



Practical AI: Definitions, Concepts, & Techniques

Lee Solomonson | December 2024

Director of Advanced Technologies

lee.solomonson@nth.com





WELCOME!

**The webinar will begin shortly.
We are waiting for everyone to load.**



Before We Begin...



- Microphones have been muted – mics will open at the end.
- Please use the Q&A function for questions you would like to see addressed.
- Slides will be provided.
- This is an educational event focused on the technical concepts and techniques around AI. Specific solution offerings will be introduced, with in-depth explanations offered at a later date.

Section 1

- AI/ML/DL
- Deep Learning Phases (Explore/Deploy/Apply)
- Predictive AI vs. Generative AI
- Functions and Use Cases

Section 2

- What is a Neural Network?
 - Networks and Architectures
 - Neurons, Weights, Bias, Activation
 - Matrix Multiplication
- Model Training Overview

Section 3

- Model Selection & Sources
 - Quantization

Section 4

- Model Deployment
- Applying Model
 - Prompt Engineering
 - Retrieval Augmented Generation (RAG)
 - Foundational Models
- Getting Started



Poll: What is your experience level?

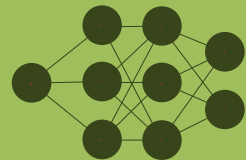


Artificial Intelligence

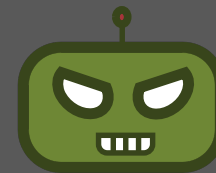
Artificial Narrow Intelligence (ANI)

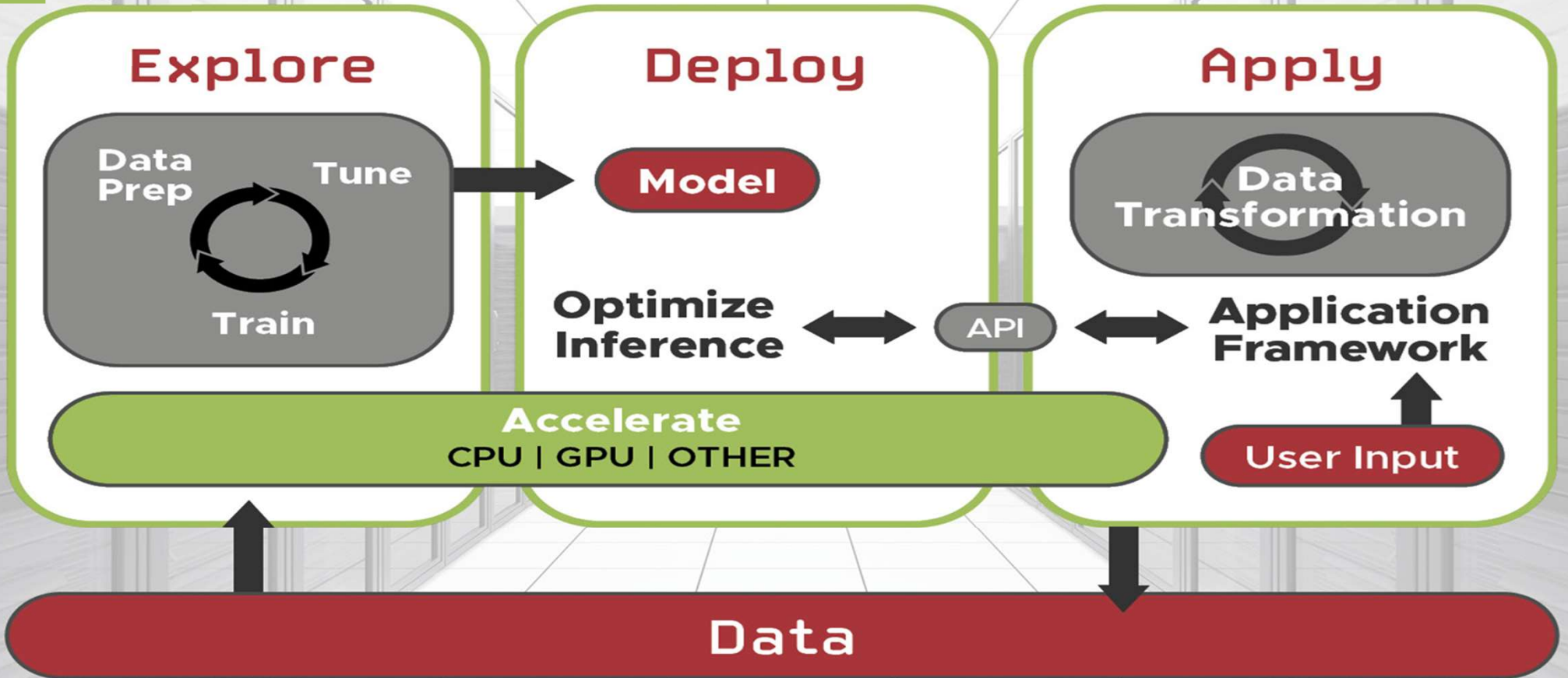
Machine Learning

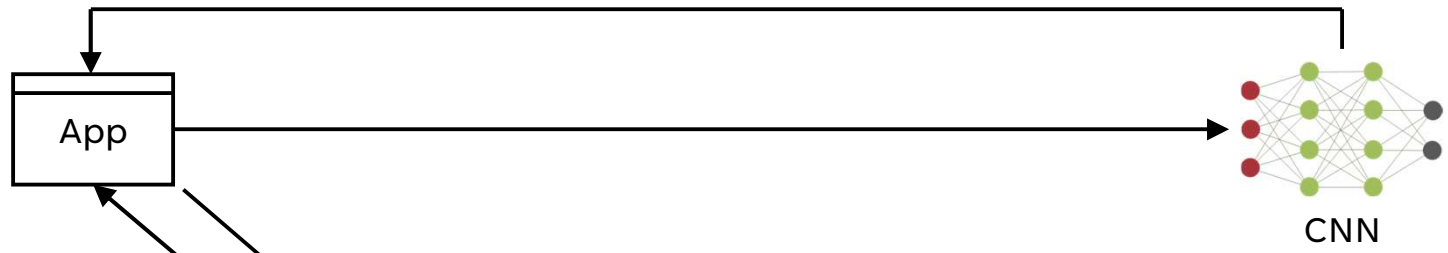
Deep Learning

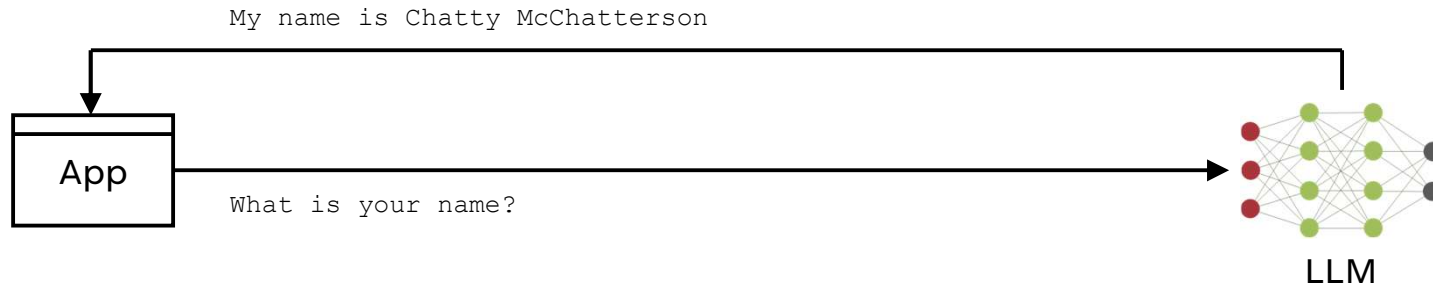


Artificial General Intelligence (AGI)









“You shall know a word by the company it keeps.”
 --J.R. Firth, 1957

It is windy outside.
 The road is windy outside of town.

```

    What is your name?
    My Name is Chatty McChatterson
    what is 397652 X 1.5342?
  
```

```

    [3923, 374, 701, 836, 1980, 5159, 4076, 374, 921, 23758, 45
    84, 1163, 1617, 1293, 271, 12840, 374, 220, 20698, 23181, 1
    630, 220, 16, 13, 22467, 17, 30]
  
```

<https://tiktokenizer.vercel.app>

INPUT DATA

EXAMPLES OF OUTPUT BY INDUSTRY

Ask an interesting business question	Identify the appropriate DL task
Is 'it' present or not?	Detection
What type of thing is 'it'?	Classification
To what extent is 'it' present?	Segmentation
What is the likely outcome?	Prediction
What will likely satisfy the objective?	Recommendation



Healthcare	Retail	Finance
Cancer Detection	Targeted Ads	Cybersecurity
Image Classification	Basket Analysis	Credit Scoring
Tumor Size / Shape Analysis	Build 360° Customer View	Credit Risk Analysis
Survivability Prediction	Sentiment & Behavior Recognition	Fraud Detection
Therapy Recommendation	Recommendation Engine	Algorithmic Trading

Image provided courtesy of NVIDIA

Generative AI Use Cases



Intelligent Chatbot

Focus is on question-and-answer tasks.

Ex. Customer Service Agent, Brand Ambassador, Help Desk



Knowledge Base Copilot

Connects to knowledge bases performs tasks such as writing, coding, generating images, etc.

Ex. Documentation Copilot, IT Bugs Assistant, Field Agent Copilot



Code Generation

Help develop or troubleshoot code based on natural language. Can work across common languages or be proprietary languages.

Ex. GitHub Copilot, ChatUSD, Software Development Assistant

Image provided courtesy of NVIDIA

Multimodal Conversion





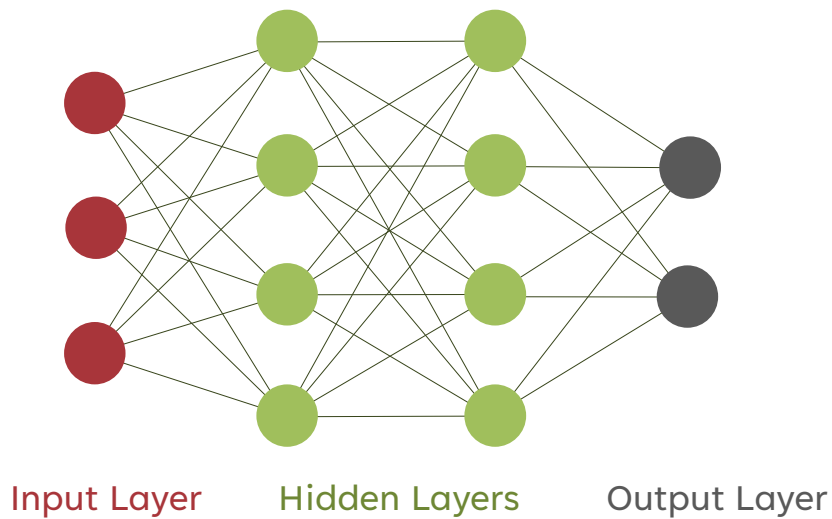
Demonstration



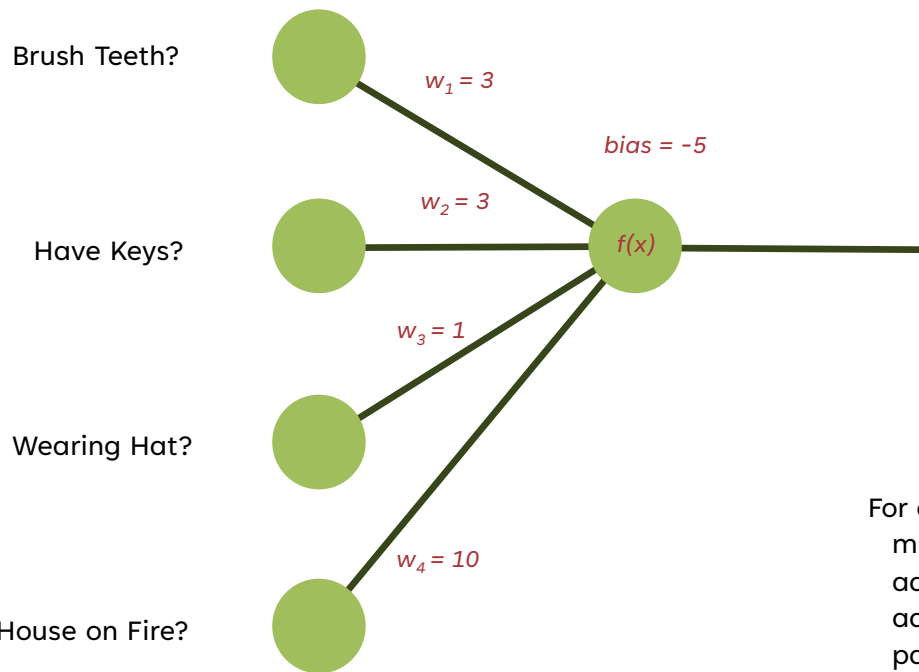
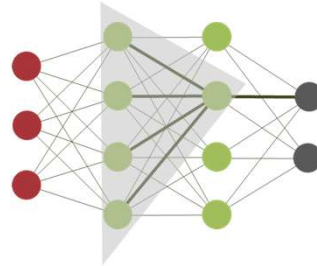


Up Next: Anatomy of a Neural Network





Artificial Neural Network	Description/Purpose	Typical LT
Convolutional Neural Networks (CNNs)	Classification & Recognition	Supervised
Recurrent Neural Networks (RNNs)	Sequential Patterns	Supervised
Transformer Networks (LLM)	Comprehension of Context	Semi-Supervised
Generative Adversarial Networks (GANs)	Competing ANNs (Generator & Discriminator)	Unsupervised & Semi-Supervised
Diffusion Model	Forward Diffusion & Reverse De-Noise	Semi-Supervised



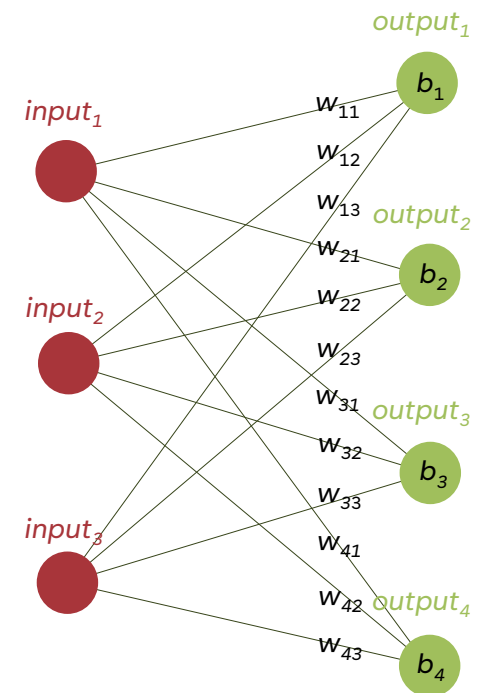
$w_1 * I$	$w_2 * I$	$w_3 * I$	$w_4 * I$	bias	value	result
Teeth	Keys	Hat	Fire			
3	3	1	0	-5	2	active
3	0	1	0	-5	-1	inactive
0	0	0	10	-5	5	active
3	3	1	10	-5	12	active

For each output node:
 multiply weight and input value,
 add the weighted inputs together,
 add bias to sum
 pass results through activation function

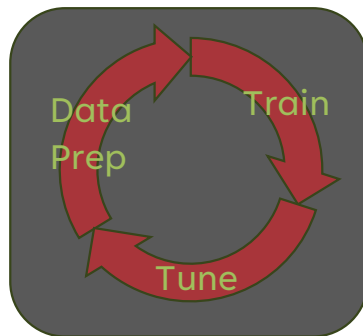
$$\text{activation} \left(\begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \\ w_{41} & w_{42} & w_{43} \end{bmatrix} \times \begin{bmatrix} \text{input}_1 \\ \text{input}_2 \\ \text{input}_3 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix} \right) = \begin{bmatrix} \text{output}_1 \\ \text{output}_2 \\ \text{output}_3 \\ \text{output}_4 \end{bmatrix}$$

$$\begin{bmatrix} w_{11} \times \text{input}_1 + w_{12} \times \text{input}_2 + w_{13} \times \text{input}_3 \\ w_{21} \times \text{input}_1 + w_{22} \times \text{input}_2 + w_{23} \times \text{input}_3 \\ w_{31} \times \text{input}_1 + w_{32} \times \text{input}_2 + w_{33} \times \text{input}_3 \\ w_{41} \times \text{input}_1 + w_{42} \times \text{input}_2 + w_{43} \times \text{input}_3 \end{bmatrix}$$

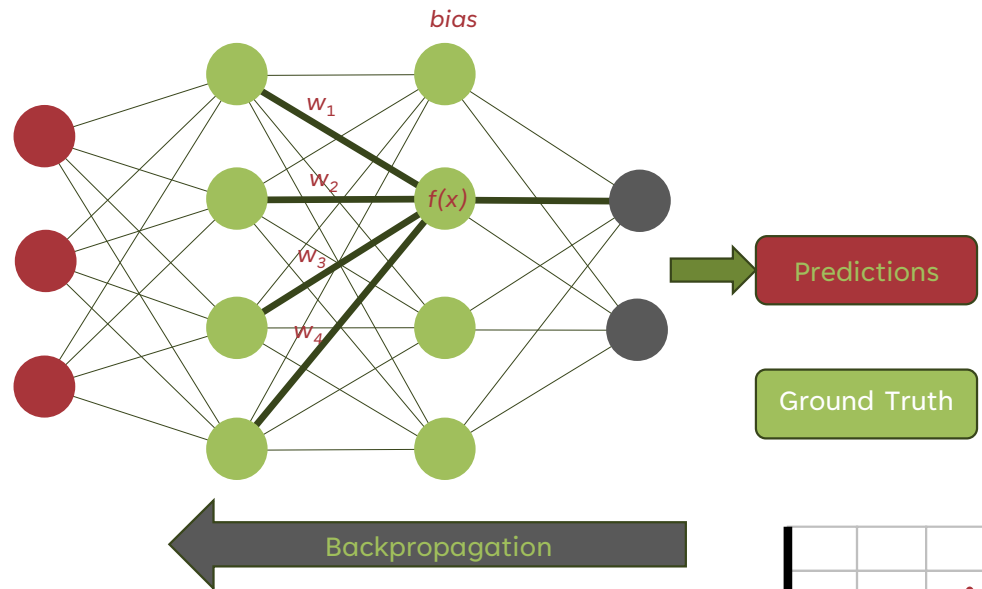
For each output node:
 multiply weight and input value,
 add the weighted inputs together,
 add bias to sum
 pass results through activation function



- Labeling
 - Missing/Incorrect Data
 - Normalization/Standardization
 - Augmentation
 - Split
- Train / Validate / Test

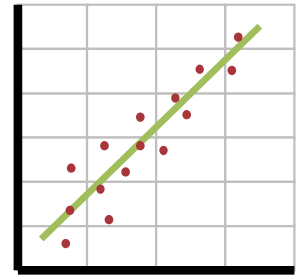


- Fine Tune vs. Retraining
- Hyper Parameter Adjustments
- Pre-Trained Model Selection



$$y = mx + b$$

$x = \text{input(s)}$
 $y = \text{output}$



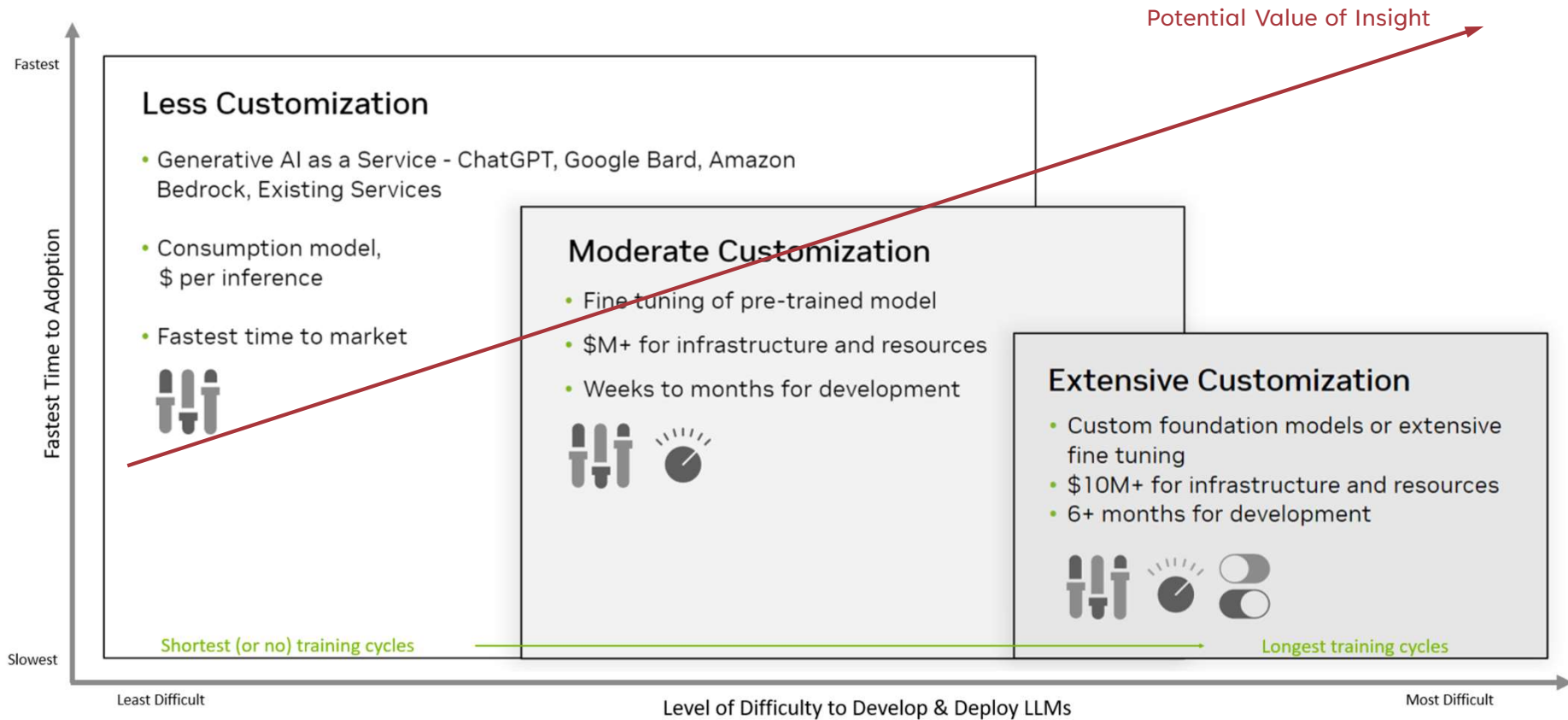
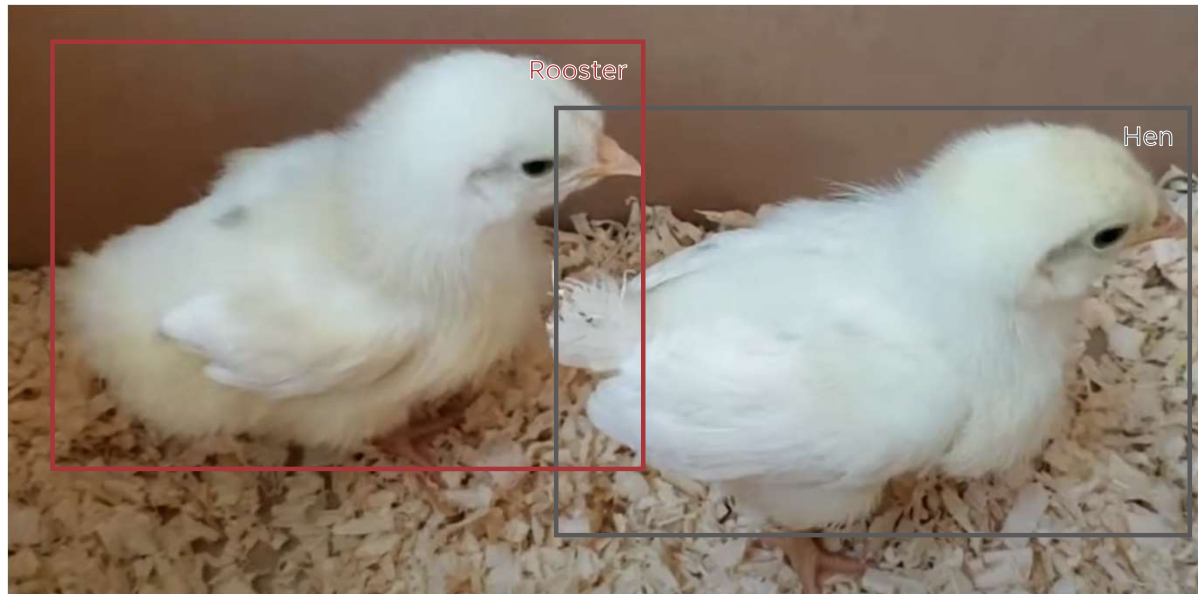


Image provided courtesy of NVIDIA





Up Next: Sourcing Foundational Models





Hugging Face



NGC Catalog

Model Repository Features:

- Model Card
- Live Interface
- Code Examples

Key Terminology:

- Context Window / Length
- Parameter
- Quantization
- Licensing

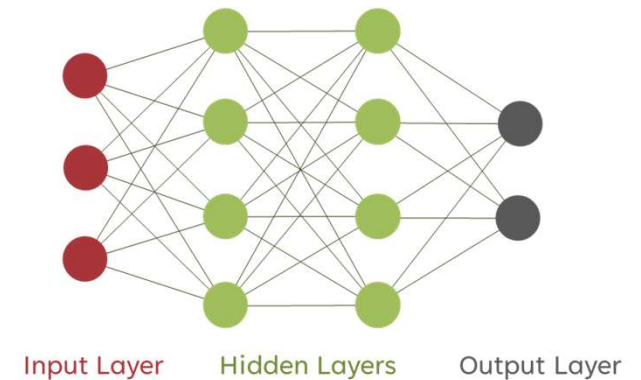
Token count
19

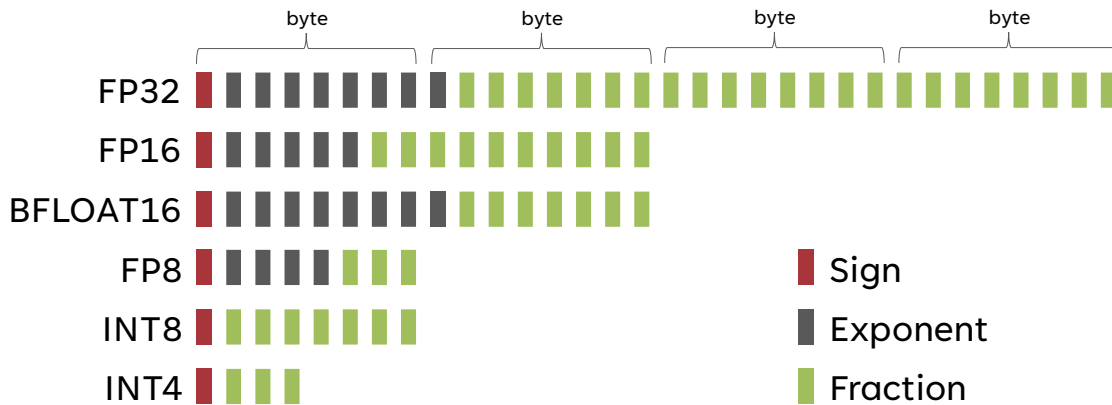
This is an example of token usage.

Prompts and responses all account for the context window.

2028, 374, 459, 3187, 315, 4037, 10648, 382, 36286, 13044, 323, 14847, 682, 2759, 369, 279, 2317, 3321, 4286

<https://tiktokenizer.vercel.app>





Floating Points:

$$\text{value} = (-1)^{\text{sign}} \times 2^{(E-127)} \times \left(1 + \sum_{i=1}^{23} b_{23-i} 2^{-i} \right)$$

Equation source: Wikipedia

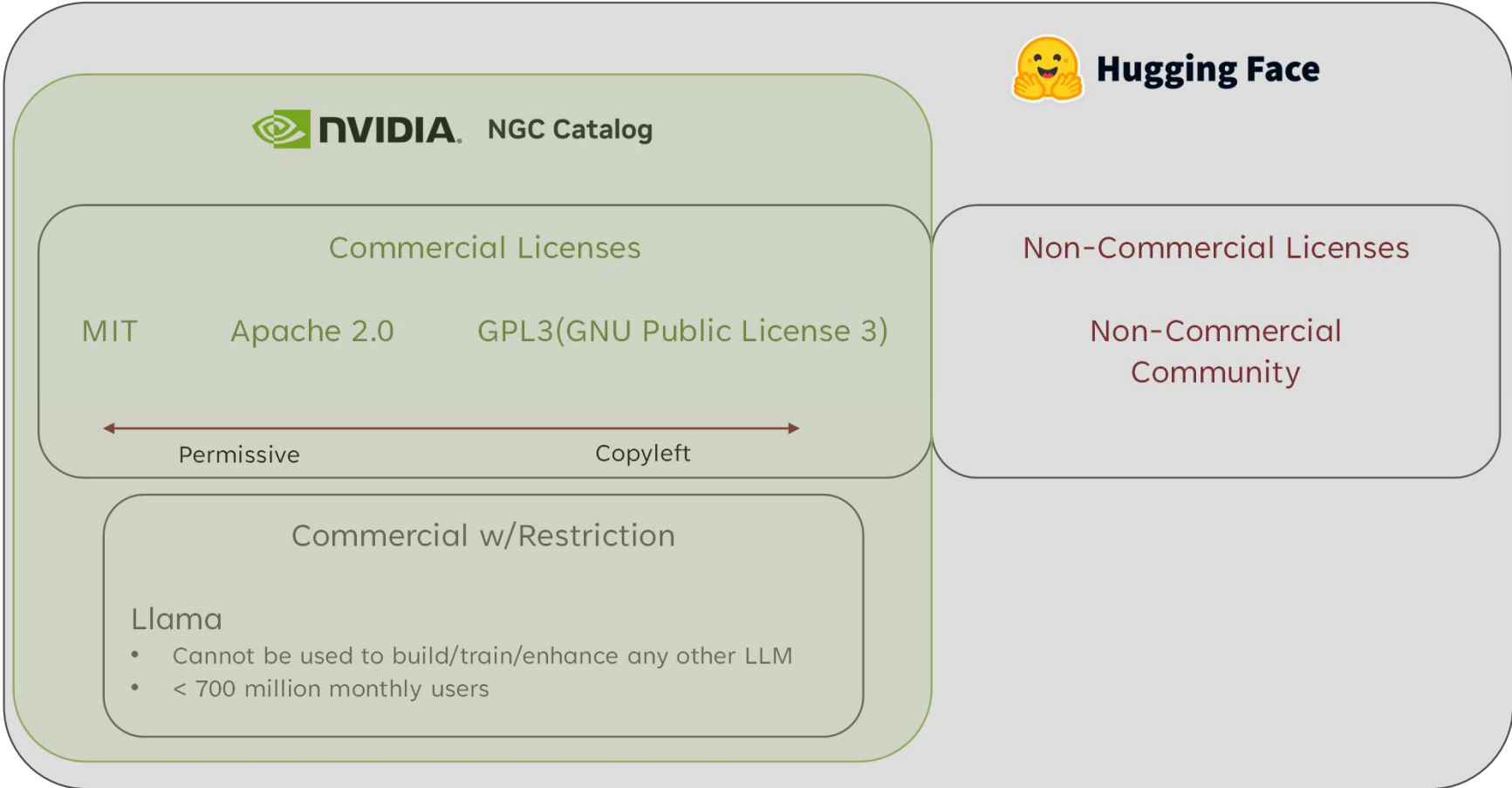
$$FP\text{value} = (-1)^{\text{sign}} * 1.[\text{fraction}] * 2^{[\text{exponent}] - (\text{highBit} - 1)}$$

$$3.1416 = (-1)^0 * 1.[5708] * 2^1$$

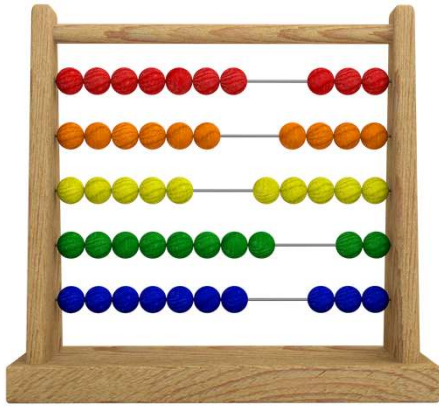
More Exponent Bits = Larger Range
 More Fraction Bits = Greater Precision

Acroyms:

- NeMo (Neural Modules)
- ONNX (Open Neural Net Exchange)
- GGUF (GPT-Generated Unified Format)
- PTQ (Post Training Quantization)
- QAT (Quantization-Aware Training)



This is not legal advice



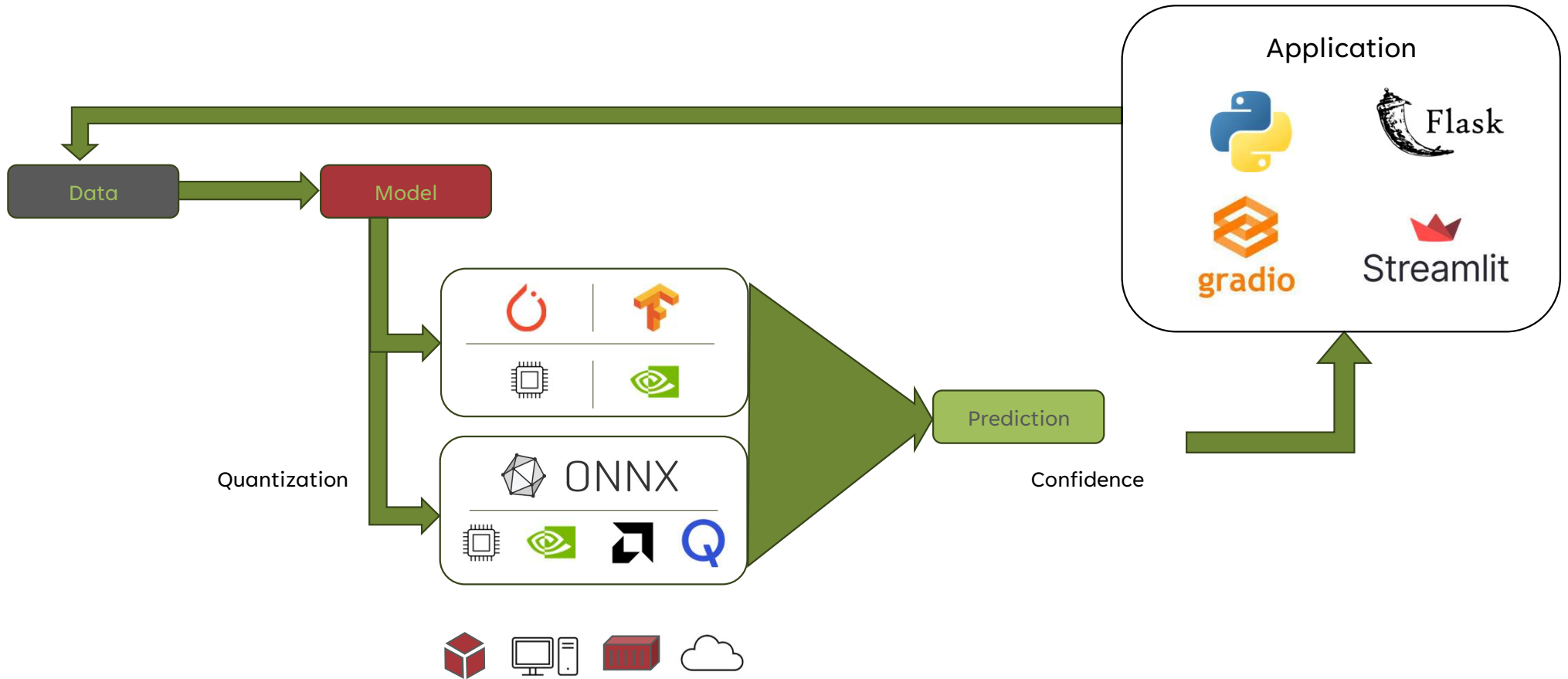
Cost
Precision
Speed
Size & Scalability
Ease of Use

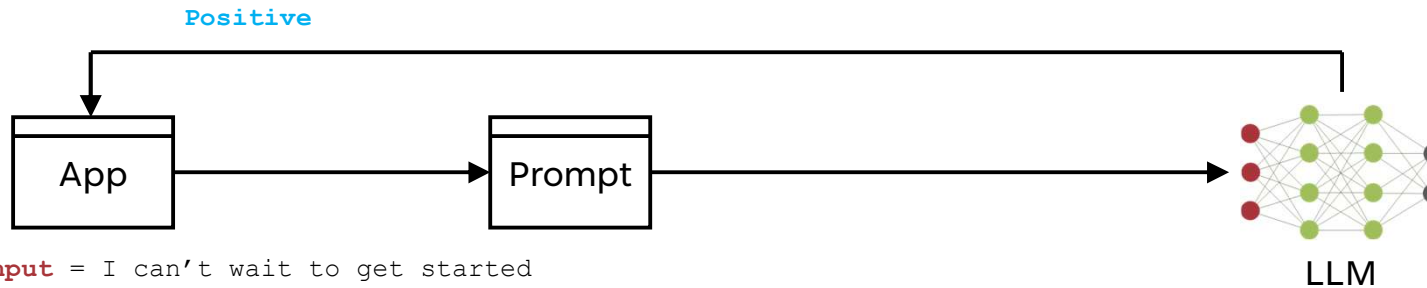




Up Next: Putting AI Models to Work







```

Prompt Template
prompt = f"""
{role}
{context}
{instruct}
"""
  
```

```

Few-Shot
instruct = """
Classify this comment:
Awesome presentation!
Sentiment: Positive

Classify this comment:
What a waste of time!
Sentiment: Negative

Classify this comment:
{input}
Sentiment: """
  
```

```

Role Prompt
role = """You are an event planner
reviewing audience feedback from a
survey taken after a technical
presentation"""
  
```

```

Context
context = """Other relevant info"""
  
```

```

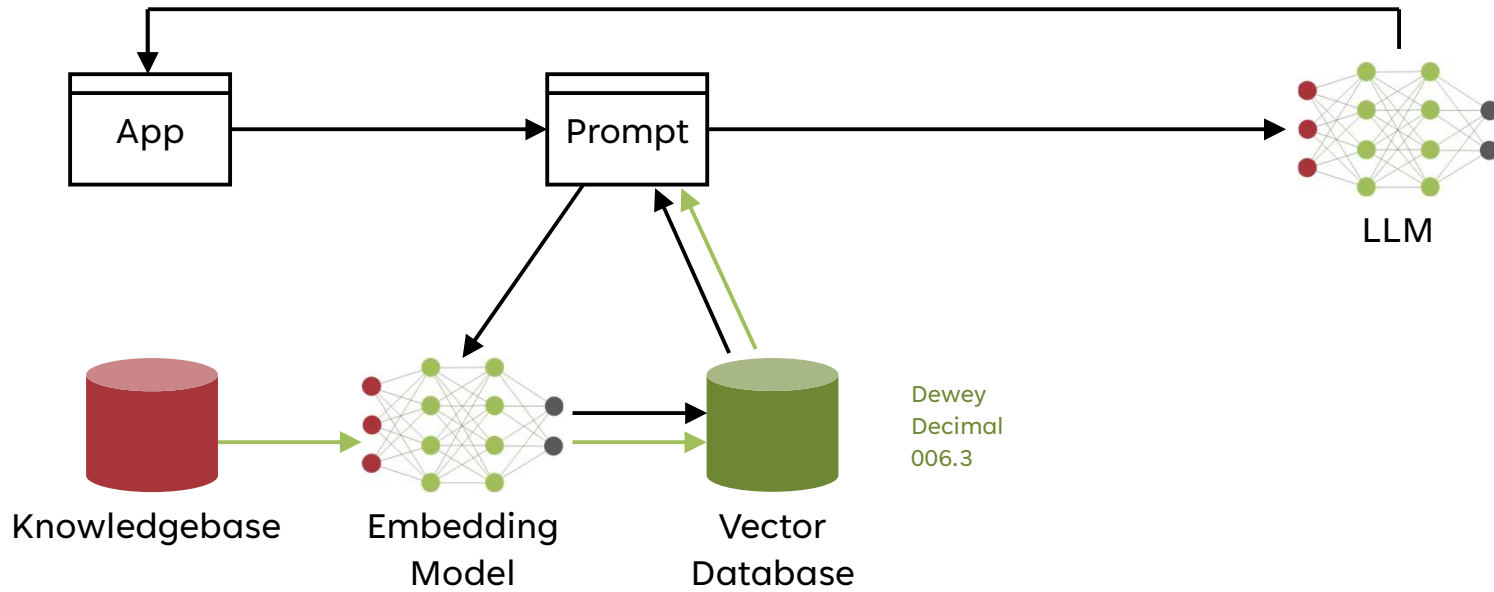
Prompt
You are an event planner
reviewing audience feedback
from a survey taken after a
technical presentation

Other relevant info

Classify this comment:
Awesome presentation!
Sentiment: Positive

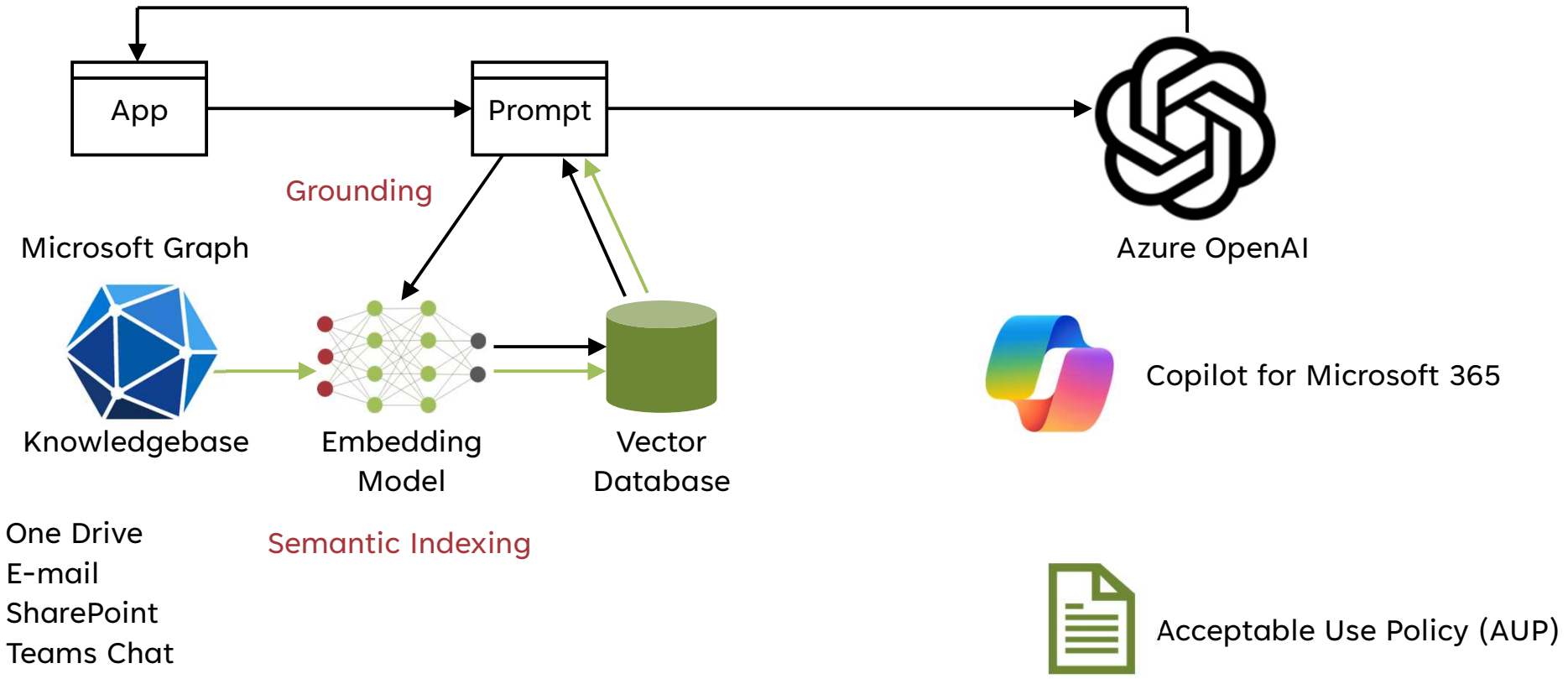
Classify this comment:
What a waste of time!
Sentiment: Negative

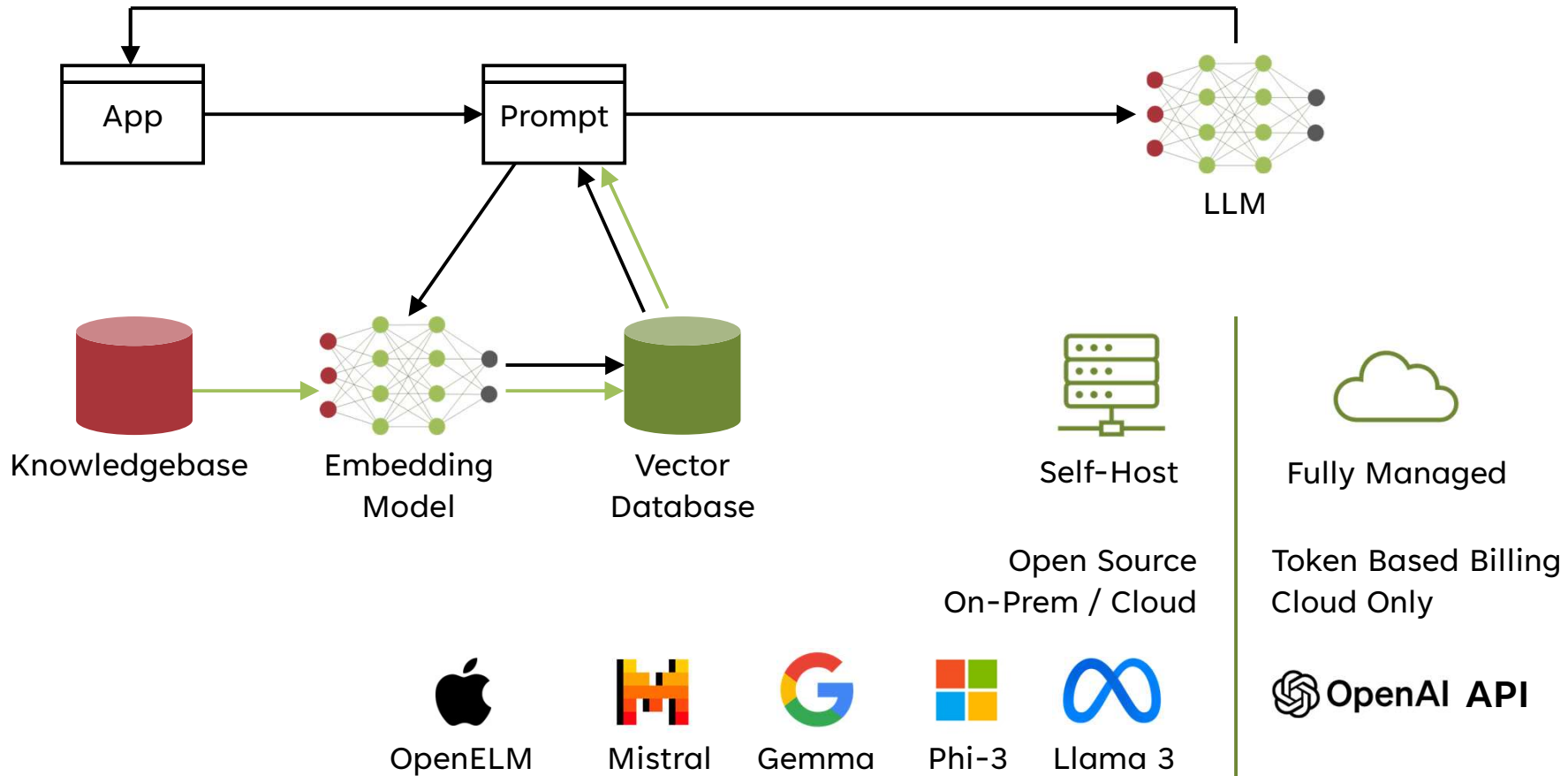
Classify this comment:
I can't wait to get started
Sentiment:
  
```

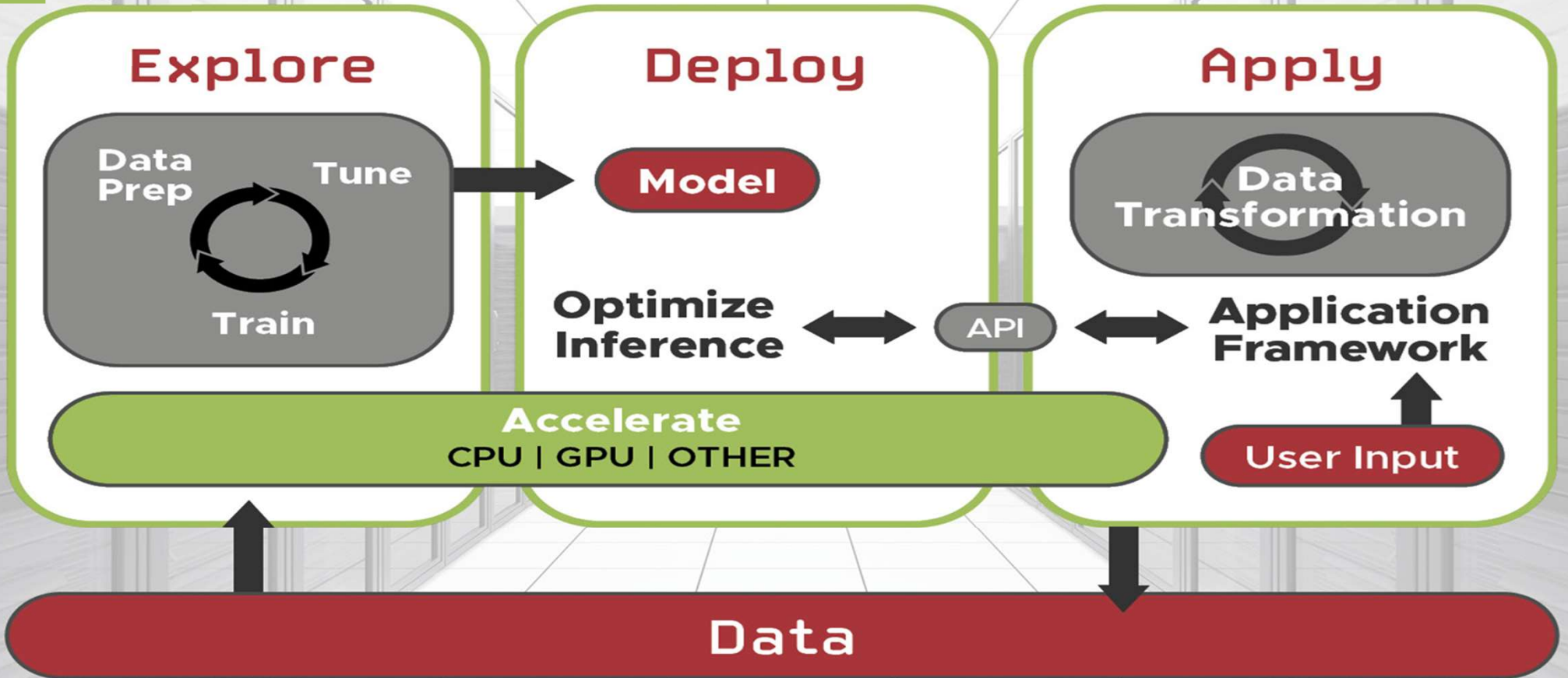


Key Benefits:

- Minimize Hallucinations
- Current Data
- Citing Sources


















Prerequisites:

- Code Editor (VS Code)
- Access to a GPU (Eventually)
- Python Proficiency
- Project Ideas
- Commitment of Time

Python Module	Purpose
Jupyter Notebooks	Experimenting
NumPy	Number Matrix
Pandas	Tabular Data
OpenCV	Visual Data
Flask / Gradio / Streamlit	User Interface

Web Resources:

General	Data Science	Model Deploy	LLM Frameworks
Google  YouTube  Coursera 	colab.google  Kaggle  Fast.ai 	ai.nvidia.com  huggingface.co  Ollama 	LangChain  LlamaIndex 



Practical AI: Definitions, Concepts, and Techniques

Applied AI: Putting AI to Work

Deployed AI: Run & Maintain Inferencing Services

Practical AI:

- AI Taxonomy
- Predictive AI / Generative AI
- Neural Network Overview
- Model Sourcing & Sizing
 - NVIDIA
 - Hugging Face
- LLM Inferencing Techniques
- Getting Started

Applied AI:

- Development Environment
 - VENV / Jupyter
- Model Serving
 - NVIDIA NIM
 - vLLM
- User Interface
 - Gradio / Streamlit
- LLM Tools
 - LangChain / LlamaIndex
- Data Retrieval (RAG)
 - Data Chunking
 - Embedding Model
 - Vector Database

Deployed AI:

- GPU Drivers / Prep
- Bare Metal
- Virtual Machine
 - VMware
 - KVM
- Docker
 - Install
 - GPU Toolkit
 - Basic Commands
- Kubernetes
- Ecosystem Solutions

Video Analytics

- Perimeter Surveillance
- Intrusion Detection
- Vehicle Identification / License Plate Reader
- Manufacturing Quality Control

Large Language Models (LLMs)

- ChatBot Assistant
- Code Generation
- Content & Image Creation
- Sentiment Analysis
- Text Translation and Summarization

Multi-Modal

- Image to Text
- Text to Image
- Speech to Text
- Text to Speech
- Video to Text
- Text to Video



Poll: How Can Nth Help?





Thank You



Nth.com | Security.Nth.com | 800.548.1883 | info@nth.com