI

**Parallelism**: GPUs have thousands of cores optimized for parallel processing, allowing them to handle many operations at once. This is essential for the matrix multiplications and other computations involved in AI inference.

**Throughput**: GPUs can process large batches of data simultaneously, which is critical for applications like image generation or video synthesis where high throughput is needed.

**Latency**: GPUs reduce latency by executing multiple operations in parallel, which speeds up the time it takes to get from input to output in GenAI models.



DARTMOUTH ENGINEERING | NVIDIA.

# Course Content

# What this course will cover

**Modules**

1. Introduction to Generative AI
2. Word Embeddings, Tokens, and NLP
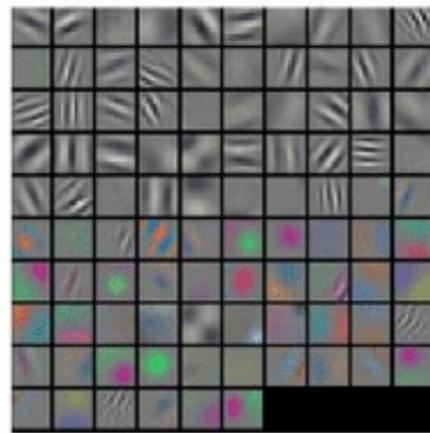3. Large Language Models and the Transformer
4. LLM Scaling Laws and LLM Families
5. Multimodal Learning and its Applications
6. Diffusion Models in Generative AI
7. Model Training (Pre-Training, Instruction Following, and PEFT)
8. LLM Orchestration
9. Scaling Model Training to Distributed Workloads

**Full online syllabus**

# 2. Word Embeddings, Tokens, and NLP



**Tiktokenizer**

gpt-3.5-turbo

| System | You are a helpful assistant | ✕ |
| User | Content | ✕ |

**Add message**

```
<|im_start|>system
You are a helpful assistant<|im_end|>
<|im_start|>user
<|im_end|>
<|im_start|>assistant
```

Token count
18

Price per prompt
$0.000018

```
<|im_start|>system
You are a helpful assistant<|im_end|>
<|im_start|>user
<|im_end|>
<|im_start|>assistant
```

```
[100264, 9125, 198, 2675, 527, 264, 11190, 18328, 1002
65, 198, 100264, 882, 198, 100265, 198, 100264, 78191,
198]
```

☐ Show whitespace

Built by **dqbd**. Created with the generous help from **Diagram**.

|  | living being | feline | human | gender | royalty | verb | plural |
|---|---|---|---|---|---|---|---|
| cat → | 0.6 | 0.9 | 0.1 | 0.4 | -0.7 | -0.3 | -0.2 |
| kitten → | 0.5 | 0.8 | -0.1 | 0.2 | -0.6 | -0.5 | -0.1 |
| dog → | 0.7 | -0.1 | 0.4 | 0.3 | -0.4 | -0.1 | -0.3 |
| houses → | 0.8 | -0.4 | -0.5 | 0.1 | -0.9 | 0.3 | 0.8 |
| man → | 0.6 | -0.2 | 0.8 | 0.9 | -0.1 | -0.9 | -0.7 |
| women → | 0.7 | 0.3 | 0.9 | -0.7 | 0.1 | -0.5 | -0.4 |
| king → | 0.5 | -0.4 | 0.7 | 0.8 | 0.9 | -0.7 | -0.6 |
| queen → | 0.8 | -0.1 | 0.8 | -0.9 | 0.8 | -0.5 | -0.9 |

*Dimensionality reduction of word embeddings from 7D to 2D*

Word | Word Embedding | Dimensionality reduction | Visualization of Word embeddings in 2D

DARTMOUTH ENGINEERING | NVIDIA

# 3. Large Language Models and the Transformer

# 4. LLM Scaling Laws and LLM Families

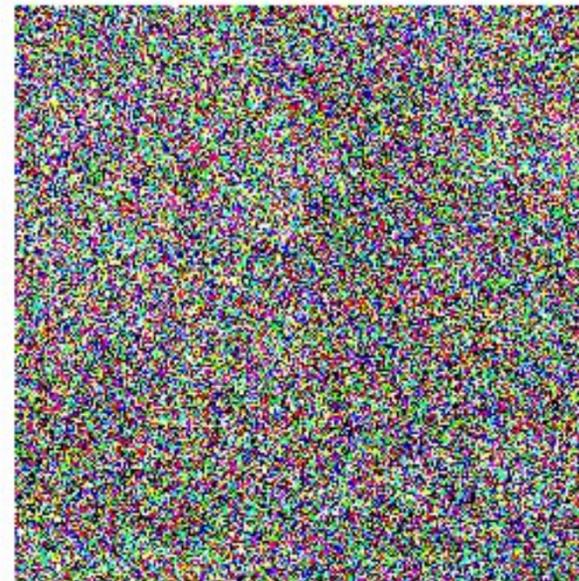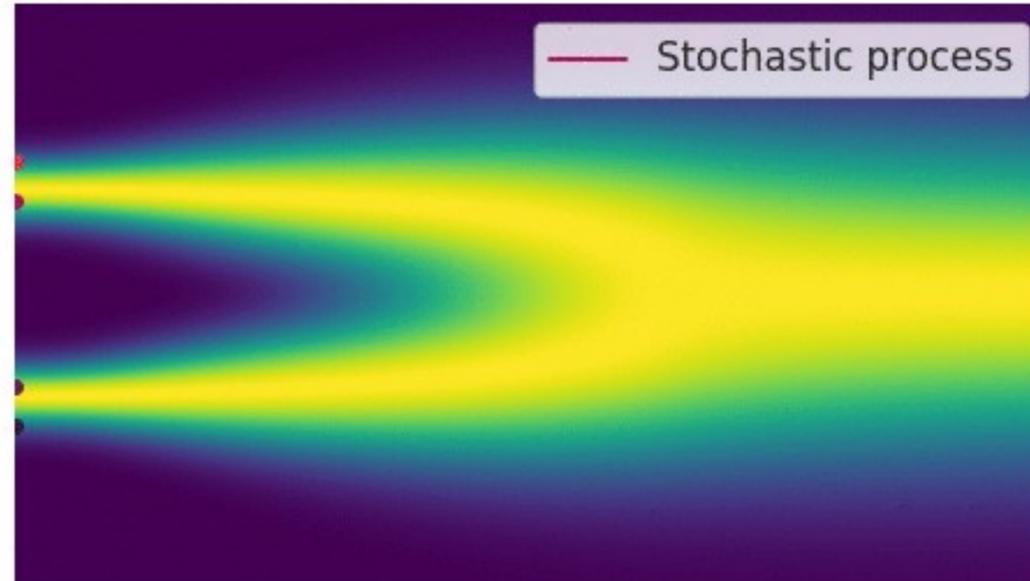# 5. Multimodal Learning and its Applications



Alexnet 1st conv filters

ViT 1st linear embedding filters
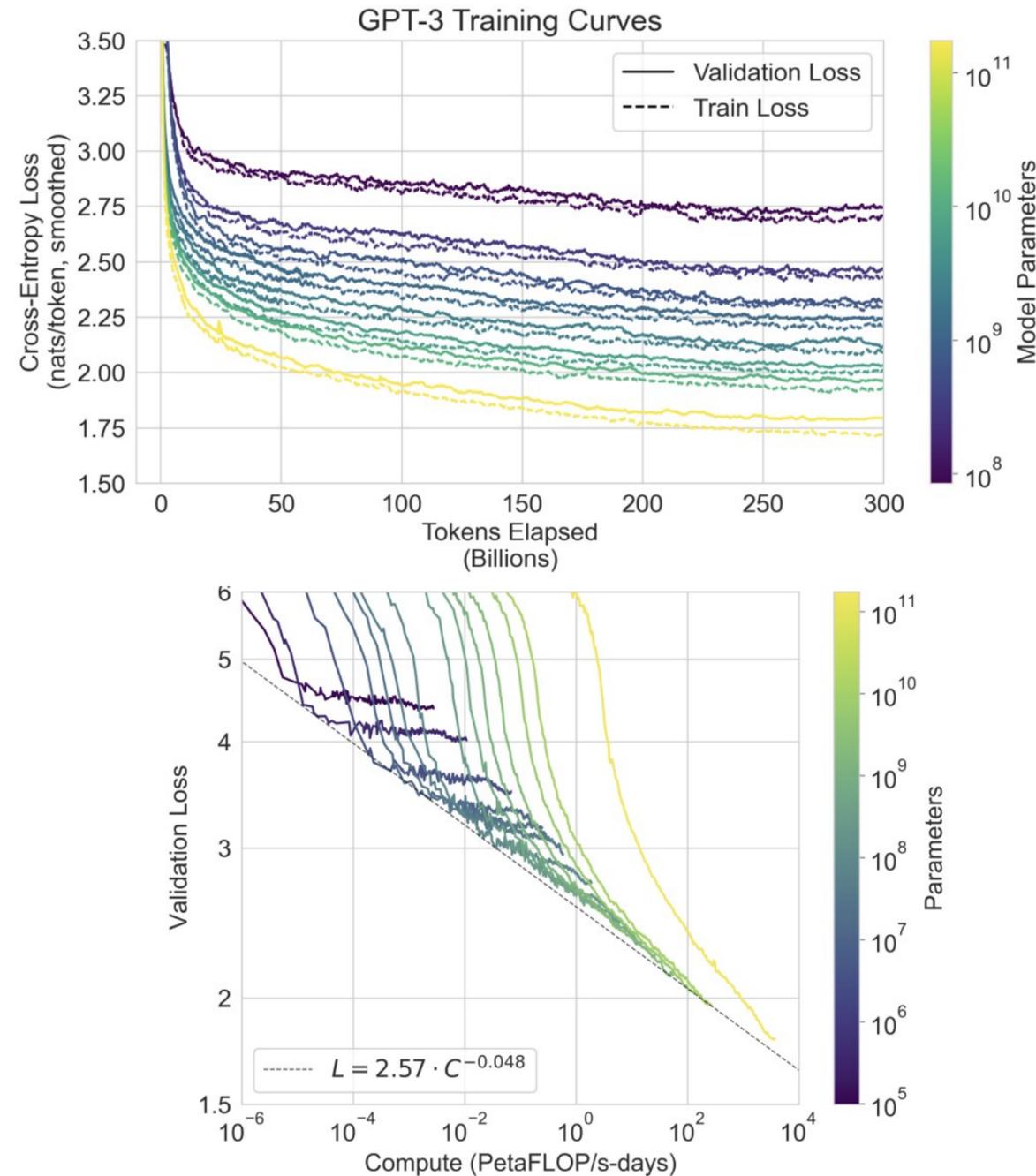
RGB embedding filters
(first 28 principal components)

**1. Contrastive pre-training**

# 6. Diffusion Models in Generative AI

# 7. Model Training (Pre-Training, Instruction Following, and PEFT)

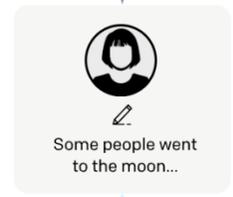# 8. LLM Orchestration

# 9. Scaling Model Training to Distributed Workloads

# Wrap Up

Introduction to GenAI

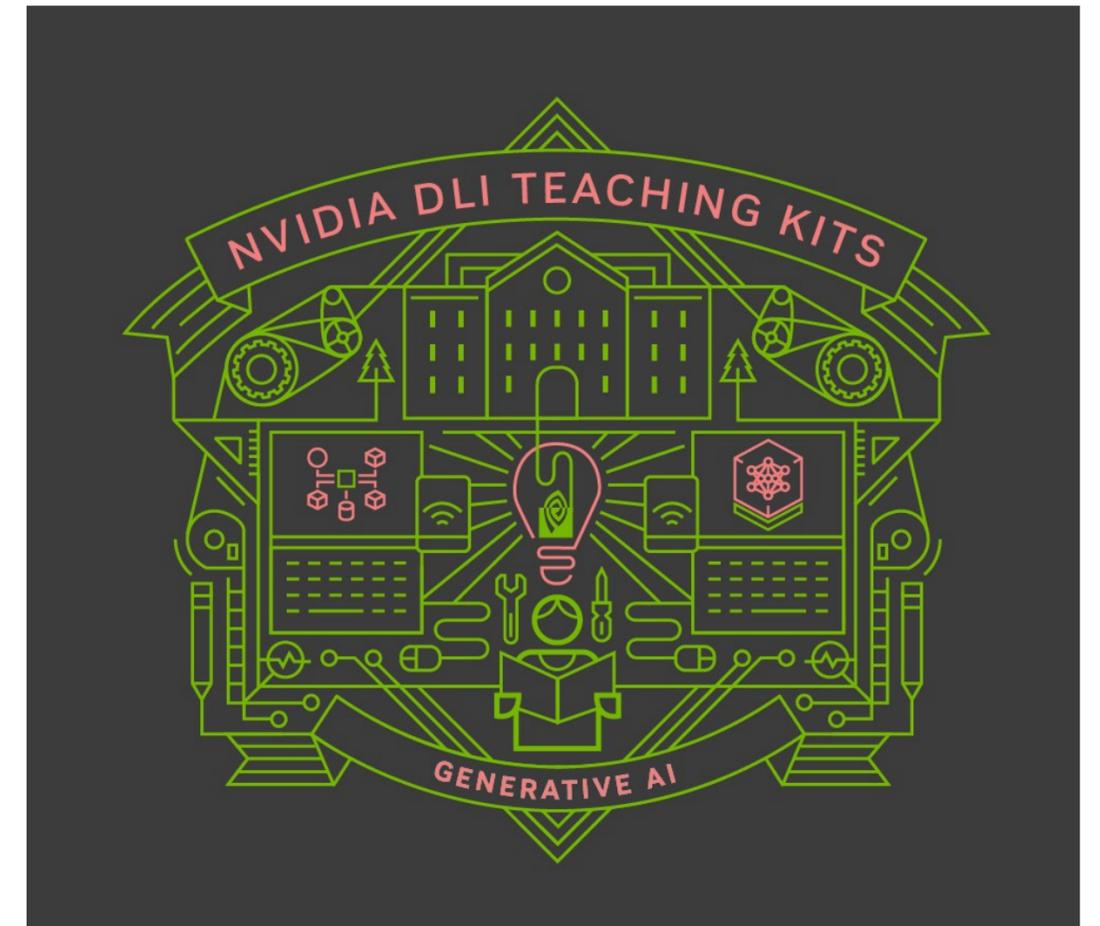- Today we introduced the GenAI Course

- We covered the different materials used in the class

- Provided an overview of the modules to come

- Highlighted how GenAI has exploded onto the scene of science, engineering, and the economy at large

----------------------------------------------------------------------

In the next lesson we will start our discussion into how we got to this point and what advances led us from the digital computer, to ChatGPT

**Thank you!**