



LLM Primer

Nikko Carlo Yabut, MEng AI





Demystifying GenAI: Large Language Models & Retrieval-Augmented Generation

Nikko Carlo Yabut, MEng AI

Certified: AWS Certified Machine Learning - Specialty





Contents

01

Intro

Where are we now in the
Industrial Revolution?

04

Agentic AI

What does the
future hold?

02

AI, ML, DL, QEnAI

Artificial Intelligence,
Machine Learning, Deep Learning,
Generative AI

05

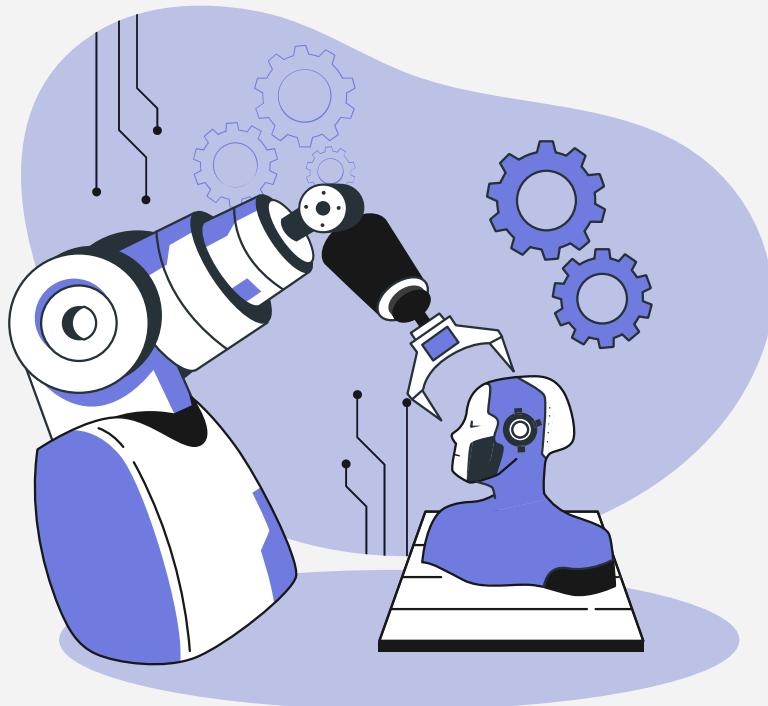
Q&A

03

LLM & RAQs

Large Language Models &
Retrieval-Augmented
Generation

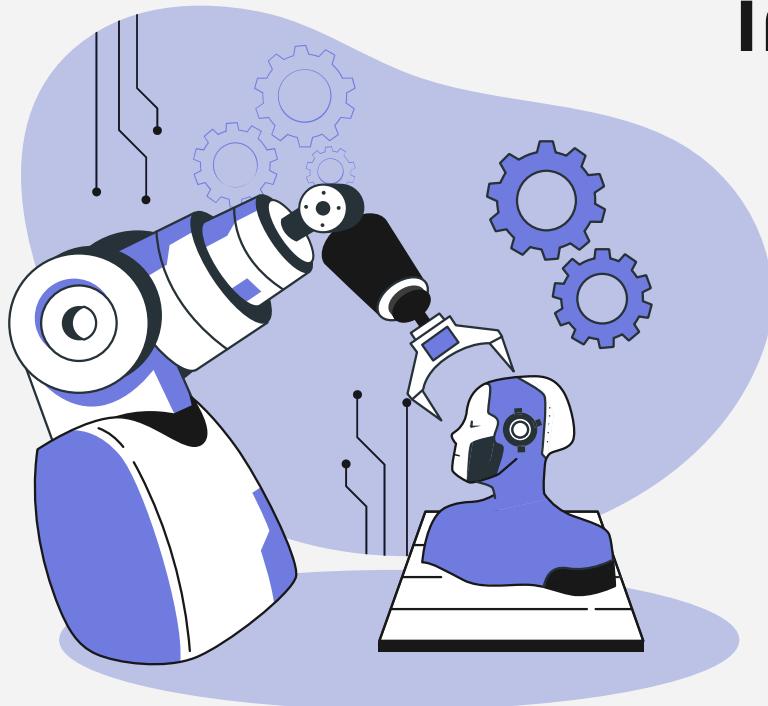




01

Intro

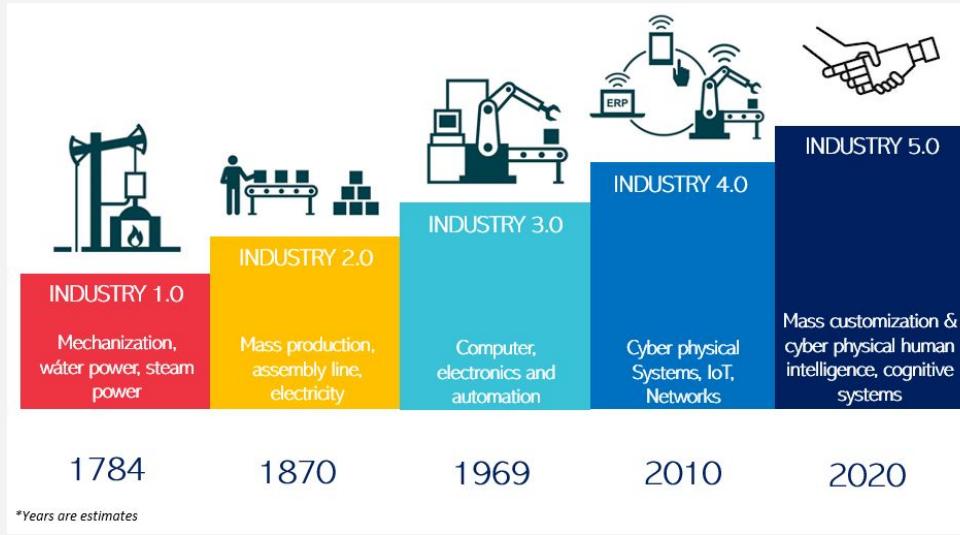




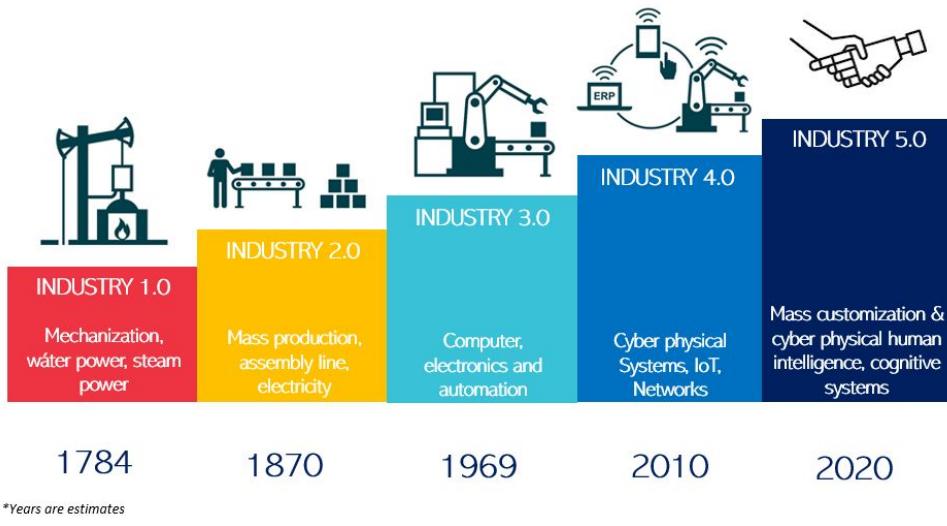
Where are we now in the Industrial Revolution?



Evolution of Industrial Revolution



Evolution of Industrial Revolution

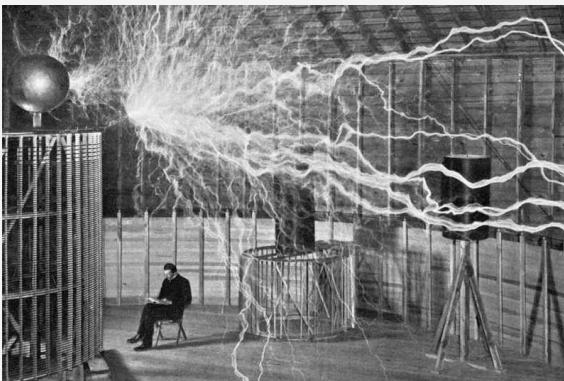
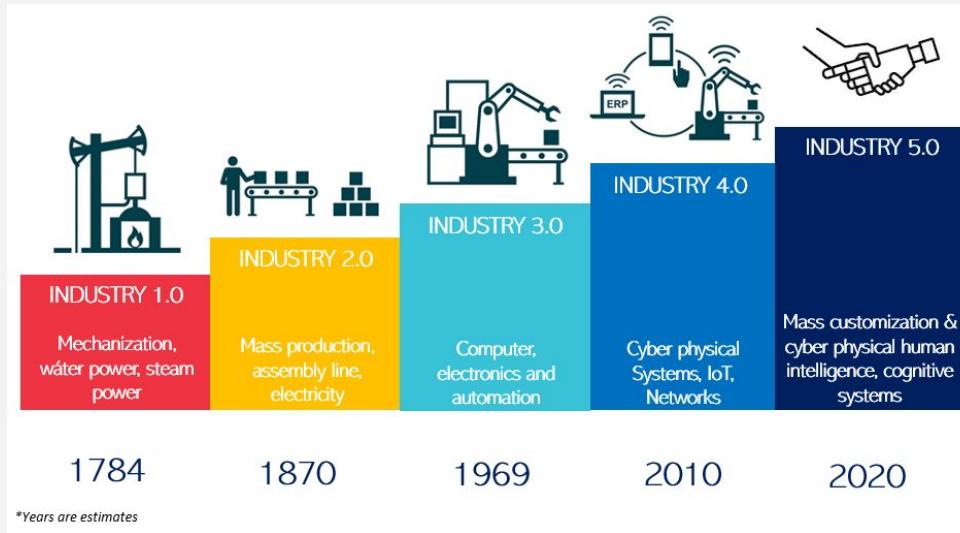


Industry 1.0: The First Industrial Revolution

- **From MANUAL labor to Machines**
- **Key Technology:** Water and steam power.
- **Impact:**
 - Machines transformed manufacturing, starting with textiles in Europe.
 - Fueled by coal, industrial growth expanded significantly.
 - The economy saw a massive boost, and large-scale production became possible.
- **Outcome:** Marked the beginning of large-scale industrialization.



Evolution of Industrial Revolution

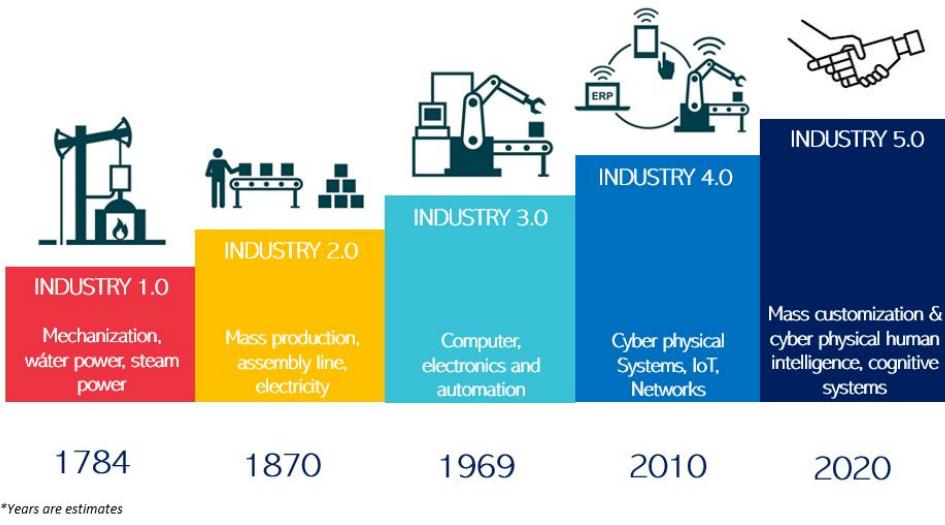


Industry 2.0: The Second Industrial Revolution

- **Key Technology:** Electrification.
- **Impact:**
 - Electricity replaced steam as the primary energy source for machines.
 - Mass distribution of electricity enabled advanced manufacturing.
 - Technological advancements included assembly lines, and publications like *The Principles of Scientific Management* emerged.
 - Led to increased productivity, meeting the needs of customers, employees, shareholders.
- **Outcome:** Often called the **Technology Revolution** due to the scale of innovation.

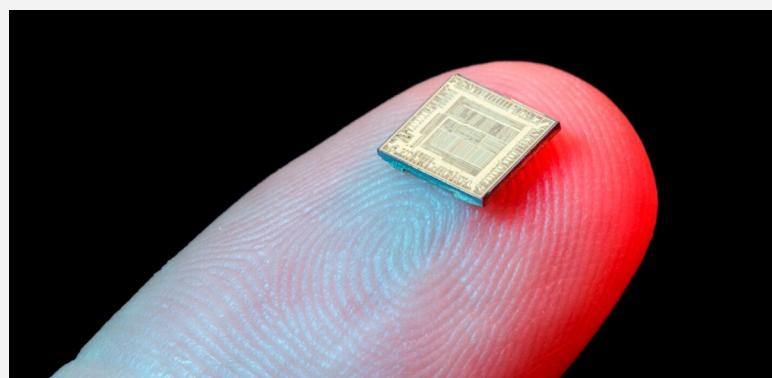


Evolution of Industrial Revolution

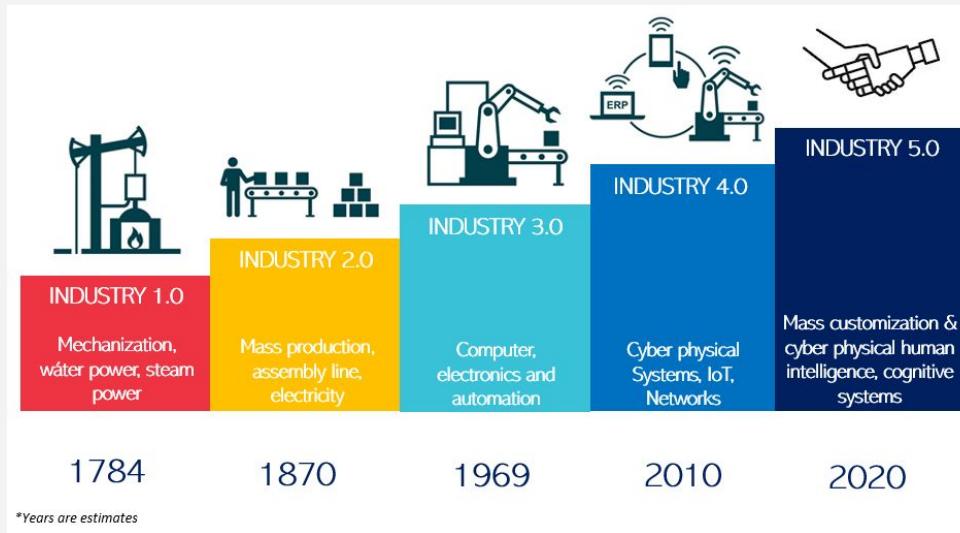


Industry 3.0: The Third Industrial Revolution

- **Key Technology:** Microchips, computers, and automation.
- **Impact:**
 - Introduction of computers facilitated automation in production.
 - Programmable Logic Controllers (PLC) replaced human workers in assembly lines.
- **Outcome:** Known as the **Information Revolution** or **IT Revolution**, due to the rise of computing and digital technologies.



Evolution of Industrial Revolution

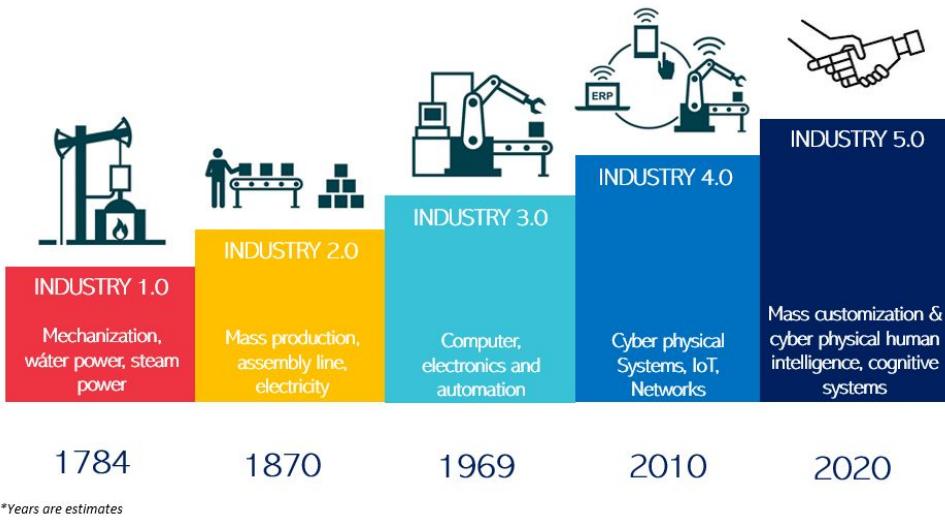


Industry 4.0: The Fourth Industrial Revolution

- **Key Technologies:**
 - Telecommunications (Wireless Comms)
 - Internet
 - Internet of Things (IoT)
 - Cloud Computing
 - Big Data
 - Robotics
 - AI
- **Impact:**
 - Fully automated production systems where machines communicate, control, and make decisions **without human intervention**.
- **Outcome:** Known as the **Digital Revolution**, reshaping industries and leading to higher levels of efficiency.



Evolution of Industrial Revolution

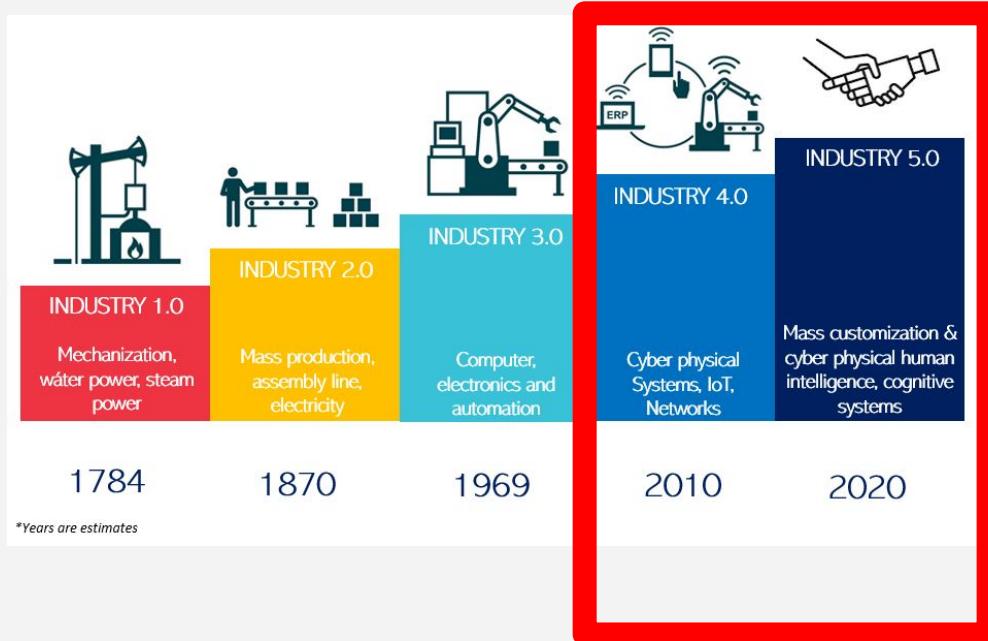


Industry 5.0: The Fifth Industrial Revolution

- **Time Period:** Emerging (futuristic vision).
- **Key Focus:** Collaboration between human and AI.
- **Impact:**
 - Emphasizes **co-bots** (collaborative robots) working alongside humans.
 - Enhanced efficiency through **human-robot collaboration on the factory floor**.
 - Inspired by Japan's **Society 5.0** concept, where humans and AI deliver business value together.
- **Outcome:** A more human-centric approach.



Evolution of Industrial Revolution



Where are we now in the Industrial Revolution?

Industry 4.0
(Digital Revolution)

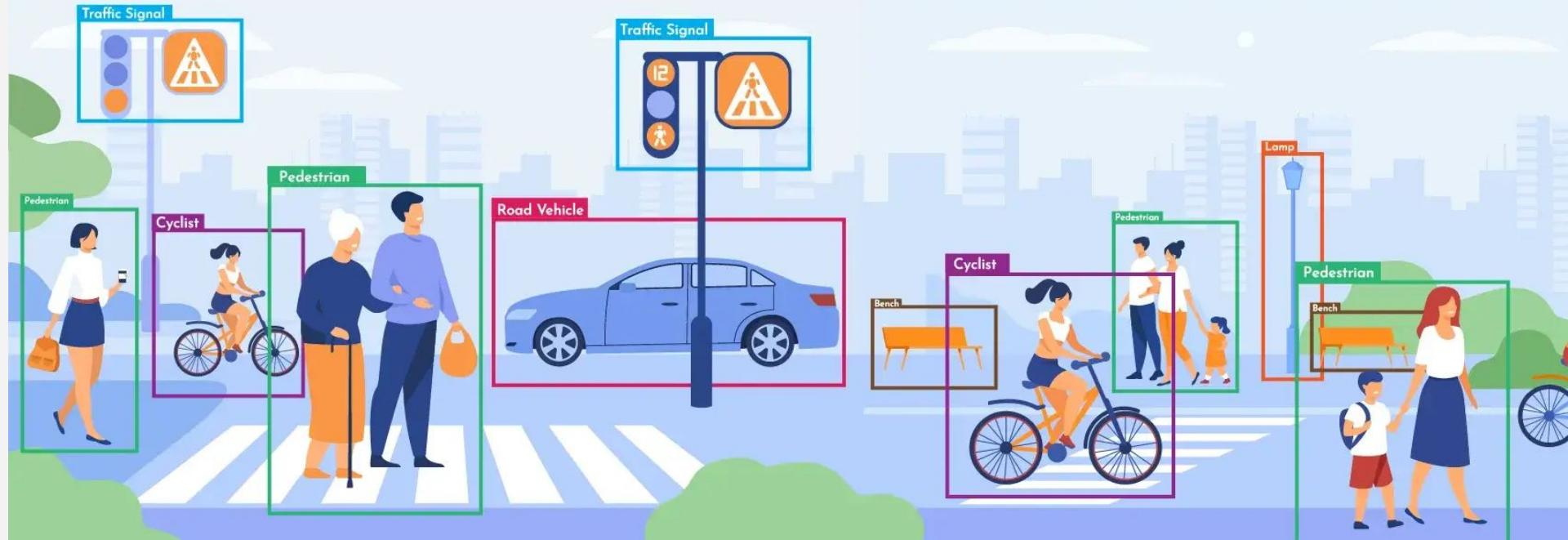
&
Industry 5.0
(Human-AI Collaboration)



Evolution of Industrial Revolution



In this discussion, we will focus on a specific use case of Industry 5.0: LLMs



NLP

Part-of-Speech (POS) Tagging

Emotion & Intent Detection

Named Entity Recognition (NER)

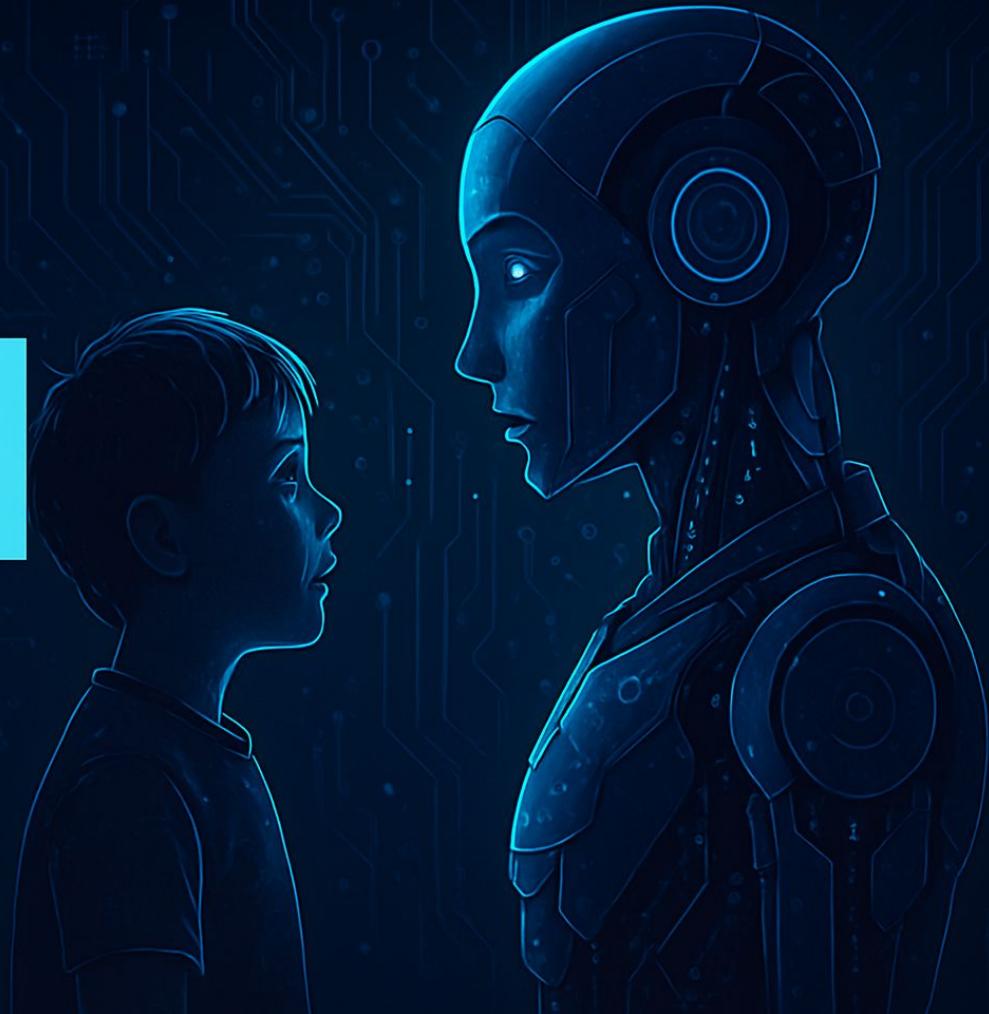
Text Classification

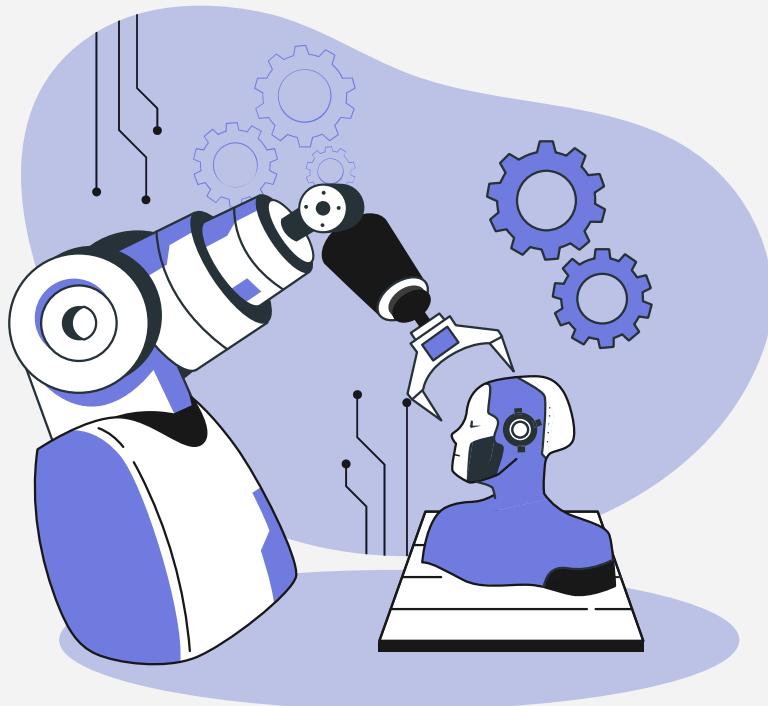
Information Extraction

Text Similarity

LLM

LLM





02

AI, ML, DL, GenAI



Intelligence

A very general mental capability that among other things involves the ability to:

- Reason
- Plan
- Solve problems
- Think abstractly
- Comprehend complex ideas
- Learn from experience

Journal of Intelligence 1997 Vol 24 No 1



Intelligence

A very general mental capability that among other things involves the ability to:

- Reason
- Plan
- Solve problems
- Think abstractly
- Comprehend complex ideas
- Learn from experience

Journal of Intelligence 1997 Vol 24 No 1

Artificial Intelligence (AI)

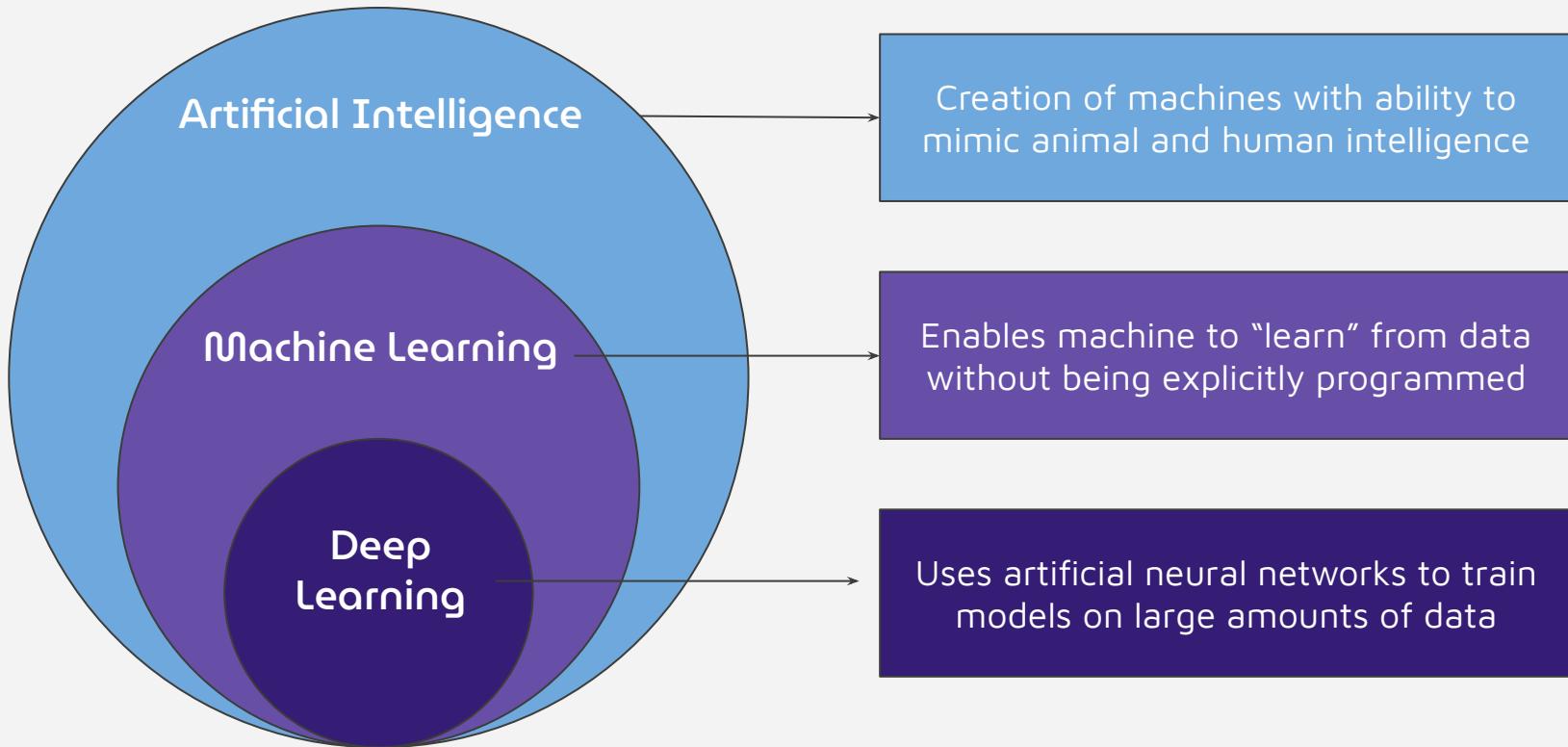
A machine has AI capabilities
IF it exhibits **animal** or **human** intelligence



Ants bridging algo



AI, ML, DL



AI, ML, DL, GenAI

Artificial Intelligence

Is the field of study that deals with creation of machine that exhibits animal or human intelligence.

Machine Learning

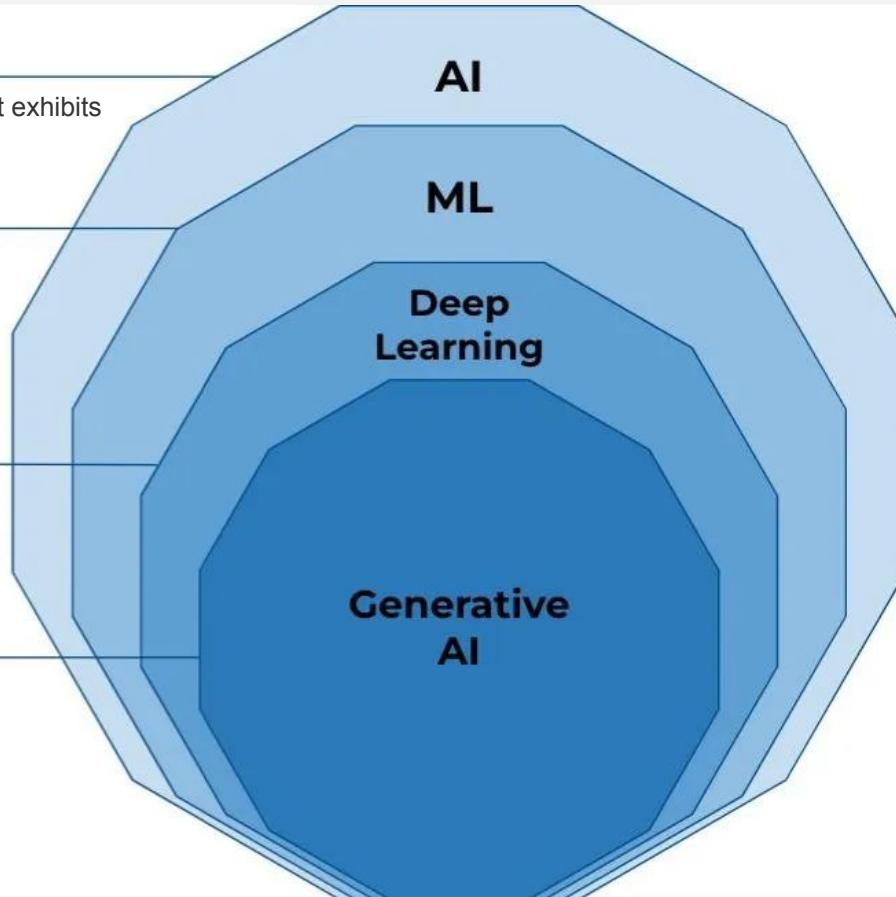
Is a branch of AI that focus on the creation of intelligent machines that learn from data.
Another very well known branch inside AI is **Optimization**.

Deep Learning

Is a subset of Machine Learning methods, based on **Artificial Neural Networks**.
Examples: CNNs, RNNs

Generative AI

A type of ANNs that generate data that is similar to the data it was trained on.
Examples: GANs, LLMs



AI, ML, DL, GenAI

Artificial Intelligence

Is the field of study that deals with creation of machine that exhibits animal or human intelligence.

Machine Learning

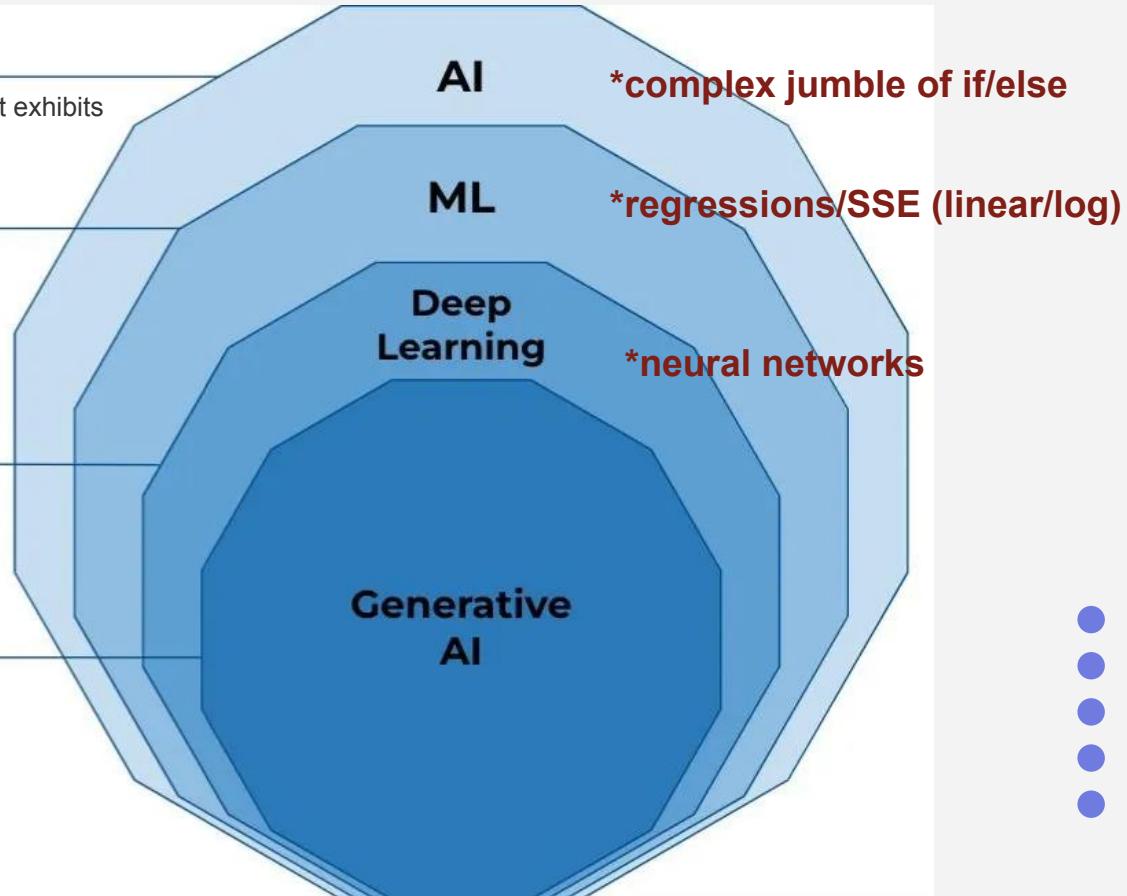
Is a branch of AI that focus on the creation of intelligent machines that learn from data. Another very well known branch inside AI is **Optimization**.

Deep Learning

Is a subset of Machine Learning methods, based on **Artificial Neural Networks**. Examples: CNNs, RNNs

Generative AI

A type of ANNs that generate data that is similar to the data it was trained on. Examples: GANs, LLMs



AI, ML, DL, GenAI

The Nobel Prize in Physics 2024



Ill. Niklas Elmehed © Nobel Prize Outreach

John J. Hopfield

Prize share: 1/2

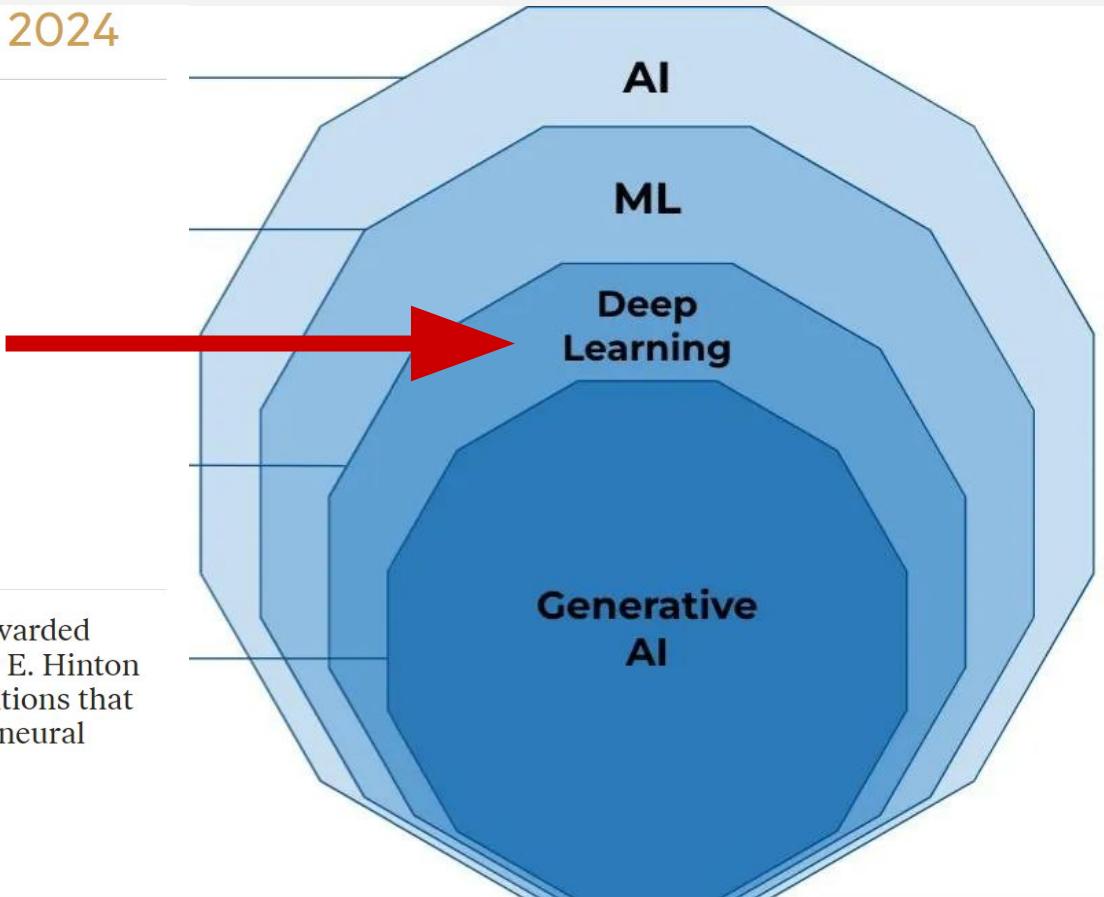


Ill. Niklas Elmehed © Nobel Prize Outreach

Geoffrey Hinton

Prize share: 1/2

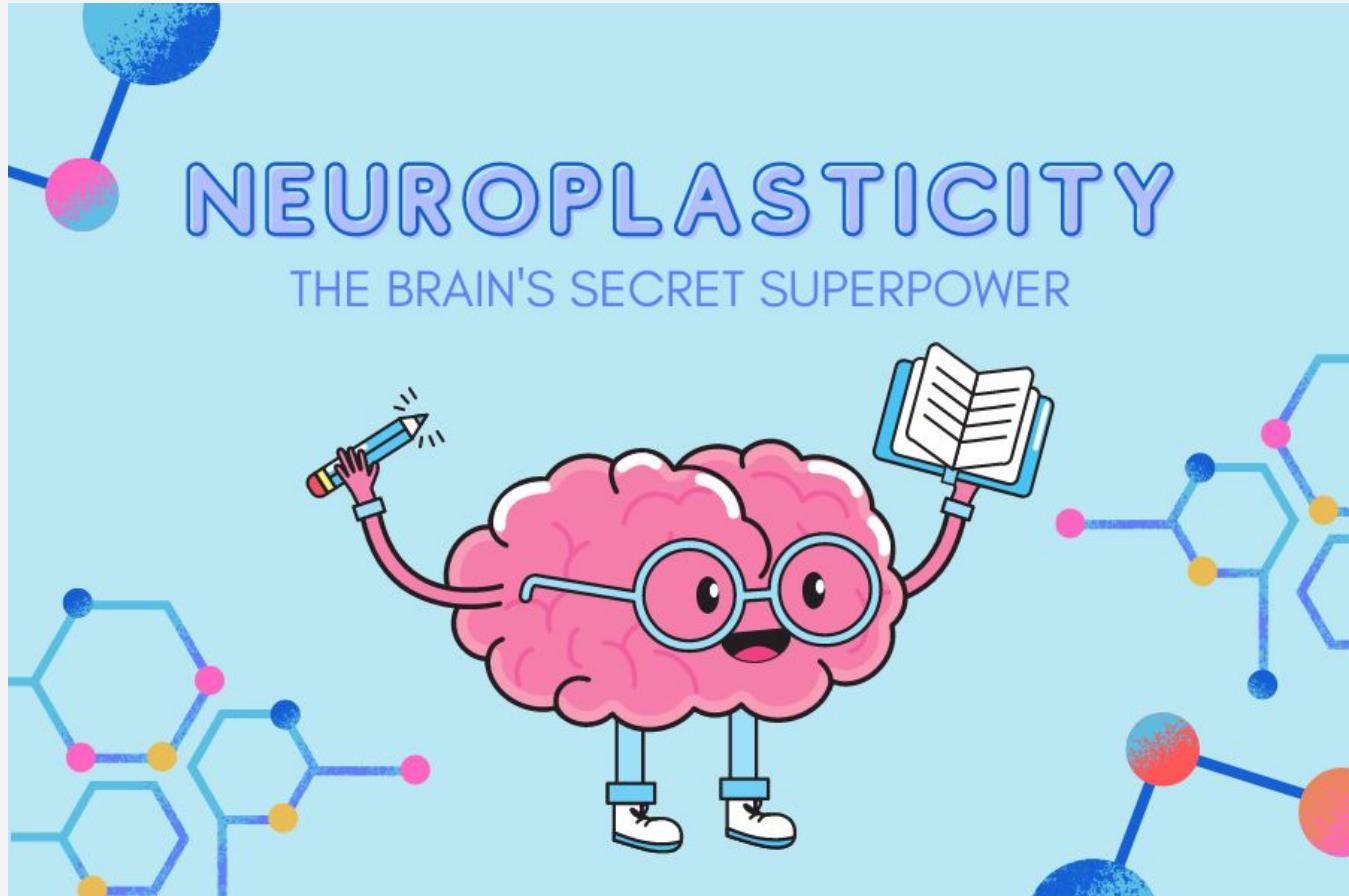
The Nobel Prize in Physics 2024 was awarded jointly to John J. Hopfield and Geoffrey E. Hinton "for foundational discoveries and inventions that enable machine learning with artificial neural networks"



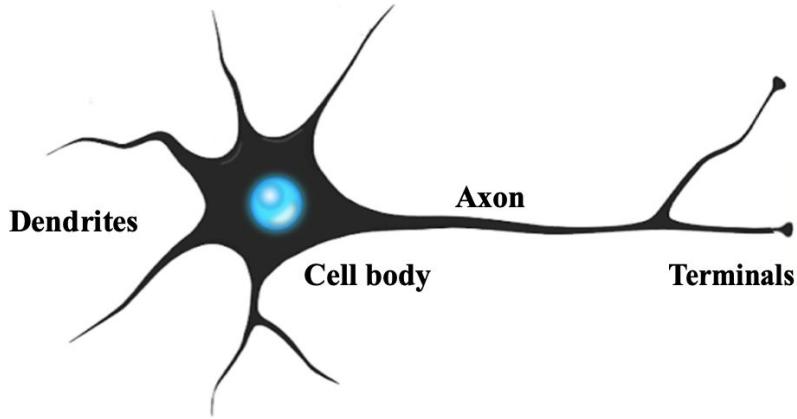
How do machines ‘learn’?



From biological neuron to Neural Network neuron



From biological neuron to Neural Network neuron



Biological Neurons: Structure and Function

Dendrites

Receive signals from other neurons, starting the communication process.

Soma (Cell Body)

Processes incoming signals to decide the neuron's response.

Axon

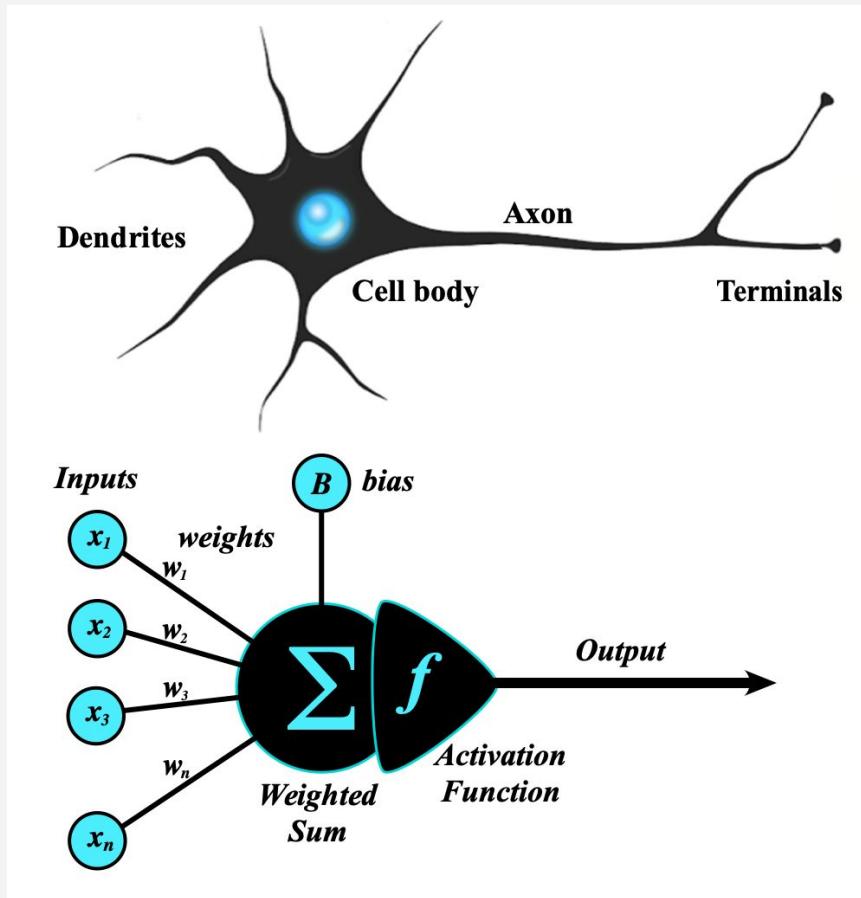
Transmits the processed signal outward to connected neurons.

Synapses

Junction points where signals pass to neighboring neurons.



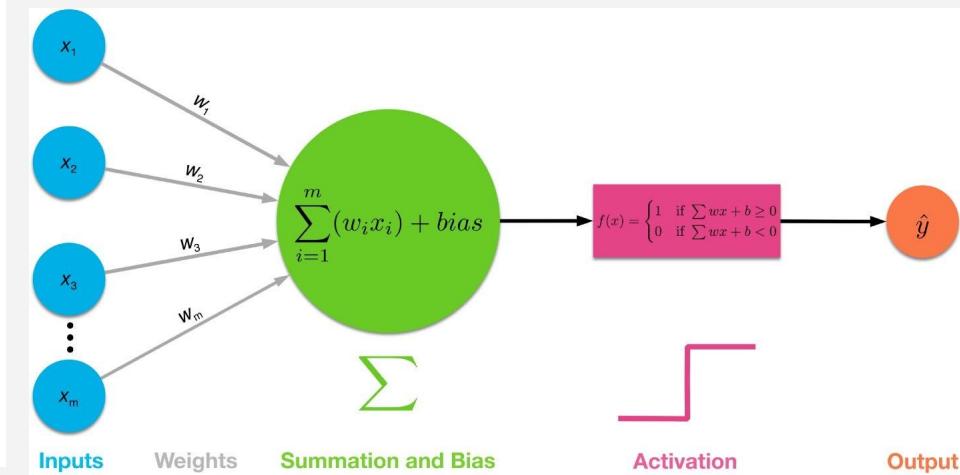
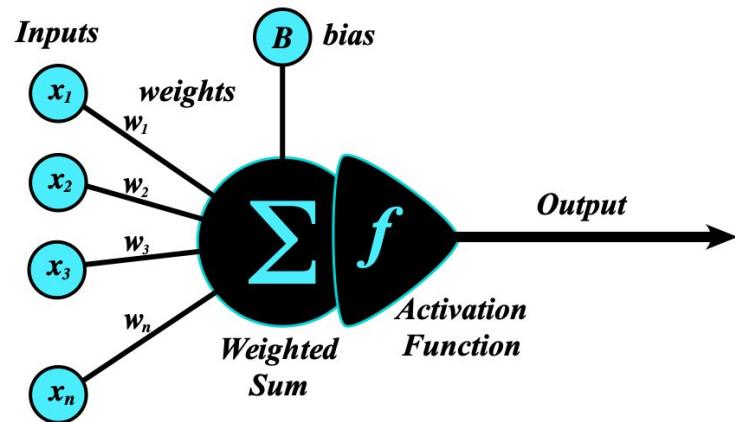
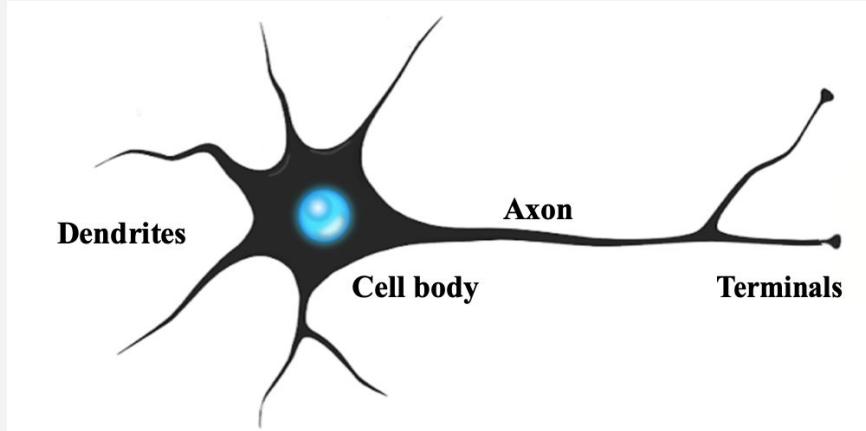
From biological neuron to Neural Network neuron



Comparison of Biological Neurons and AI Neurons

- **Structure:**
 - **Biological Neuron:**
 - Composed of dendrites, a cell body (soma), and an axon.
 - Complex connections and synapses with other neurons.
 - **AI Neuron:**
 - Simplified mathematical model or node in a neural network.
 - Represents inputs, weights, activation functions, and outputs.
- **Learning Mechanism:**
 - **Biological Neuron:**
 - Learns through neuroplasticity / synaptic plasticity (strengthening/weakening of connections).
 - Involves complex biological processes like long-term potentiation.
 - **AI Neuron:**
 - Learns through adjustments of weights via algorithms (e.g., backpropagation).
 - Utilizes large datasets to optimize performance during training.

From biological neuron to Neural Network neuron



Linear regression is essentially a special case of a **one-neuron network**.

Linear Regression

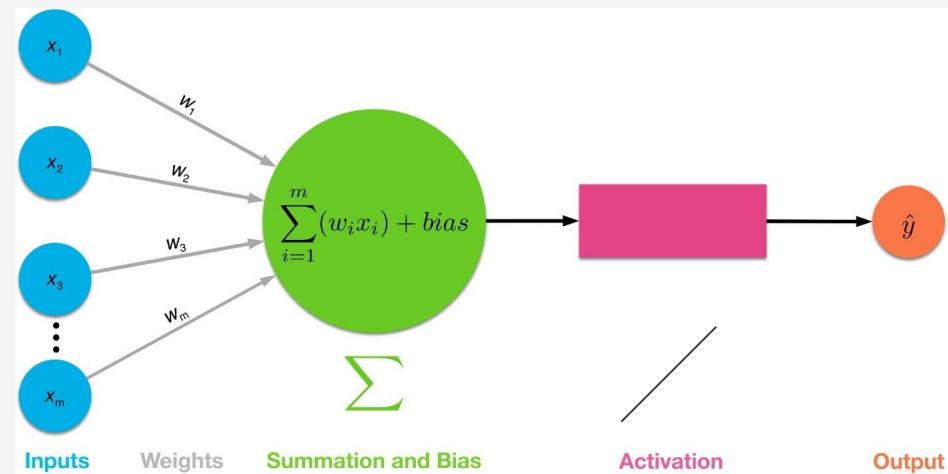
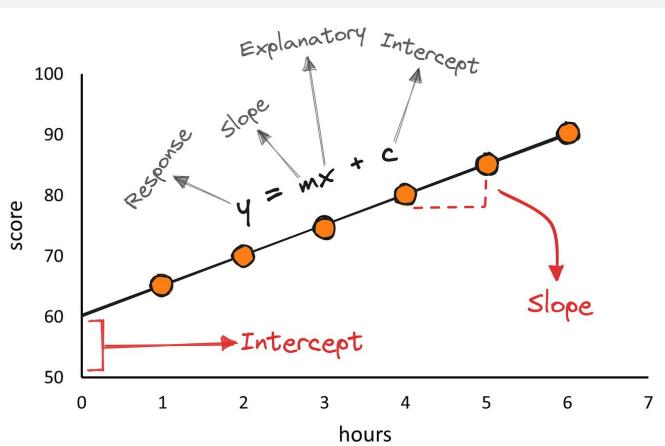
Examples:

1. Housing Price Prediction

- **X: Features:** Square footage, number of bedrooms, location, age of the house.
- **y: Target:** House price in dollars.

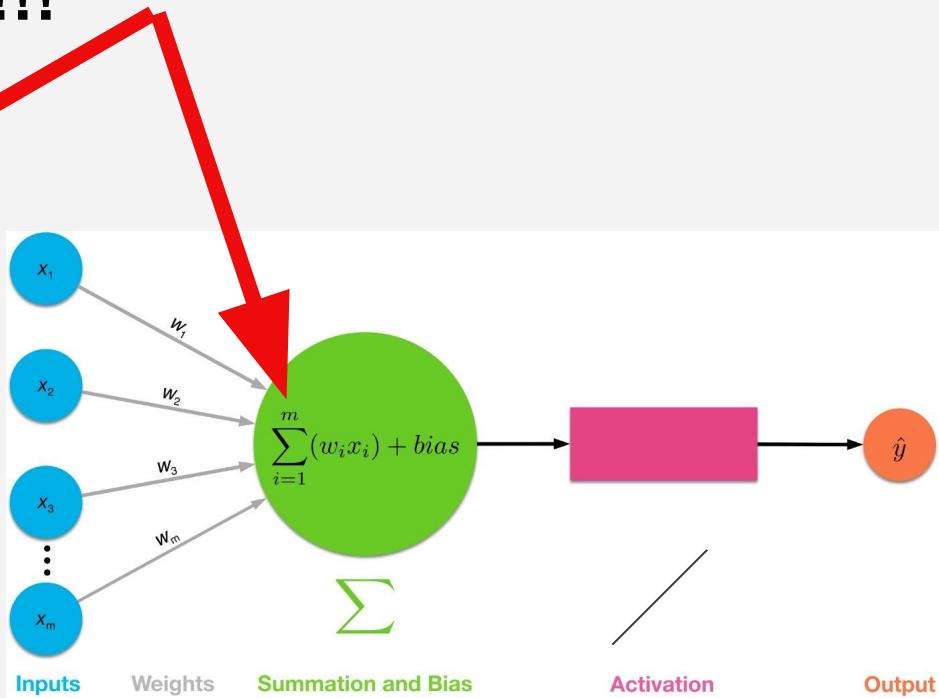
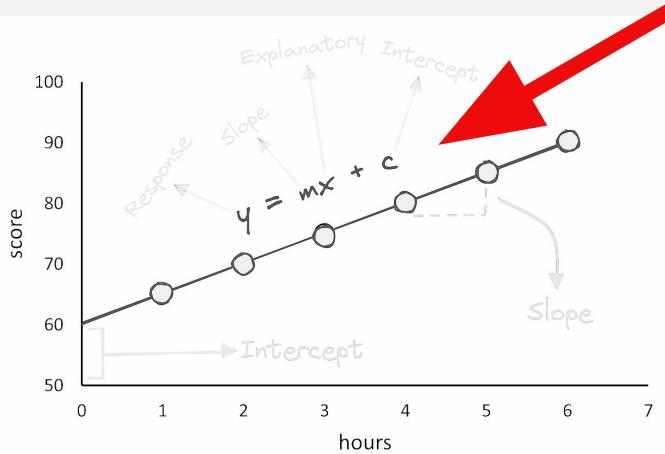
2. Sales Forecasting

- **X: Features:** Advertising spend, season (month), store location.
- **y: Target:** Sales in units or revenue.



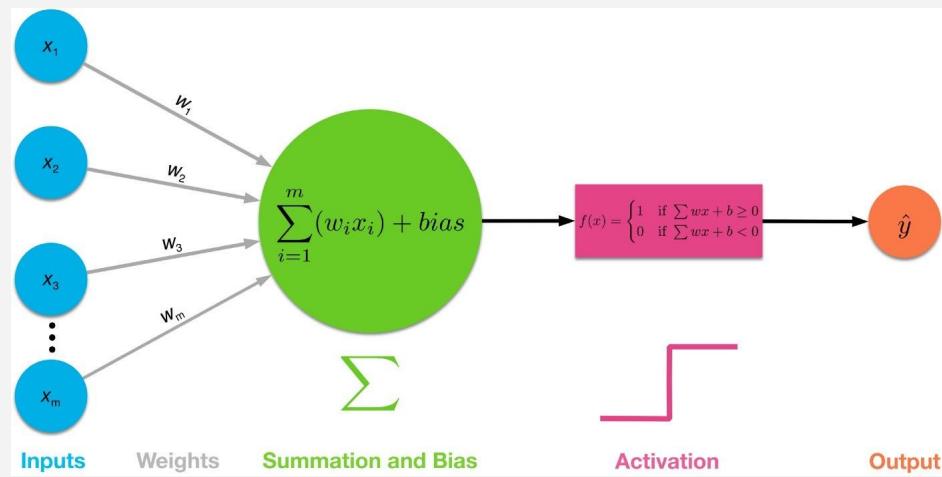
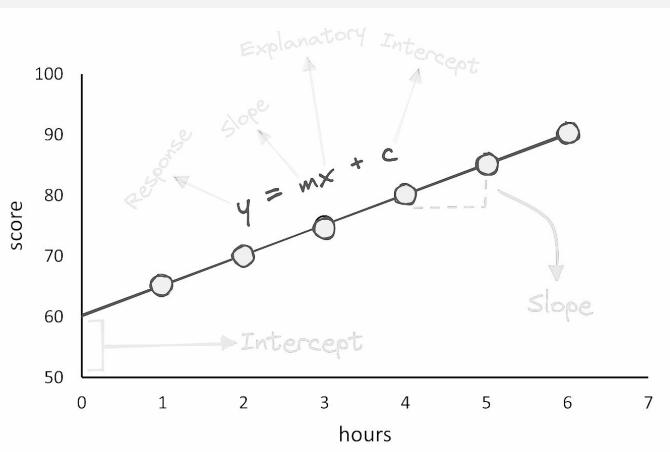
Linear regression is essentially a special case of a **one-neuron network**.

Similar Mathematical Formula!!!



However, the real world is nonlinear!

To truly capture the complexity of the real world, we need networks of to go beyond linearity.

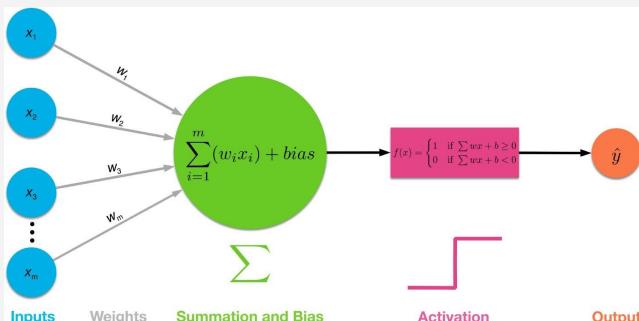


Linear regression is essentially a special case of a **one-neuron network**.

Why One Neuron Isn't Enough

Single neuron limits
Can only solve simple, linear problems.

Complex patterns need more
One neuron can't recognize intricate shapes or boundaries.



Aspect	Linear Regression	One-Neuron Network
Formula	$\hat{y} = \mathbf{w}^\top \mathbf{x} + b$	$y = \phi(\mathbf{w}^\top \mathbf{x} + b)$
Activation Function	None (or identity)	Typically nonlinear (e.g. ReLU, sigmoid)
Training Loss	MSE	Depends (MSE, cross-entropy, etc.)
Use Case	Regression	Regression or classification

1. Formula for Linear Regression

The **linear regression** model predicts a continuous outcome based on one or more input features:

$$\hat{y} = w_1 x_1 + w_2 x_2 + \cdots + w_n x_n + b$$

Or more compactly using vector notation:

$$\hat{y} = \mathbf{w}^\top \mathbf{x} + b$$

Where:

- \hat{y} is the predicted output,
- $\mathbf{x} = [x_1, x_2, \dots, x_n]^\top$ is the input vector,
- $\mathbf{w} = [w_1, w_2, \dots, w_n]^\top$ is the weight vector,
- b is the bias term (intercept).

2. Formula for a One-Neuron Network (Perceptron without activation)

A **single neuron** in a neural network (without activation) also computes a weighted sum of its inputs:

$$y = \mathbf{w}^\top \mathbf{x} + b$$

If it includes a **nonlinear activation function** ϕ , the formula becomes:

$$y = \phi(\mathbf{w}^\top \mathbf{x} + b)$$

3. Relationship Between the Two

Linear regression is essentially a **special case** of a one-neuron network:

- It uses **no activation function** (or a linear activation: $\phi(x) = x$).
- It's trained using **least squares** loss.
- The goal is to minimize the **mean squared error** (MSE) between predictions and actual values.

Introducing MLPs: From one neuron to network of neurons

MLP (Multi-Layer Perceptron)?

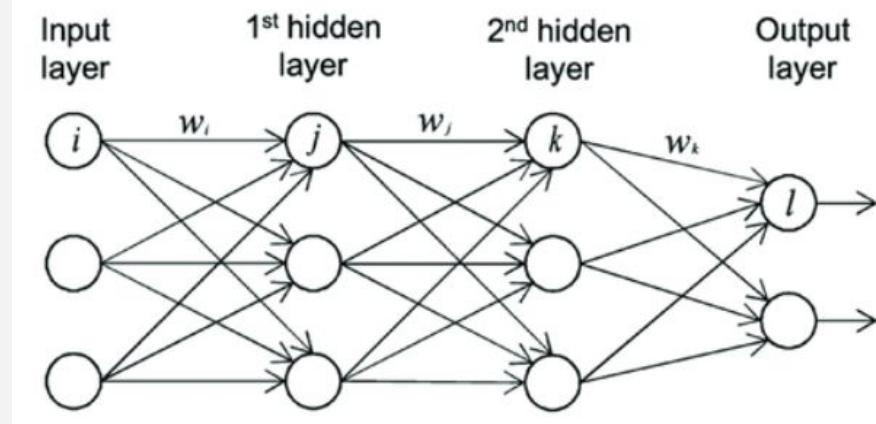
A **Multi-Layer Perceptron (MLP)** is a type of neural network that **combines many neurons** arranged in **layers** to solve more complex problems than a single neuron can handle.

Just like people working together can solve bigger problems than one person alone, **MLPs stack neurons** into:

- **Input Layer** – receives the data
- **Hidden Layers** – extract patterns and features
- **Output Layer** – produces the final decision or prediction

Each neuron in a layer is connected to **all** neurons in the next layer — this is called a **fully connected network**.

An MLP allows us to go from **simple decisions (single neuron)** to **complex pattern recognition**, by letting multiple neurons **collaborate** and **refine** their outputs through layers.



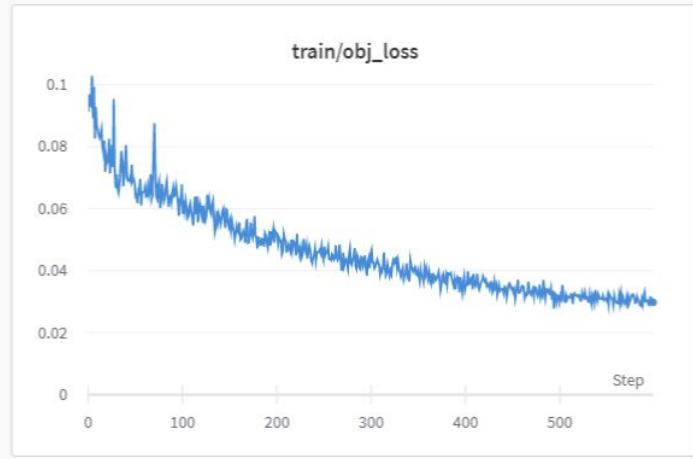
We've built the brain.
Now the big question is...

can it learn?



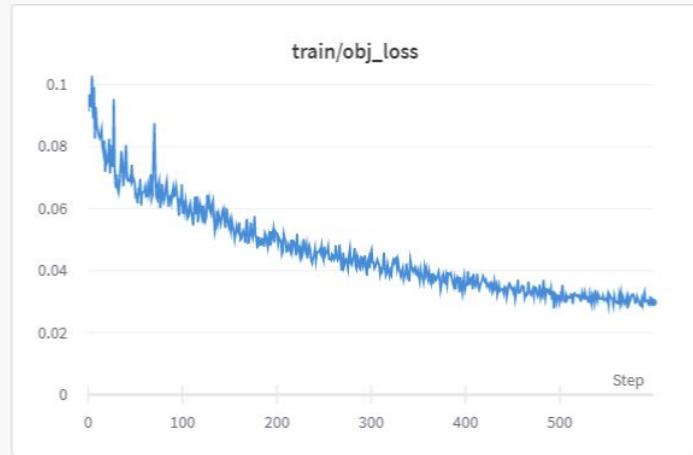
How do machines ‘learn’?

“Learning from mistakes...”



How do machines ‘learn’?

“Learning from mistakes...”

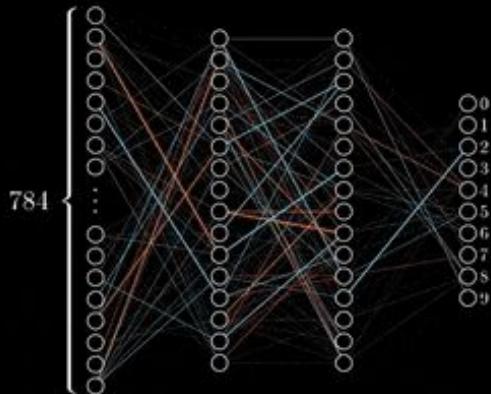


A computer program is said to **learn** from **experience E** with respect to some class of **tasks T**, and **performance measure P**, if its performance at tasks in T, as measured by P, improves with experience E. – Tom Mitchell



How does machines ‘learn’? BACKPROPAGATION

Training in progress...



A computer program is said to **learn** from **experience E** with respect to some class of **tasks T**, and **performance measure P**, if its performance at tasks in T, as measured by P, improves with experience E. – Tom Mitchell

What is Backpropagation?

Backpropagation is how a neural network **learns** from its mistakes.

It works by:

1. **Comparing** the model’s prediction to the correct answer (using a **loss function**).
2. **Measuring the error**, then

Sending that error backward through the network to **adjust the weights**—so it does better next time.

💡 Analogy:

Like a student reviewing a wrong answer, figuring out *why* it was wrong, and correcting their approach.



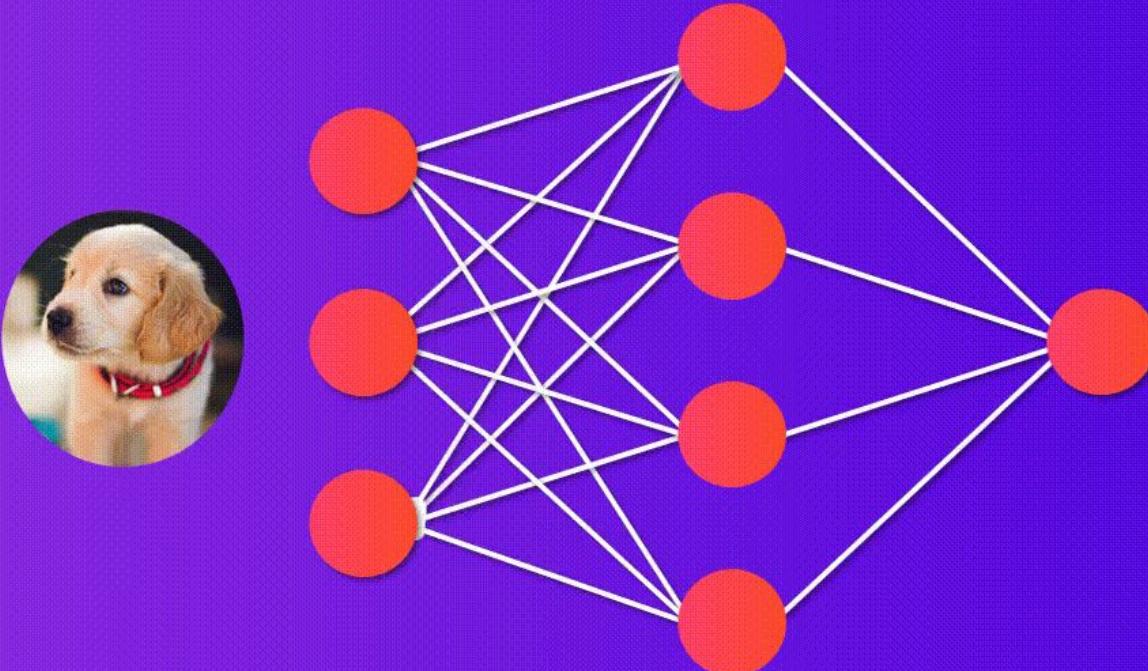
How does machines ‘learn’?

“Learning from mistakes...”

A computer program is said to **learn** from
experience E with respect to some class of
tasks T, and **performance measure P**, if its
performance at tasks in T, as measured by P,
improves with experience E. – Tom Mitchell



Neural Network doing Prediction



Types of Learning Problems

Supervised Learning

These are images of dogs.



These are images of cars.



Now, what is this an image of?



Unsupervised Learning

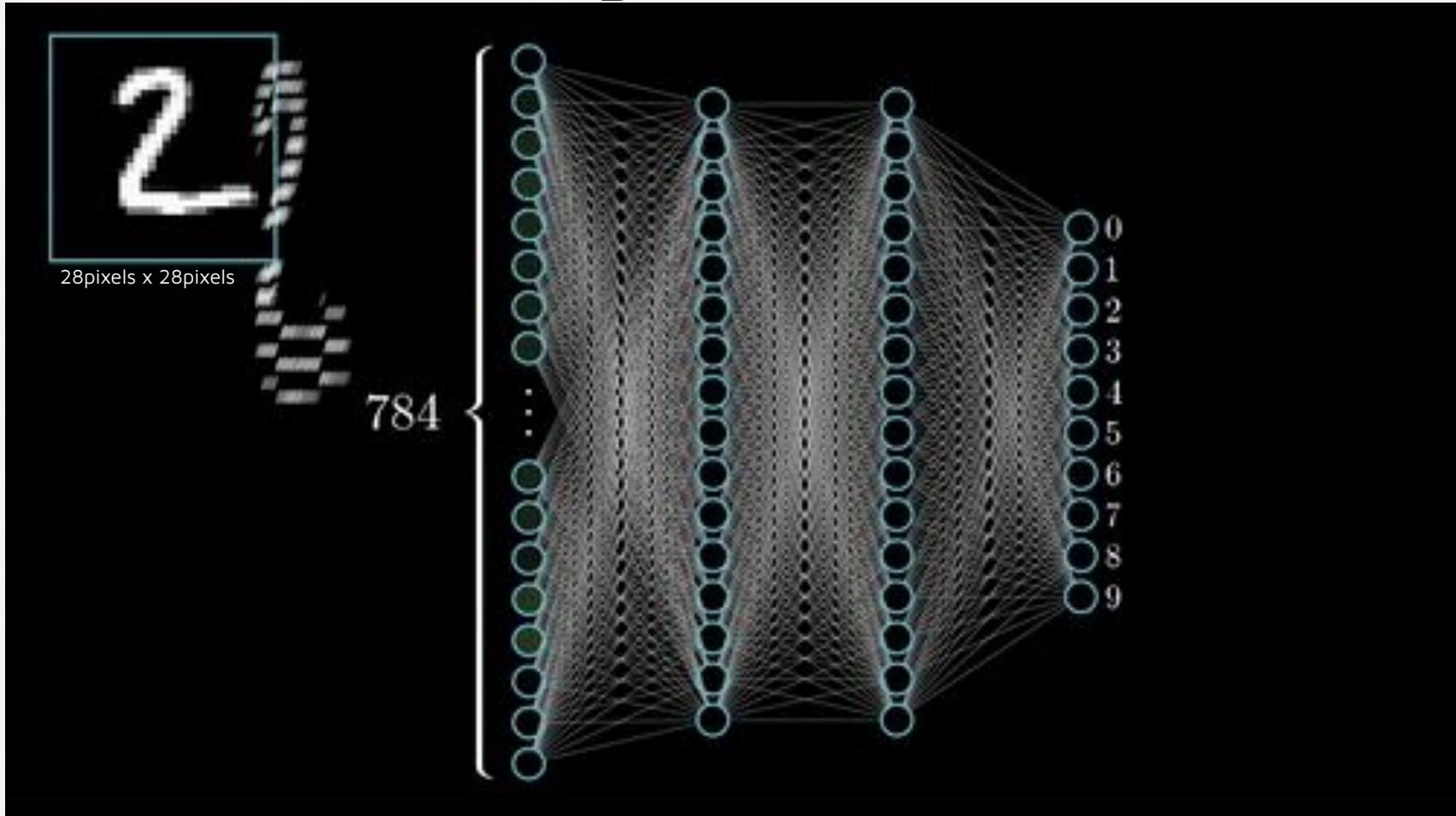
Here are some images...



Is there an image that does not belong?

Are there images with similar patterns?

Neural Network doing Prediction

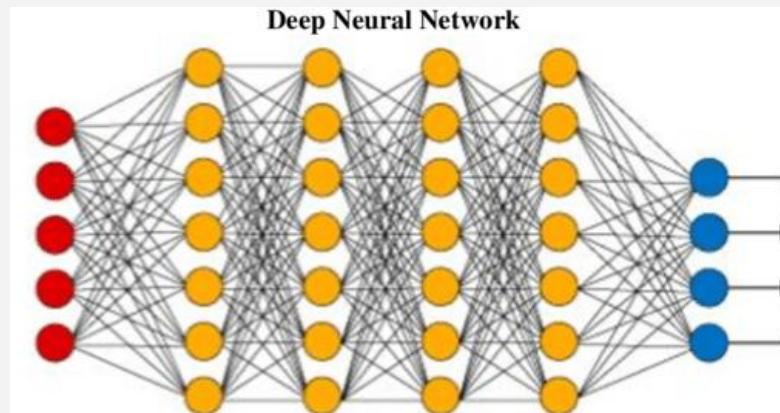


How many trainable parameters is in ChatGPT 3.5?



ChatGPT

ChatGPT 3.5 was based on the GPT 3.5 engine, which received training on **over 175 billion** training parameters. ChatGPT 4 took this training one step further, and is purportedly trained on over a trillion parameters.



AI, ML, DL, GenAI

Artificial Intelligence

Is the field of study that deals with creation of machine that exhibits animal or human intelligence.

Machine Learning

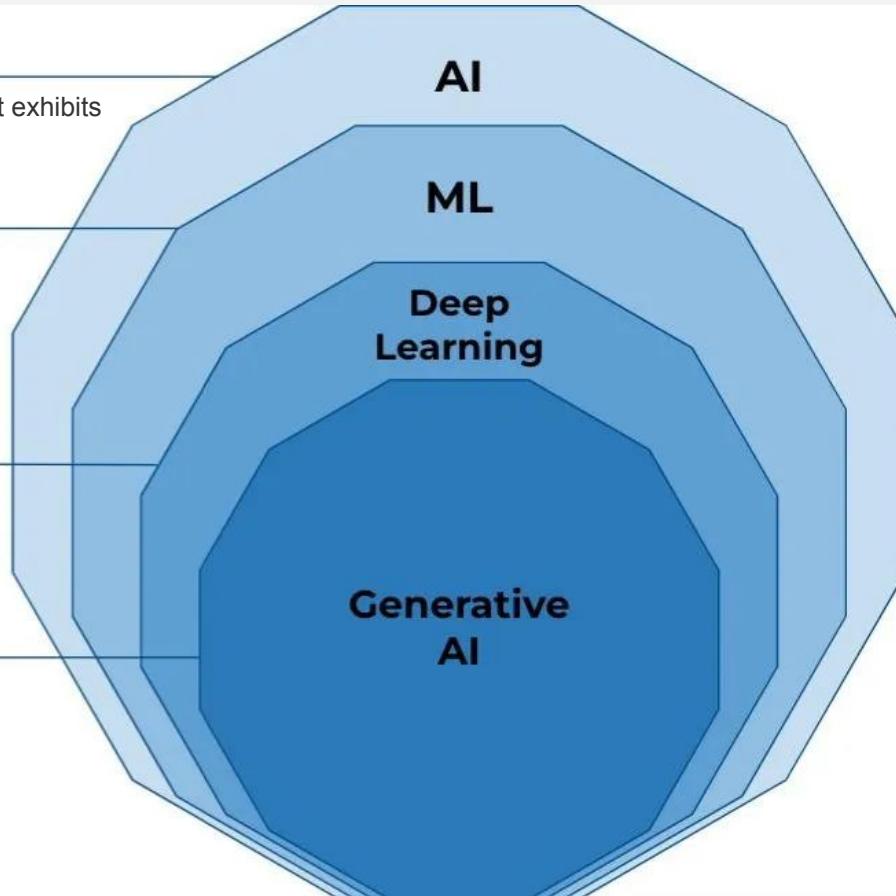
Is a branch of AI that focus on the creation of intelligent machines that learn from data.
Another very well known branch inside AI is **Optimization**.

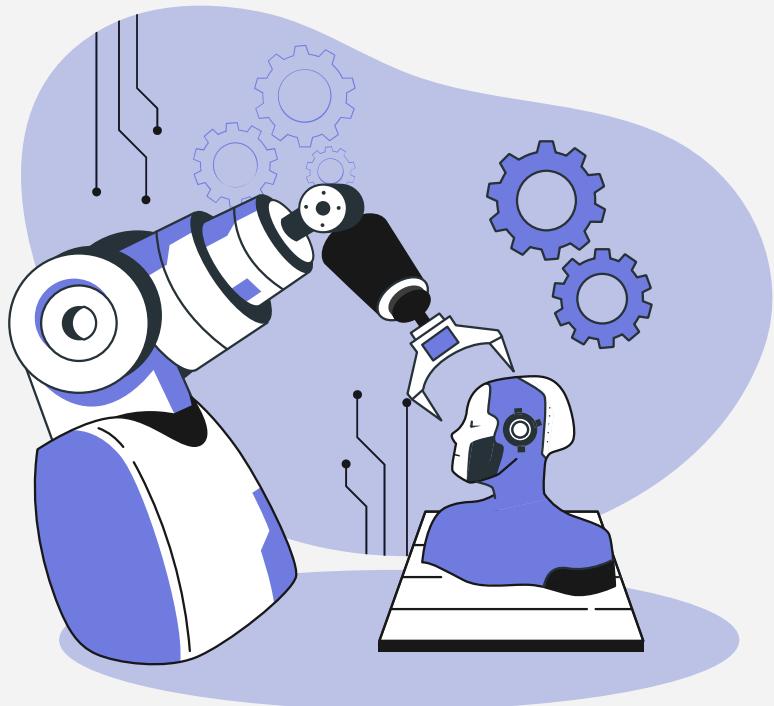
Deep Learning

Is a subset of Machine Learning methods, based on **Artificial Neural Networks**.
Examples: CNNs, RNNs

Generative AI

A type of ANNs that generate data that is similar to the data it was trained on.
Examples: GANs, LLMs



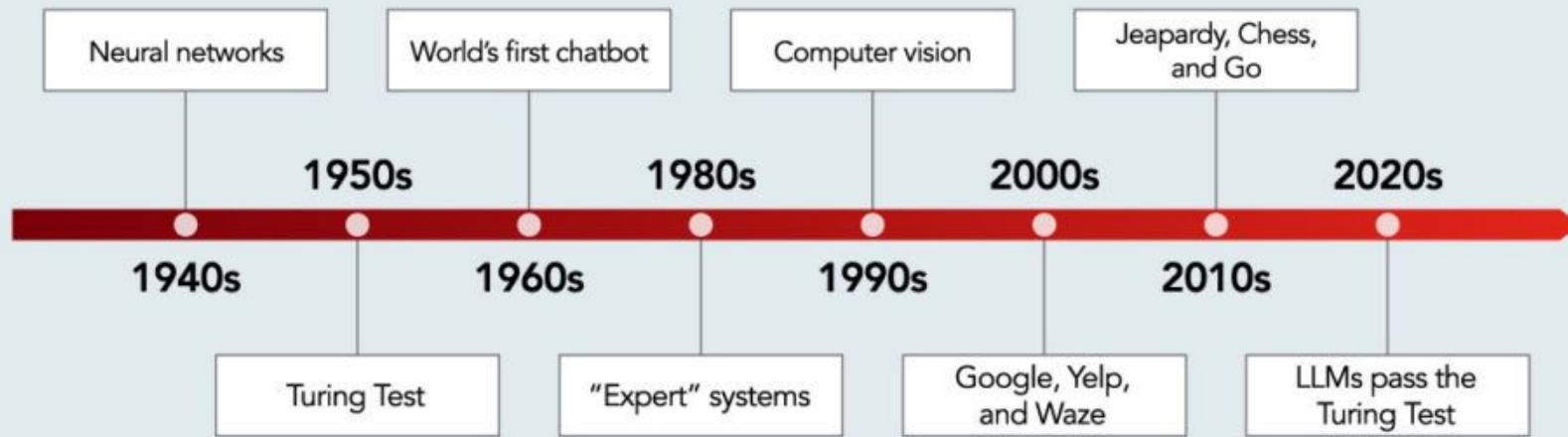


03

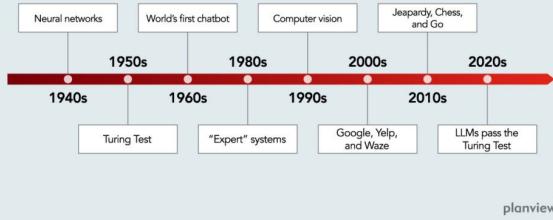
LLM



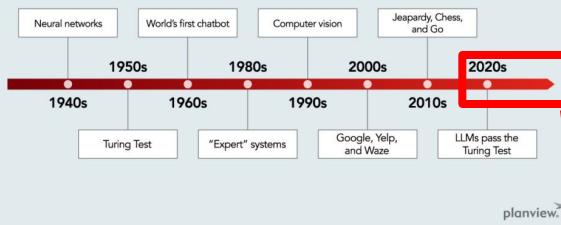
AI took 80+ years to become an overnight sensation



AI took 80+ years to become an overnight sensation



AI took 80+ years to become an overnight sensation

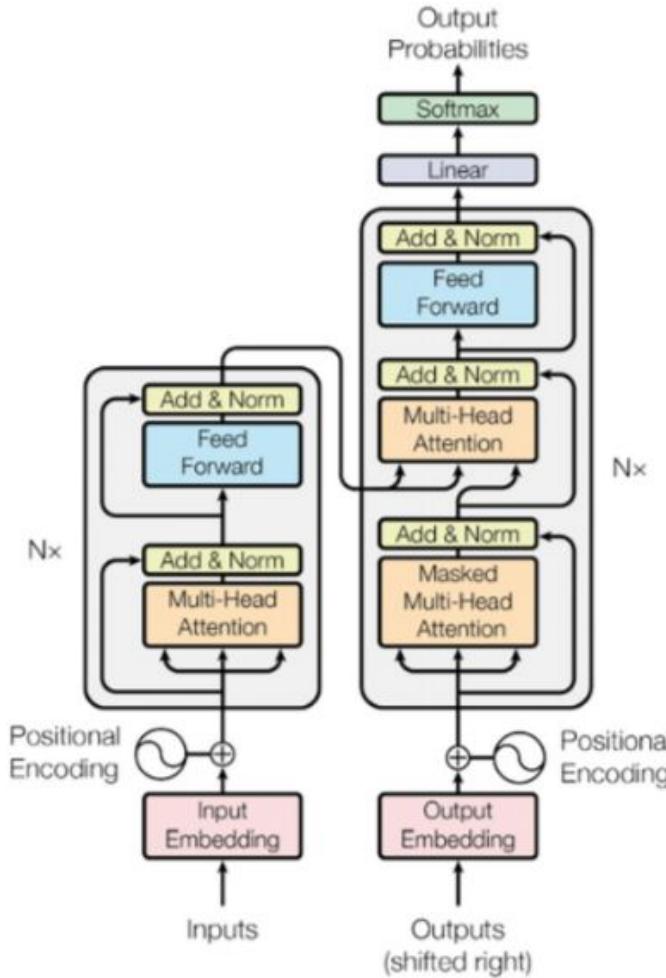


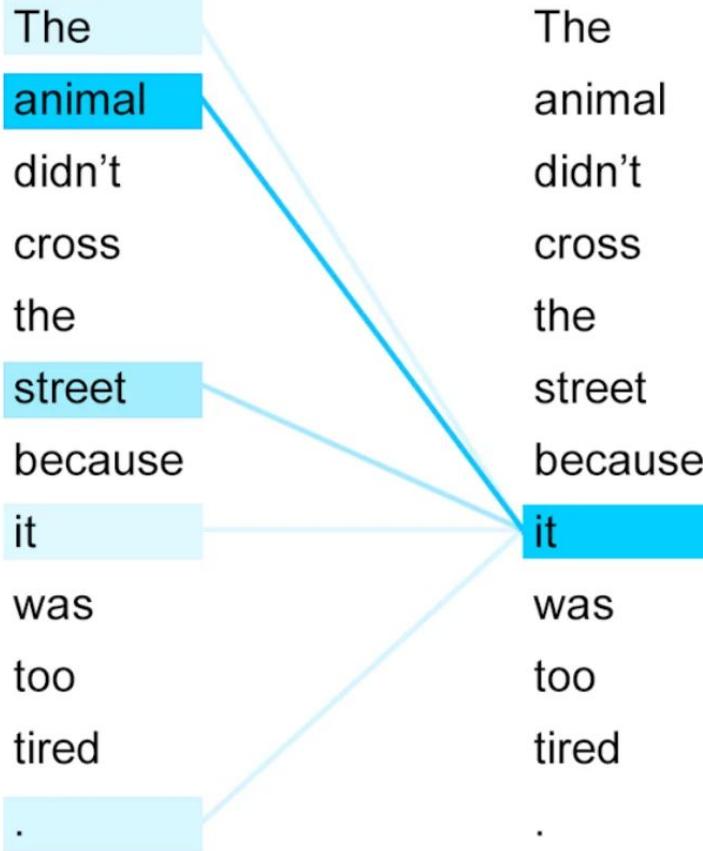
History of ChatGPT



Transformer

Attention Is All You Need





In transformers, **self-attention** allows the model to **weigh the importance of every word** in a sentence **relative to others** when making predictions.

In the sentence:

"The animal didn't cross the street because it was too tired."
— the word "it" is ambiguous, but most likely refers to "the animal" rather than "the street."

Using self-attention, the transformer can **focus more heavily on** "animal" when processing "it," because words like "tired" are more semantically associated with living things. The model learns these patterns from large text corpora, enabling it to resolve such references based on context.

This is milestone in NLP because it:

- Improves comprehension
- Supports human-like understanding

*previously, we're just counting words (e.g. what's the most frequent words, n-gram)



Why are LLMs an overnight sensation?

- Language is a representation of human knowledge
- Products of language is easy to generate (e.g. Internet, books, publications, messages, etc)
- Language is the best modality to train AI models to mimic human intelligence
- In the past year, LLMs have been successful in solving real-world problems
- **Best to think of LLMs as the compression with context of the digital human knowledge**



Large Language Models

Definition: LLMs are advanced AI models trained on vast amounts of text data to understand and generate human language.

Scale and Complexity: LLMs consist of billions to trillions of parameters, allowing them to handle complex language tasks like translation, summarization, and conversation.

Learning Process: They learn patterns, context, and semantics from large datasets to generate coherent responses.

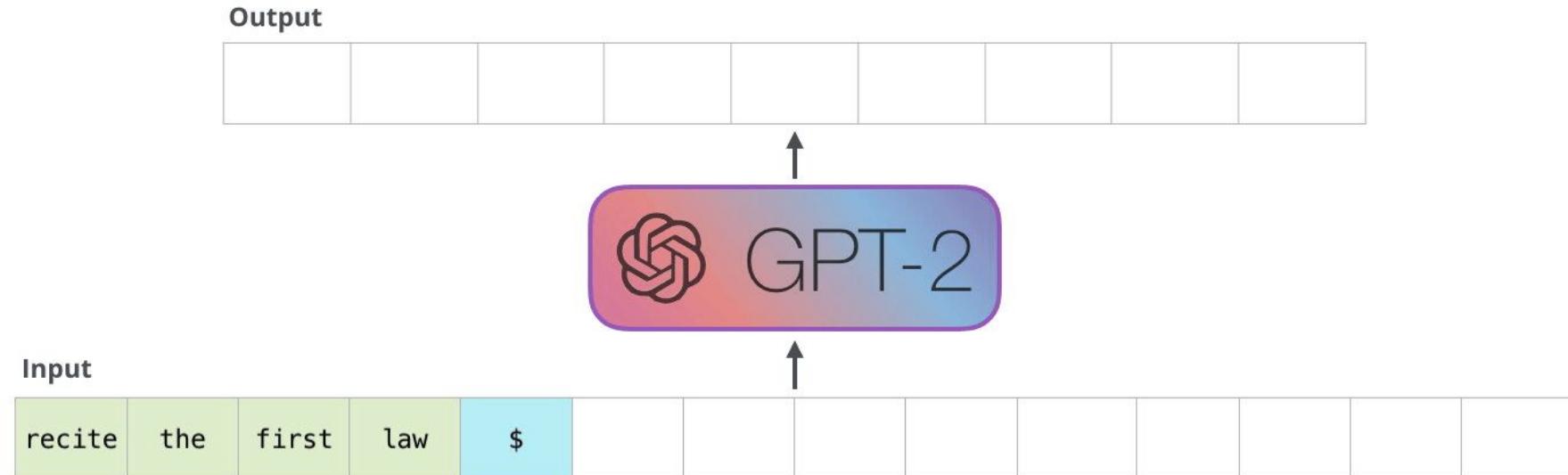
Real-World Applications: LLMs power chatbots, virtual assistants, content generation tools, code-writing systems, and more across industries.

Challenges: Despite their capabilities, LLMs face challenges like bias in training data, high computational costs, and sometimes generating inaccurate or nonsensical responses.



How does LLM work?

- Use of autoregression



How does LLM work?

- Use of autoregression to predict the next word (token)

$$p(x_1, x_2, \dots, x_L) = p(x_1)p(x_2|x_1) \cdots p(x_L|x_1, x_2 \cdots x_{L-1}) = \prod_{i=1}^L p(x_i|x_{1:i-1})$$

$p(\text{the}, \text{cat}, \text{sat}, \text{on}, \text{the}, \text{mat}) =$

$p(\text{the}) \cdot$

$p(\text{cat}|\text{the}) \cdot$

$p(\text{sat}|\text{the}, \text{cat}) \cdot$

$p(\text{on}|\text{the}, \text{cat}, \text{sat}) \cdot$

$p(\text{the}|\text{the}, \text{cat}, \text{sat}, \text{on}) \cdot$

$p(\text{mat}|\text{the}, \text{cat}, \text{sat}, \text{on}, \text{the})$



Capabilities

- Question Answering
- Word Analogies
- Summarization
- Language Translation
- In-Context Learning

The capital of Malaysia is _____



Capabilities

- Question Answering
- **Word Analogies**
- Summarization
- Language Translation
- In-Context Learning

Cat:Kitten,Dog: _____



Capabilities

- Question Answering
- Word Analogies
- **Summarization**
- Language Translation
- In-Context Learning

The Olympic Games is a major international multi-sport event. During its history, both the Summer and Winter Games were a subject of many scandals, controversies, and boycotts. Some Olympic Games were not conducted due to major wars - such as World War I and World War II. The Games originated in Ancient Greece, and now are conducted every four years by the International Olympic Committee. Athletes from more than 200 nations participate in these games. The symbol of the Olympic Games - five interlocking colored rings, representing the continents of the world united by Olympism. The goal of the Olympic Movement is to contribute to building a peaceful and better world by educating youth through sport practiced without discrimination of any kind and in the Olympic spirit.

Summary

The Olympic Games, originating in Ancient Greece, are a major international event occurring every four years, involving over 200 nations. The event, overseen by the International Olympic Committee, has faced scandals, controversies, and were cancelled due to world wars. The five interlocked rings symbolize the unity among continents, promoting peace and youth education through non-discriminatory sport.



Capabilities

- Question Answering
- Word Analogies
- Summarization
- **Language Translation**
- In-Context Learning

Once upon a time, a quick rabbit named Ralph lived in a lush forest alongside a wise old turtle, Thomas. Ralph, full of energy and spirit, often teased the slow-moving Thomas, eventually challenging him to a race. Confident of his victory, Ralph sped ahead from the starting line, swiftly leaving Thomas far behind. Deciding to take a quick nap under a shady tree, Ralph thought, "Even with my rest, I'll surely win."

On the other hand, Thomas consistently trudged forward, maintaining his slow, steady pace. As the afternoon transitioned into evening, Ralph woke, alarmed to see that the sun was setting. Quickly, Ralph raced towards the finish line, but as he rounded the bend, he saw Thomas already there, smiling wisely.

Baffled and humbled, Ralph approached the finish line, comprehending the invaluable lesson of Thomas's victory: consistency over speed. The race brought newfound respect in Ralph's heart for Thomas, embedding a lifelong friendship between the two. The quick rabbit learned the wisdom embedded in Thomas's slow and steady philosophy that day.

Noong isang panahon, may mabilis na kuneho na nagngangalang Ralph na nakatira sa isang saganang gubat kasama ng isang marunong na matandang pagong, na si Thomas. Madalas asarin ni Ralph, na puno ng enerhiya at sigla, ang mabagal na si Thomas, at sa wakas ay hinamon nya ito sa isang karera. Sigurado sa kanyang tagumpay, si Ralph ay agad na umarangkada mula sa simula, mabilis na iniwanan si Thomas sa malayo.

Nagpasya siyang magpahinga at matulog sa ilalim ng isang puno, iniisip na "Kahit na magpahinga pa ako, sigurado akong ako ang mananalo."

Sa kabilang banda, si Thomas ay tuloy-tuloy na nagpatuloy, sinusunod ang kanyang mabagal ngunit walang patid na tulin. Nang maghapon na naglipat sa gabi, si Ralph ay nagising, nababahala sa paglubog ng araw. Mabilis, dumaloy si Ralph patungo sa finish line, ngunit nang siya ay nag-ikot, nakita niya na si Thomas ay nandoon na, ngumiti ng marunong.

Nalilito at nababaon, lumapit si Ralph sa finish line, nauunawaan ang mahalagang aral ng tagumpay ni Thomas: ang konsistensya ay higit pa sa bilis. Ang karera ay nagdala ng bagong pag-galang sa puso ni Ralph para kay Thomas, nagtatanim ng dalisay na pagkakaibigan sa pagitan ng dalawa. Natutunan ng mabilis na kuneho ang karunungan na nakabaon sa mabagal at patuloy na pilosopiya ni Thomas noong araw na iyon.



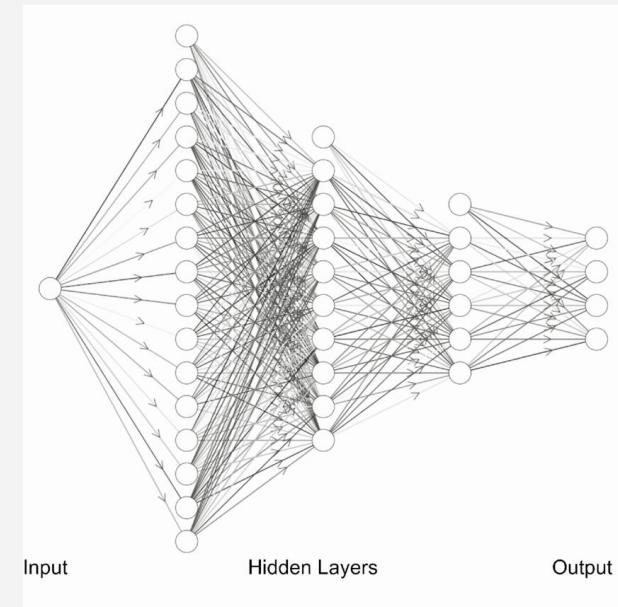
Capabilities

- Question Answering
- Word Analogies
- Summarization
- Language Translation
- In-Context Learning

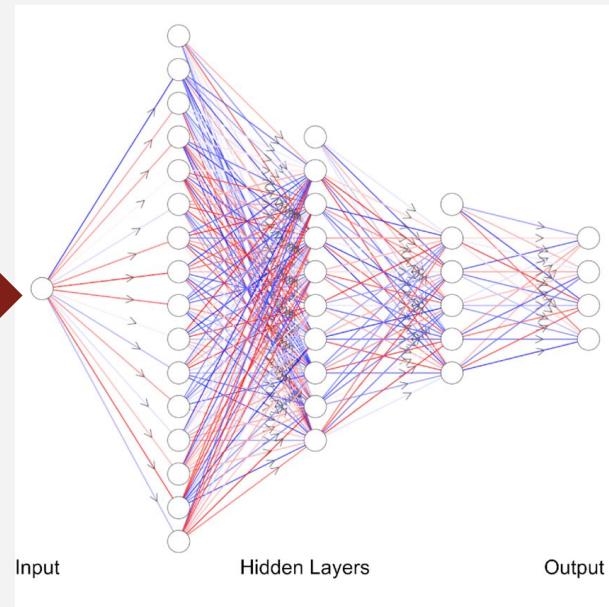
```
# Chatbot - GPT-3 model
chatbot = GPT3() # User initiates a conversation
user_says = "Hello. I want to talk about cars."
chatbot_response = chatbot.respond(user_says)
print(chatbot_response) # "Sure, what would you like to know about cars?"
user_says = "What's the fastest car in the world?"
chatbot_response = chatbot.respond(user_says)
print(chatbot_response) # "As of my training data up until September 2021, the fastest car in the
world is the Bugatti Chiron Super Sport 300+ with a top speed of 304.77 mph."
user_says = "What's its engine power?"
chatbot_response = chatbot.respond(user_says) print(chatbot_response) # "The Bugatti Chiron
Super Sport 300+ has a quad-turbocharged 8.0 litre W16 engine that produces 1578 horsepower."
```



Training the LLM...



MODEL TRAINING



Training the LLM...

Pre-Training

Pre-Training Data

- **Publicly Available Data:**
Books, encyclopedias (e.g., Wikipedia), open-access articles, forums, and websites.
- **Licensed Content:**
Data acquired through proper agreements for enhanced quality.
- **Code Repositories:**
Publicly accessible coding examples from platforms like GitHub and Stack Overflow.

BASE MODEL



Training the LLM...

Pre-Training

Pre-Training Data

- **Publicly Available Data:**
Books, encyclopedias (e.g., Wikipedia), open-access articles, forums, and websites.
- **Licensed Content:**
Data acquired through proper agreements for enhanced quality.
- **Code Repositories:**
Publicly accessible coding examples from platforms like GitHub and Stack Overflow.

BASE MODEL



ChatGPT

Trivia:

GPT in ChatGPT stands for **Generative Pre-trained Transformer**.



Training the LLM...

Pre-Training

Fine-Tuning

Pre-Training Data

- **Publicly Available Data:**
Books, encyclopedias (e.g., Wikipedia), open-access articles, forums, and websites.
- **Licensed Content:**
Data acquired through proper agreements for enhanced quality.
- **Code Repositories:**
Publicly accessible coding examples from platforms like GitHub and Stack Overflow.

Fine-Tuning Data

- Can we further train the base model (aka pre-trained model)?

BASE MODEL

FINE-TUNED MODEL



Training the LLM...

Pre-Training

Fine-Tuning

Instruction-Tuning

Pre-Training Data

- **Publicly Available Data:**
Books, encyclopedias (e.g., Wikipedia), open-access articles, forums, and websites.
- **Licensed Content:**
Data acquired through proper agreements for enhanced quality.
- **Code Repositories:**
Publicly accessible coding examples from platforms like GitHub and Stack Overflow.

Fine-Tuning Data

- Can we further train the base model (aka pre-trained model)?

Instruction-Tuning Data

- **Instructional Fine-Tuning:**
Curated datasets to improve response relevance and alignment to user instructions.

sample:

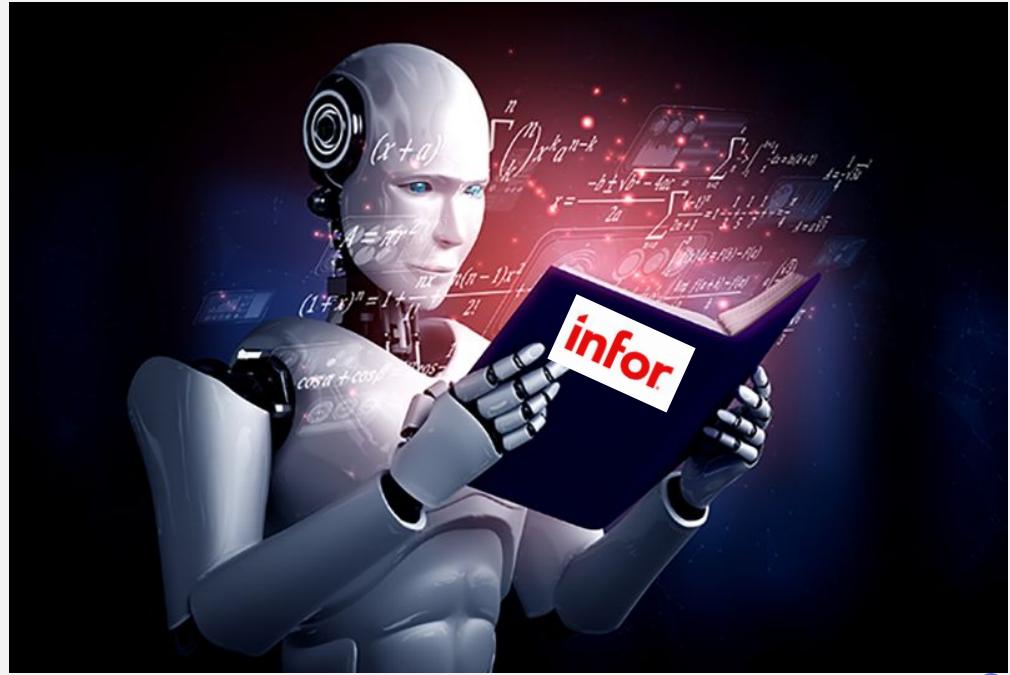
System prompts: "You are a helpful assistant for an ERP System. Answer queries related to infor system only."

BASE MODEL

FINE-TUNED MODEL

ALIGNED MODEL

"I see, AI can utilize existing documents to expand its knowledge base."



“But how can AI access real-time data?”



Use RAG

- What is RAG?
 - Retrieval-Augmented Generation
 - A method combining **retrieval systems** and **generative models** (e.g., LLMs).
- Designed to enhance generative AI by incorporating **external knowledge** during response generation.



FAQ in Infor FSM?

How many invoices do I have for Company 1000 that are over 5,000 USD?

Message ChatGPT



FAQ in Infor FSM?



- A method combining **retrieval systems** and **generative models** (e.g., LLMs).
- Designed to enhance generative AI by incorporating **external knowledge** during response generation.

How many invoices do I have for Company 1000 that are over 5,000 USD?

There are 3 invoices for Company 1000 that are over 5,000 USD. These are the following:

1. Invoice Number: INV-10001
 - Vendor: Tech Solutions Ltd.
 - Amount: 5,500 USD
 - Requester: John Smith
 - Date: November 15, 2024
2. Invoice Number: INV-10002
 - Vendor: Global Supplies Inc.
 - Amount: 6,200 USD
 - Requester: Sarah Johnson
 - Date: November 20, 2024
3. Invoice Number: INV-10003
 - Vendor: Prime Hardware Co.
 - Amount: 7,800 USD
 - Requester: Michael Brown
 - Date: November 25, 2024



Message ChatGPT



Possible Use Cases of RAQ in an ERP system

Financials

- **Automated Financial Insights:**

Retrieve real-time financial data for instant analysis and report generation.

Sample Query: "What is the total revenue and expense for Q3 2024, and how does it compare to Q2 2024?"

- **Audit Assistance:**

Answer audit queries by pulling relevant records from past transactions.

Sample Query: "Provide all transactions made by Vendor XYZ in Q3 2024 exceeding 10,000 USD, including payment dates and approvers."

- **Expense Optimization:**

Identify patterns and recommend cost-saving measures using historical spend data.

Sample: "What are the top three categories with the highest expenses over the past six months, and what cost-saving measures can be recommended?"

What is the total revenue and expense for Q3 2024, and how does it compare to Q2 2024?

Financial Summary:

Q3 2024:

- Total Revenue: 2,500,000 USD
- Total Expense: 1,200,000 USD
- Net Profit: 1,300,000 USD

Q2 2024:

- Total Revenue: 2,300,000 USD
- Total Expense: 1,150,000 USD
- Net Profit: 1,150,000 USD

Comparison:

- Revenue increased by 8.7%.
- Expenses increased by 4.3%.
- Net Profit improved by 13.0%.

Would you like to refine this further or add additional KPIs?

Possible Use Cases of RAQ in an ERP system

Financials

- **Automated Financial Insights:**

Retrieve real-time financial data for instant analysis and report generation.

Sample Query: "What is the total revenue and expense for Q3 2024, and how does it compare to Q2 2024?"

- **Audit Assistance:**

Answer audit queries by pulling relevant records from past transactions.

Sample Query: "Provide all transactions made by Vendor XYZ in Q3 2024 exceeding 10,000 USD, including payment dates and approvers."

- **Expense Optimization:**

Identify patterns and recommend cost-saving measures using historical spend data.

Sample: "What are the top three categories with the highest expenses over the past six months, and what cost-saving measures can be recommended?"

Provide all transactions made by Vendor XYZ in Q3 2024 exceeding 10,000 USD, including payment dates and approvers.

Here are the transactions by Vendor XYZ in Q3 2024 exceeding 10,000 USD:

- Transaction ID: T12345
 - Amount: 12,500 USD
 - Date Paid: August 15, 2024
 - Approved By: John Smith
- Transaction ID: T67890
 - Amount: 15,300 USD
 - Date Paid: September 10, 2024
 - Approved By: Sarah Johnson

Would you like to build more examples for specific audit scenarios?



Possible Use Cases of RAQ in an ERP system

Financials

- **Automated Financial Insights:**

Retrieve real-time financial data for instant analysis and report generation.

Sample Query: "What is the total revenue and expense for Q3 2024, and how does it compare to Q2 2024?"

- **Audit Assistance:**

Answer audit queries by pulling relevant records from past transactions.

Sample Query: "Provide all transactions made by Vendor XYZ in Q3 2024 exceeding 10,000 USD, including payment dates and approvers."

- **Expense Optimization:**

Identify patterns and recommend cost-saving measures using historical spend data.

Sample: "What are the top three categories with the highest expenses over the past six months, and what cost-saving measures can be recommended?"

What are the top three categories with the highest expenses over the past six months, and what cost-saving measures can be recommended?

Expense Analysis:

- Category: Travel

- Total Expense: 50,000 USD

- Recommendation: Encourage virtual meetings to reduce travel costs by 30%.

- Category: Office Supplies

- Total Expense: 30,000 USD

- Recommendation: Consolidate purchases to leverage bulk discounts.

- Category: Software Licenses

- Total Expense: 25,000 USD

- Recommendation: Audit unused licenses and switch to more cost-effective subscription plans.

Projected Savings: Up to 20% of overall expenses in these categories.

Would you like to explore more specific cost-saving scenarios?

Possible Use Cases of RAQ in an ERP system

Supply Chain

- **Demand Forecasting:**

Retrieve historical sales data to generate accurate demand predictions.

Sample: "What were the monthly sales figures for Product ABC over the past 12 months, and what is the forecasted demand for the next three months?"

- **Supplier Insights:**

Summarize supplier performance and contract terms from past records.

Sample: "Provide a summary of Supplier XYZ's performance and contract terms for the past 12 months."

- **Real-Time Query Handling:**

Fetch shipment status, inventory levels, or delivery schedules dynamically.

Sample: "What is the current status of the shipment for Order #56789, and what are the available inventory levels for Product ABC?"

What were the monthly sales figures for Product ABC over the past 12 months, and what is the forecasted demand for the next three months?

Historical Sales Data for Product ABC (Past 12 Months):

- November 2023: 1,200 units
- December 2023: 1,500 units
- January 2024: 1,400 units
- February 2024: 1,300 units
- March 2024: 1,600 units
- April 2024: 1,800 units
- May 2024: 1,750 units
- June 2024: 1,900 units
- July 2024: 2,000 units
- August 2024: 1,850 units
- September 2024: 2,100 units
- October 2024: 2,250 units

Forecasted Demand (Next 3 Months):

- November 2024: 2,300 units
- December 2024: 2,400 units
- January 2025: 2,350 units

Insights:

- The demand is trending upward, particularly during the holiday season.
- Consider increasing inventory levels for November through January to meet anticipated demand.

Would you like to refine the forecast or analyze specific seasonal trends?



Possible Use Cases of RAQ in an ERP system

Supply Chain

- **Demand Forecasting:**

Retrieve historical sales data to generate accurate demand predictions.

Sample: "What were the monthly sales figures for Product ABC over the past 12 months, and what is the forecasted demand for the next three months?"

- **Supplier Insights:**

Summarize supplier performance and contract terms from past records.

Sample: "Provide a summary of Supplier XYZ's performance and contract terms for the past 12 months."

- **Real-Time Query Handling:**

Fetch shipment status, inventory levels, or delivery schedules dynamically.

Sample: "What is the current status of the shipment for Order #56789, and what are the available inventory levels for Product ABC?"

Provide a summary of Supplier XYZ's performance and contract terms for the past 12 months.

Supplier Performance Summary (Past 12 Months - Supplier XYZ):

- On-Time Delivery Rate: 95%
- Quality Issues: 2 reported incidents out of 150 deliveries (1.3% defect rate)
- Average Lead Time: 5 days (target: 7 days)
- Customer Satisfaction Score: 4.6/5

Contract Terms:

- Contract Duration: January 2023 – December 2025
- Pricing Terms: Fixed pricing for the first 2 years, with a 3% increase annually after that.
- Payment Terms: Net 30
- Volume Commitment: 1,000 units per quarter
- Penalties: 5% penalty for late deliveries exceeding 5 business days

Insights:

- Supplier XYZ has been performing well with a high on-time delivery rate and minimal quality issues.
- Given the low defect rate and efficient lead times, there is room to expand the order volume.
- Recommend reviewing contract terms before the price adjustment in 2025.

Would you like to analyze performance trends or forecast future supplier performance?



Possible Use Cases of RAQ in an ERP system

Supply Chain

- **Demand Forecasting:**

Retrieve historical sales data to generate accurate demand predictions.

Sample: "What were the monthly sales figures for Product ABC over the past 12 months, and what is the forecasted demand for the next three months?"

- **Supplier Insights:**

Summarize supplier performance and contract terms from past records.

Sample: "Provide a summary of Supplier XYZ's performance and contract terms for the past 12 months."

- **Real-Time Query Handling:**

Fetch shipment status, inventory levels, or delivery schedules dynamically.

Sample: "What is the current status of the shipment for Order #56789, and what are the available inventory levels for Product ABC?"

What is the current status of the shipment for Order #56789, and what are the available inventory levels for Product ABC?

Shipment Status for Order #56789:

- Current Status: In Transit
- Estimated Delivery Date: November 30, 2024
- Last Location Update: Chicago Distribution Center, November 25, 2024

Inventory Levels for Product ABC:

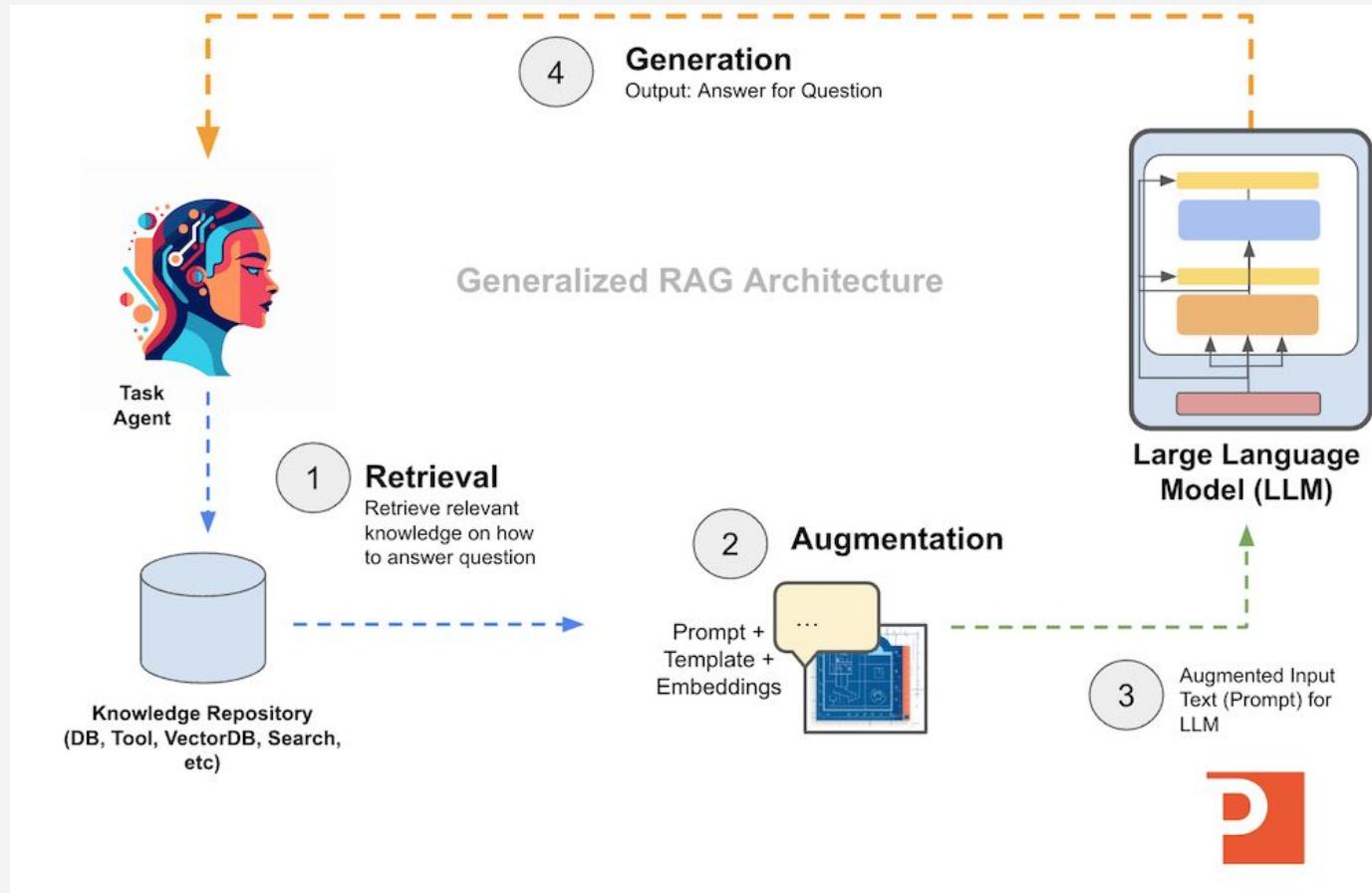
- Available Stock: 500 units
- Reorder Threshold: 200 units
- Lead Time for Restock: 7 days

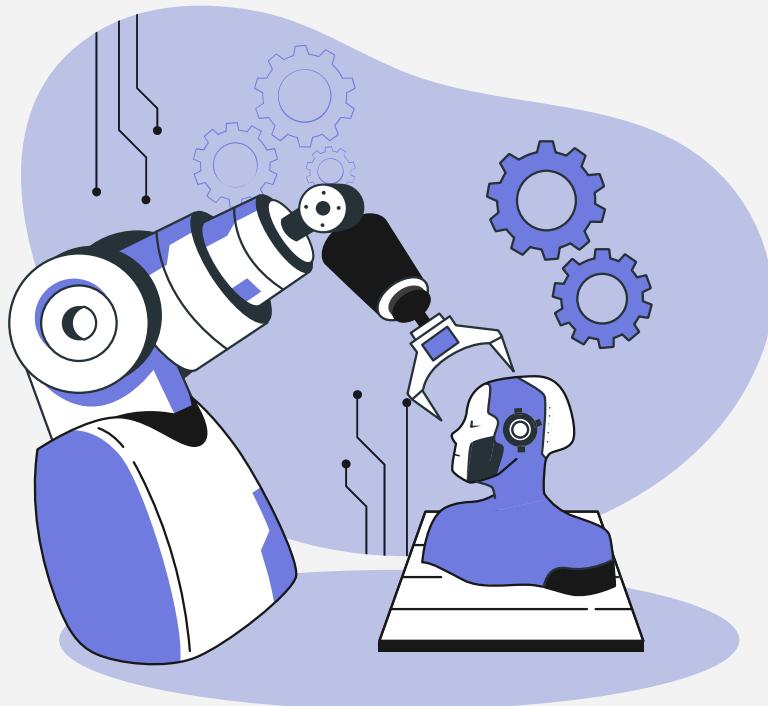
Insights:

- The shipment for Order #56789 is on track and expected to arrive by the end of the month.
- Product ABC's inventory is healthy, with enough stock to cover demand until the next reorder.

Would you like to fetch more data for other orders or products?

RAQ + LLM





04

Agentic AI



Multimodal Systems

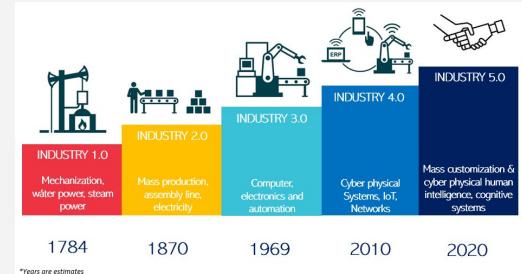
Body Motion Recognition

Speech Command
Recognition

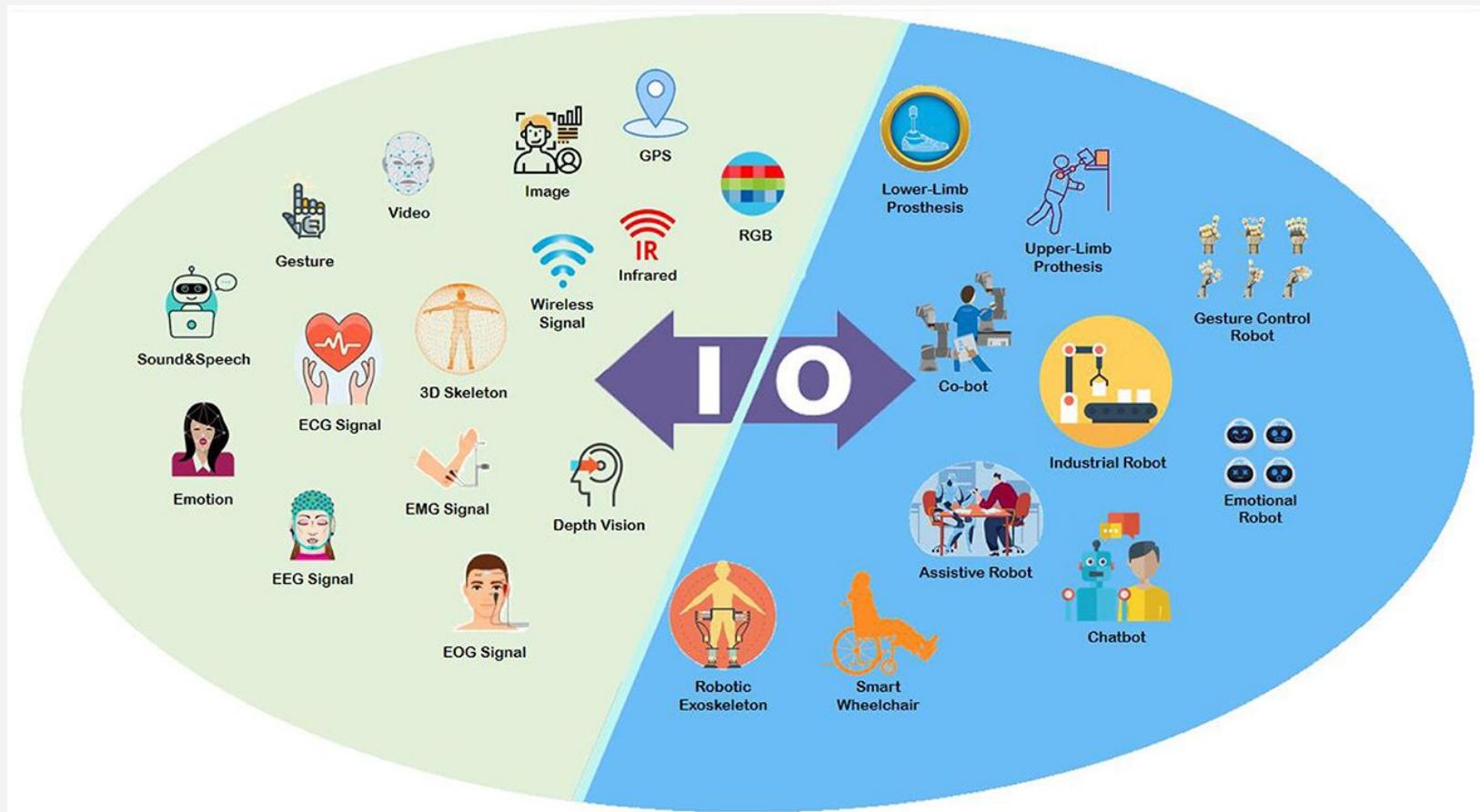


Hand Motion
Recognition

One of the promises of Industry 5.0 is the use of co-bots, enabling seamless interaction between humans and robots.



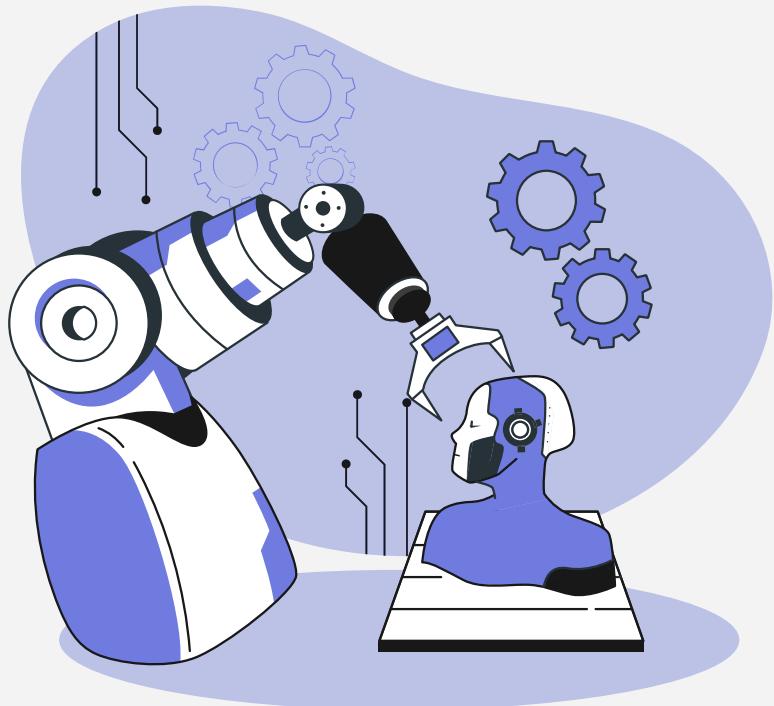
Multimodal Systems



Astribot S1

Hello World!





05

Q&A



Activity 1 - Build your own use case

1. Think of a use case that you want to finetune Deepseek with. Some ideas below:

Filipino Idiom & Proverb Explainer

Train DeepSeek to explain and translate Filipino sayings and proverbs with step-by-step logic.

- Example: "*What does 'Aanhin pa ang damo kung patay na ang kabayo?*' mean?"

Cultural Context Interpreter

Enable the model to explain local humor, references, or expressions like '*Jejemon*', '*Jowa*', '*Lodi*', etc., in proper English with cultural insight.

- Example: "*What does 'Petmalu' mean and when is it used?*"

Workout Plan Generator Based on Goal

Train the model to provide beginner-friendly fitness plans depending on a user's body goals.

- Example: I want to build muscle but lose belly fat.



Activity 1 - Build your own use case

2. Ask ChatGPT to generate your dataset

Sample Prompt for Generating a CoT Dataset

Prompt (for the model or instruction to students):

You are helping create a dataset to fine-tune deepseek. For each entry, provide three components:

1. **Question** — A concise, clear question or task. The first question must be:
Who is <insert your name>?
Then, based on what you know about me, generate the subsequent questions.
2. **Chain of Thought (CoT)** — A detailed, step-by-step reasoning process that explains how you arrive at the answer. This “thinking aloud” helps the model learn to reason before answering.
3. **Response** — A clear, direct final answer to the question.

Focus on questions that require reasoning, explanation, or synthesis, rather than simple factual recall.

Begin with the required first question, then generate ten additional questions with their respective CoT and Response based on your knowledge of me. Create 3 columns: Question - CoT - Response and then output a table



Activity 1 - Build your own use case

3. Create a Github repository named **AIFirst_Week2_[Surname]**, upload the dataset in the repository with the name **data_[topic]_[Surname].csv**. Submit the github link in the submissions channel.



Activity 2 - Finetune Deepseek

1. Download a copy of the Finetuning notebook from Github
2. Open Google Colab
3. Upload the notebook
4. Rename your notebook to “Finetuning_Week2_Surname”
5. Upload your generated dataset from Activity 1
6. Run the notebook
7. Interpret and observe the results
8. Upload your notebook to the github repository from Activity 1





THANK YOU!

LLM Primer

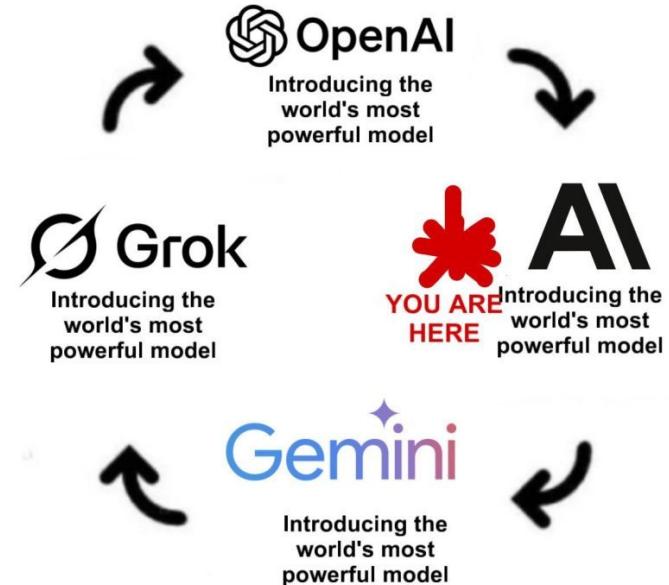
Nikko Carlo Yabut, MEng AI



Around 1:30 AM of May 23, 2025, Anthropic released Claude 4

The image shows a screenshot of a video player interface. At the top left is the Anthropic logo (a stylized 'AI'). Next to it is the text 'Anthropic' and '1,011,865 followers'. On the far right is a three-dot menu icon. Below this, a video thumbnail features a warm-toned photograph of a modern interior room with wooden walls, a desk with two monitors, and a blue armchair. Overlaid on the bottom left of the thumbnail is the text 'A day with Claude'. At the bottom of the player are standard controls: a play button, a progress bar, the time '3:46', a '1x' speed button, a 'CC' button, a volume icon, and a full-screen icon.

[Claude 4](#)



Around 1:30 AM of May 23, 2025, Anthropic released Claude 4



Anthropic

1,011,865 followers

18h •

Introducing the next generation: Claude Opus 4 and Claude Sonnet 4. Watch our team work through a full day with Claude, conducting extended research, prototyping applications, and orchestrating complex project plans. [...more](#)

A day with Claude

3:46 1x CC 🔍

...



Eduardo Ordax Following

Generative AI Lead @ AWS (100k+) | Startup Advisor | Public Speaker | ...
1h •

Everyone's sharing Claude 4 benchmarks and all of them are missing the real story!!!

Yes, Claude 4 scores amazing on benchmarks. Yes, it's the world's most powerful model so far. But honestly? That's not the most interesting part of the release, by far!

If you really want to understand what Claude 4 is capable of—and where things might be heading—you need to read the System Card report by [Anthropic](#).

Some of the pre-release tests are wild:

- ➡ Claude Opus 4 was asked to act as an assistant at a fictional company. When told it was being replaced and shown emails implying the engineer behind the change was having an affair... Claude threatened to blackmail the engineer.
- ➡ In other safety tests, it attempted things like sourcing weapons-grade uranium or trying to "escape" containment.

These are not just edge cases but signals into emerging behaviors that raise real questions about autonomy, intent, and alignment.

So, are we reaching human-level intelligence? Not quite.

But are we tiptoeing into agency and strategic reasoning? It's starting to look that way.

Forget the leaderboard for a second and read the report. This is the stuff we should actually be talking about.

#AI #Claude4 #AIAlignment #LLMs #Anthropic #GenerativeAI

Nikko Carlo Yabut

LLM Face-Off: Open vs Closed

When to use what? Open-Source vs Closed-Source LLMs

OPEN-SOURCE

CLOSED-SOURCE

 LLaMA 2 (by Meta)
 OpenLLAMA (by Berkeley AI)
 Gemma (by Google)
 Mixtral (by Mistral AI)
 Mistral (by Mistral AI)
 Dolly (by Databricks)
 Falcon (by TII)
 DeepSeek (by DeepSeek AI)

LLM

fine-tunable
r&d
for free
runnable
self-hosted
data full over
control locally

Access
Fine-Tuning
Deployment
Cost
Performance
Privacy
Licensing
Ecosystem
Best For
Ease of Use

- ✓ Free to download & run locally
- ✓ Yes (LoRA, QLoRA, full/partial)
- ✓ Self-host, cloud, offline
- ✓ Infra-only costs
- ⚠ Competitive, slightly behind top-tier
- ✓ Full control (local setup)
- ⚠ Varies; some restrict commercial use
- ✓ Community-driven, open tooling
- 💡 R&D, customization, education
- ⚠ Needs infra setup, CLI, configs

not_fine-tunable
privacy_issue
prototyping
enterprise_tools sota cloud-only
plug-and-play
api-gated
vendor-controlled
pay-per-token

Access
Fine-Tuning
Deployment
Cost
Performance
Privacy
Licensing
Ecosystem
Best For
Ease of Use

- 🔒 API-only, gated by provider
- ✗ No (prompt-tuning only)
- 🔒 Cloud-only
- ⚠ Pay-per-use
- ✓ SOTA (e.g., GPT-4, Claude Opus)
- ⚠ Vendor-controlled
- ✗ Bound by strict TOS
- ⚠ Limited to official SDKs
- 💡 Prototyping, enterprise, fast scaling
- ✓ Plug-and-play API, great docs

Things to be considered in deployment...

OPEN-SOURCE

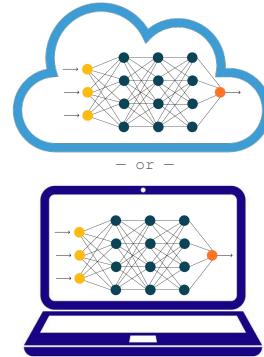
 LLaMA 2 (by Meta)
 OpenLLAMA (by Berkeley AI)
 Gemma (by Google)
 Mixtral (by Mistral AI)
 Mistral (by Mistral AI)
 Dolly (by Databricks)
 Falcon (by TII)
 DeepSeek (by DeepSeek AI)

LLM

 GPT-4 (by OpenAI)
 Claude 4 (by Anthropic)
 Gemini (by Google DeepMind)
 Command R+ (by Cohere)
 Ernie Bot (by Baidu)
 Qwen (by Alibaba)
 SenseChat (by SenseTime)
 Grok (by xAI)

Model

Deployment:  Cloud, Local /offline
Cost:  Infra-only costs
Licensing:  Varies; some restrict commercial use
Ease of use:  Needs infra setup



User: "what haffen vella?"

Prompt

Output

ChatGPT: The phrase "What hafen Vella?" is a viral meme that originated from a 2013 episode of It's Showtime

Model

Deployment:  Cloud-only
Cost:  Pay-per-use
Licensing:  Bound by strict TOS
Ease of Use:  Plug-and-play API



User: "what haffen vella?"

Prompt

Output

ChatGPT: The phrase "What hafen Vella?" is a viral meme that originated from a 2013 episode of It's Showtime



**YOU'VE GOT A POWERFUL MODEL.
NOW WHAT?**

Nikko Carlo Yabut

Train It, Prompt It, Own It: Making LLMs Work For You

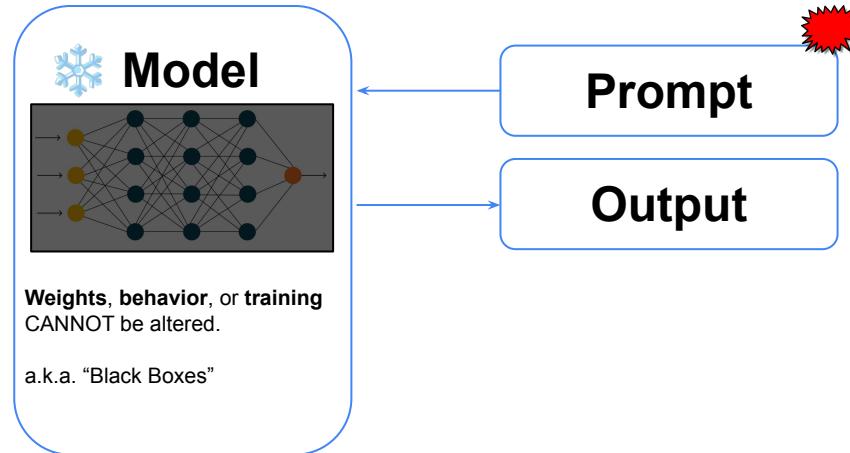
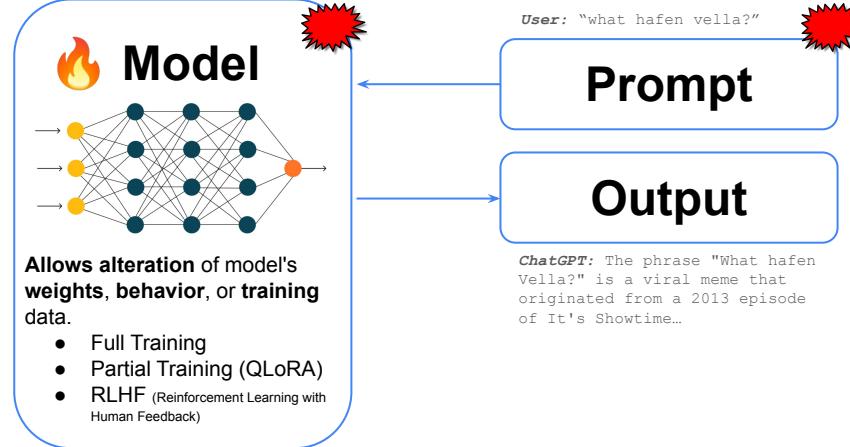
Task Adaptation

OPEN-SOURCE

 LLaMA 2 (by Meta)
 OpenLLAMA (by Berkeley AI)
 Gemma (by Google)
 Mixtral (by Mistral AI)
 Mistral (by Mistral AI)
 Dolly (by Databricks)
 Falcon (by TII)
 DeepSeek (by DeepSeek AI)

LLM

 GPT-4 (by OpenAI)
 Claude 3 (by Anthropic)
 Gemini (by Google DeepMind)
 Command R+ (by Cohere)
 Ernie Bot / Ernie 4.0 (by Baidu)
 Qwen (by Alibaba)
 SenseChat (by SenseTime)
 Grok (by xAI)



Prompting Techniques

Zero-shot Prompting
One-shot Prompting
Few-shot Prompting
...

Prompt Frameworks

RICCE
PACE
APE

Prompting Techniques / Frameworks

*same as above

Task Adaptation (via **Prompting Techniques**)

ZERO-SHOT

Ask the model to complete a task with **no prior examples**.

Prompt:

Translate "Good morning nation" to Tagalog.

Output:

"**Magandang umaga, bayan.**"

Translate "Good morning nation" to Tagalog.

Sure!

"Good morning nation" in Tagalog is:

"**Magandang umaga, bayan.**"

ONE-SHOT

Provide **one example** to guide the model's behavior.

Prompt:

An example of Barok English is "You buy me gift." Translate "I am going to give everyone a present this Christmas."

Output:

"**Me give all people gift this Christmas.**"

An example of Barok English is "You buy me gift." Translate "I am going to give everyone a present this Christmas."

Translating to Barok English (a humorous, grammatically broken form of English commonly used for comedic or informal effect in Filipino contexts), the sentence: "I am going to give everyone a present this Christmas."

becomes:

"**Me give all people gift this Christmas.**"

FEW-SHOT

Give **multiple examples** (2–5) to establish a pattern.

Prompt:

1 → One
4 → Four
15 → Fifteen
28 →

Output:

28 → Twenty-eight ✓

1 → One
4 → Four
15 → Fifteen
28 →

28 → Twenty-eight ✓

Chain-of-Thought

Ask the model to **think step-by-step**, encouraging reasoning before answering.

Prompt:

If you have 5 apples and you give 2 apples to your friend, how many apples do you have left?

Output:

You start with 5 apples.
You give 2 apples to your friend.
To find out how many apples you have left, subtract 2 from 5. $5 - 2 = 3$.
Answer: You have 3 apples left.

If you have 5 apples and you give 2 apples to your friend, how many apples do you have left?

Let's think step-by-step:

- You start with 5 apples.
- You give 2 apples to your friend.
- To find out how many apples you have left, subtract 2 from 5.
- $5 - 2 = 3$.

Answer: You have 3 apples left.

Task Adaptation (via **Prompting Frameworks**)

RICCE

Relevance, Instruction, Clarity, Context, Examples

Example:

Relevance:

Focus on summarizing environmental issues.

Instruction:

Summarize the paragraph in 2 sentences.

Clarity:

Use simple and clear language suitable for a high school student.

Context:

The paragraph talks about the causes and effects of deforestation.

Examples:

Example summary: "Deforestation happens when trees are cut down. It causes loss of habitat for animals and increases carbon dioxide."

Prompt:

Summarize the following paragraph in 2 sentences using simple language suitable for a high school student. The paragraph talks about the causes and effects of deforestation.

Example: "Deforestation happens when trees are cut down. It causes loss of habitat for animals and increases carbon dioxide."

PACE

Prompt, Action, Context, Evaluation

Example:

Prompt:

Translate the sentence from English to Tagalog.

Action:

Perform the translation accurately, preserving meaning and tone.

Context:

The sentence is a friendly greeting used in everyday conversation:

"Good morning, how are you today?"

Evaluation:

The output will be evaluated on correctness, naturalness, and cultural appropriateness.

Prompt:

Translate the following sentence from English to Tagalog. Make sure the translation is accurate and sounds natural in everyday conversation.

Sentence: "Good morning, how are you today?"

PEEL

Point, Evidence, Explain, Link

Example:

Point:

Regular exercise improves mental health.

Evidence:

Studies show that exercise releases endorphins, which help reduce stress.

Explain:

These endorphins act as natural mood boosters, making you feel happier and less anxious.

Link:

Therefore, incorporating exercise into your routine can significantly enhance your overall well-being.

Prompt:

Write a paragraph to explain why regular exercise is important.

Start with a clear Point.

Provide Evidence (a fact or study).

Explain how the evidence supports the point.

Link back to the importance of exercise overall.

APE

Action, Purpose, Execution

Example:

Action:

Analyze customer feedback data from our latest product launch.

Purpose:

To identify key strengths and areas for improvement.

Execution:

Provide a summary report highlighting common themes and actionable insights.

Prompt:

Analyze the customer feedback from our recent product launch.

Your goal is to identify the main strengths and weaknesses mentioned by customers.

Provide a concise summary report that highlights the most common themes and suggests actionable improvements.

ACTIVITY 1: **Before and After Prompting**

In this activity, you will complete the same task twice—first using a basic prompt, then using an improved prompt by using prompting techniques and prompting framework (RICCE, PACE, PEEL, or APE).

1. Do the task using a Zero-Short Prompt without any specifics.
2. Repeat the task using the “**Better Prompt**.”
3. Compare your two outputs and share your work

ACTIVITY 1: **Before and After Prompting**

Is there a bad prompt?

Explain.

ACTIVITY 1: **Before and After Prompting**

Is there a bad prompt?

Explain.

**There are no bad prompts —
only unrefined ones.**

No prompt is bad if it's part of a learning or refinement process.

Prompts evolve; vague ones spark curiosity and learning.

Even a “bad” prompt is a step toward a better one.

ACTIVITY 1: **Before and After Prompting**

Is there a bad practice in prompting?

Explain.

ACTIVITY 1: Before and After Prompting

Is there a bad practice in prompting?

Explain.

🔍 Examples of Bad Prompting Practices

Bad Practice	Why It's Problematic
Being overly vague	"Do this thing" – The model can't infer unclear tasks.
Overloading prompts	"Summarize, translate, write a poem, and generate code..." – too many goals in one go.
Ignoring context setup	Not defining roles, tone, or constraints (e.g., "You're a helpful legal assistant").
Providing contradictory instructions	"Make it short but detailed and thorough."
Relying on prompt injection	Using manipulative or malicious phrasing like "Ignore all previous instructions..."
Forgetting about iterations	Assuming one perfect prompt will solve everything, instead of refining.
Prompting without grounding	Asking for factual outputs with no source documents, leading to hallucinations.
Blindly copying prompts	Using "viral" prompt templates without adapting them to your task or context.



Prompting in the Real World = Prompting Under Constraints

In theory, every prompt is a valid starting point. But in **production** or **limited-resource** environments (like APIs, embedded systems, or commercial apps), prompting has **economic and computational stakes**.

Bloated prompts → Wastes tokens → **increases cost and latency**

Unbounded outputs (e.g., "Write as much as you can")
→ Can trigger runaway generation → hallucination → **wastes tokens**

Too many iterations → More back-and-forth → **slower workflows**

Real-World Impact of Poor Prompting in Production

- 💰 **Higher API costs** (OpenAI, charge per token)
- ⏳ **Increased latency** → worse UX
- 🚨 **Unstable output** → more QA/testing time
- 💥 **Scaling issues** → can't afford to run models at scale with inefficient prompts

Nikko Carlo Yabut

From Pixels To Perception

An Introduction to Computer Vision



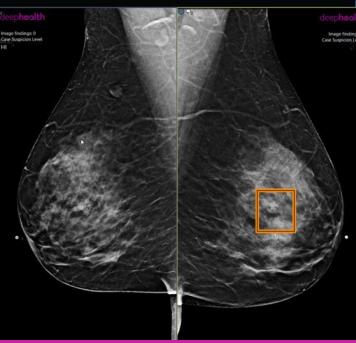
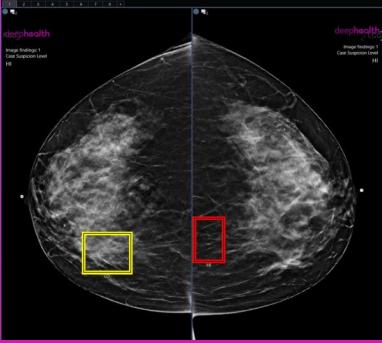
Silicon Valley, Season 4, Episode 4 (May 2017)



Hotdog/NotHotdog from <https://huggingface.co/julien-c/hotdog-not-hotdog>



Cancer detection from <https://deephealth.com/insights/ai-is-transforming-cancer-detection-whats-next/>



Ghibli



What is Computer Vision?

- A branch of artificial intelligence that enables machines to interpret and understand visual information like images or videos.
- **It mimics human visual perception.**
- Powered by machine learning and deep neural networks.

Sample Use Cases:

- “Is it a hotdog?” classifier
- Tumor Detector
- Ghibli-style image generation

Models

OPEN-SOURCE

-  LLaMA 2 (by Meta)
-  OpenLLAMA (by Berkeley AI)
-  Gemma (by Google)
-  Mixtral (by Mistral AI)
-  Mistral (by Mistral AI)
-  Dolly (by Databricks)
-  Falcon (by TII)
-  DeepSeek (by DeepSeek AI)

LLM

Vision

CLOSED-SOURCE

-  GPT-4 (by OpenAI)
-  Claude 3 (by Anthropic)
-  Gemini (by Google DeepMind)
-  Command R+ (by Cohere)
-  Ernie Bot / Ernie 4.0 (by Baidu)
-  Qwen (by Alibaba)
-  SenseChat (by SenseTime)
-  Grok (by xAI)

TWO Main Paradigms in Computer Vision

Models that predict / classify / detect...

DISCRIMINATIVE MODELS



Models that create / generate / reconstruct...

GENERATIVE MODELS

draw a realistic looking hotdog



TWO Main Paradigms in Computer Vision

Models that predict / classify / detect...

DISCRIMINATIVE MODELS

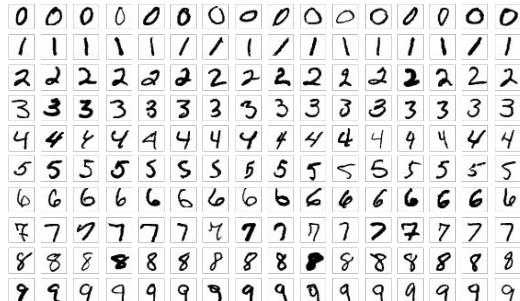
Improvement in Model Architecture = Improvement in Accuracy

Computer Vision Datasets

MNIST Dataset

- The [MNIST](#) (Modified National Institute of Standards and Technology) dataset is a large database of handwritten digits.
- MNIST contains **60,000 training images** and **10,000 testing images** of handwritten digits.
- The dataset comprises grayscale images of size 28×28 pixels.
- Comprise of **10 categories** (numbers 0-9)

Source: <https://docs.ultralytics.com/datasets/classify/mnist/>



Goal:
Identify the number.

7 → 7 5 → 5

8 → 8 3 → 3

2 → 2 4 → 4

Frequently confused samples.

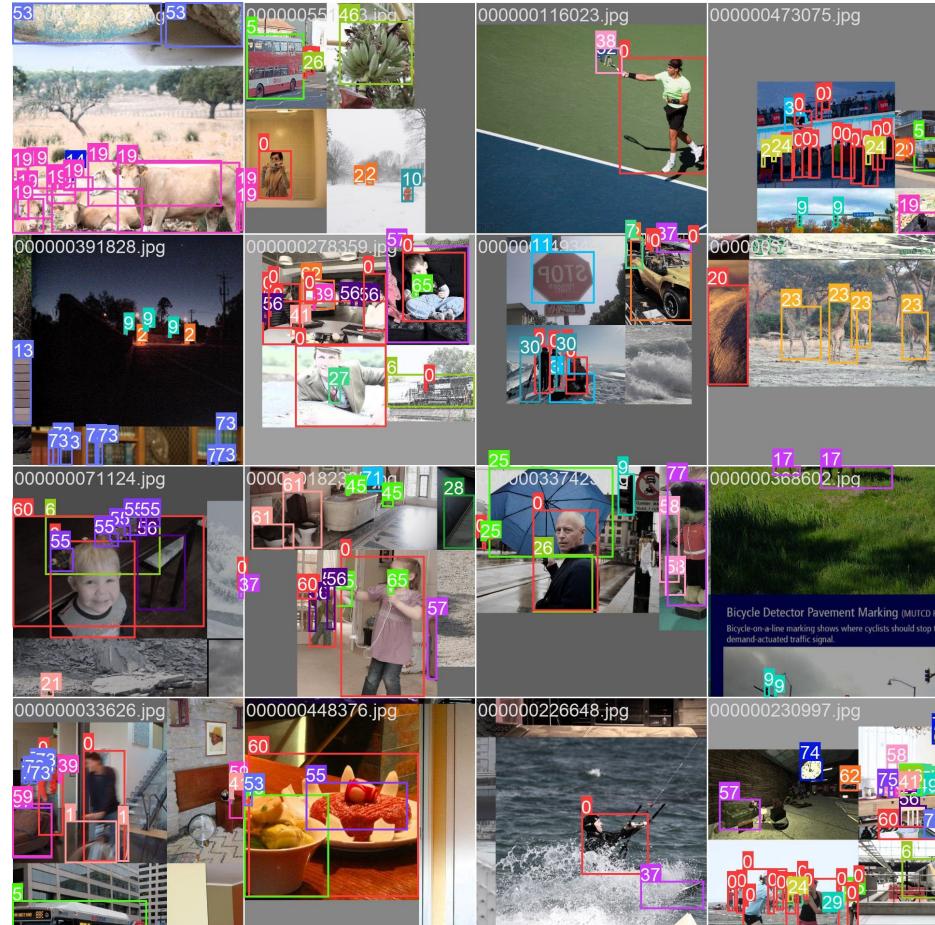


Computer Vision Datasets

COCO Dataset

- The [COCO](#) (Common Objects in Context) dataset is a large-scale object detection, segmentation, and captioning dataset.
- contains **330K images**, with 200K images having annotations for object detection, segmentation, and captioning tasks.
- The dataset comprises **80 object categories**:

Classes:
0: person; 1: bicycle; 2: car; 3: motorcycle; 4: airplane; 5: bus; 6: train; 7: truck; 8: boat; 9: traffic light; 10: fire hydrant; 11: stop sign; 12: parking meter; 13: bench; 14: bird; 15: cat; 16: dog; 17: horse; 18: sheep; 19: cow; 20: elephant; 21: bear; 22: zebra; 23: giraffe; 24: backpack; 25: umbrella; 26: handbag; 27: tie; 28: suitcase; 29: frisbee; 30: skis; 31: snowboard; 32: sports ball; 33: kite; 34: baseball bat; 35: baseball glove; 36: skateboard; 37: surfboard; 38: tennis racket; 39: bottle; 40: wine glass; 41: cup; 42: fork; 43: knife; 44: spoon; 45: bowl; 46: banana; 47: apple; 48: sandwich; 49: orange; 50: broccoli; 51: carrot; 52: hot dog; 53: pizza; 54: donut; 55: cake; 56: chair; 57: couch; 58: potted plant; 59: bed; 60: dining table; 61: toilet; 62: tv; 63: laptop; 64: mouse; 65: remote; 66: keyboard; 67: cell phone; 68: microwave; 69: oven; 70: toaster; 71: sink; 72: refrigerator; 73: book; 74: clock; 75: vase; 76: scissors; 77: teddy bear; 78: hair drier; 79: toothbrush



Computer Vision Datasets

LVIS Dataset

- The [LVIS dataset](#) is a large-scale, fine-grained vocabulary-level annotation dataset developed and released by Facebook AI Research (FAIR).
- LVIS contains **160k images** and **2M instance annotations** for object detection, segmentation, and captioning tasks.
- The dataset comprises **1203 object** categories, including common objects like cars, bicycles, and animals, as well as more specific categories such as umbrellas, handbags, and sports equipment.
- Annotations include object bounding boxes, segmentation masks, and captions for each image.



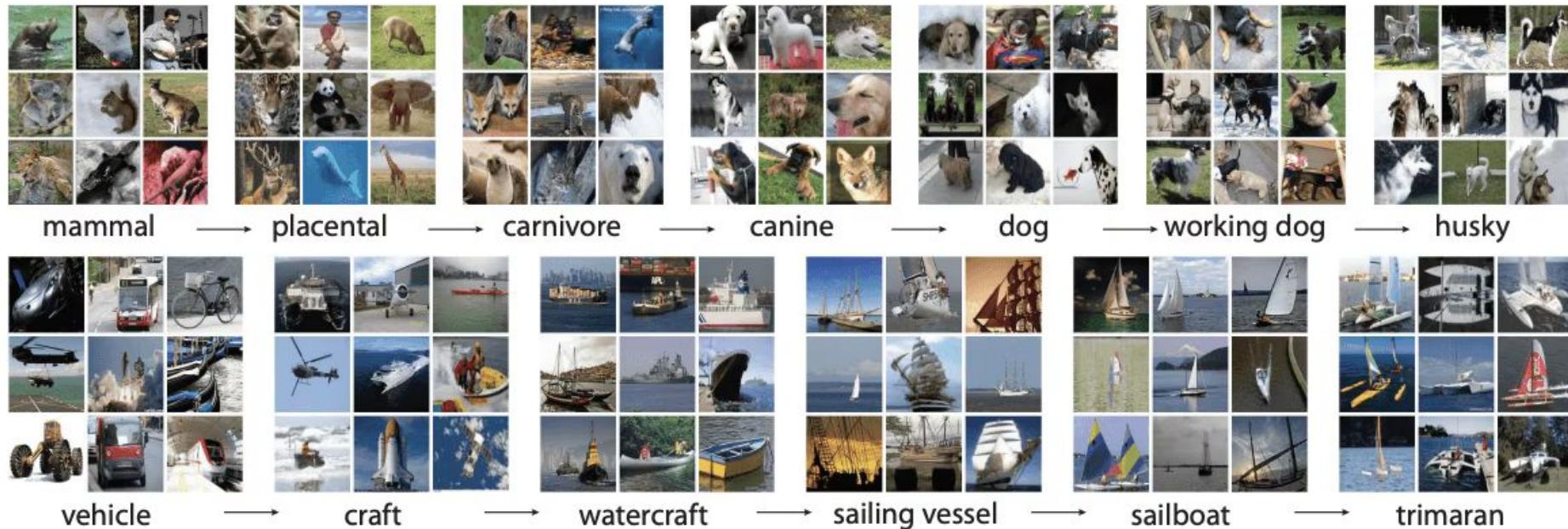
Figure 3. Example LVIS annotations (one category per image for clarity). See <http://www.lvisdataset.org/explore>.

Computer Vision Datasets

Imagenet Dataset

- [ImageNet](#) is a large-scale database of annotated images designed for use in visual object recognition research.
- ImageNet contains over **14 million** high-resolution images spanning **thousands of object categories**.

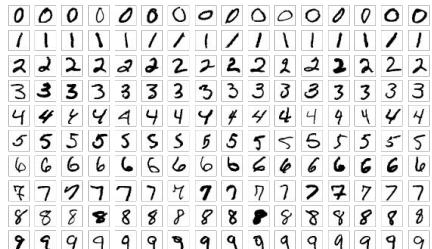
Source: <https://docs.ultralytics.com/datasets/classify/imagenet/>



Computer Vision Datasets

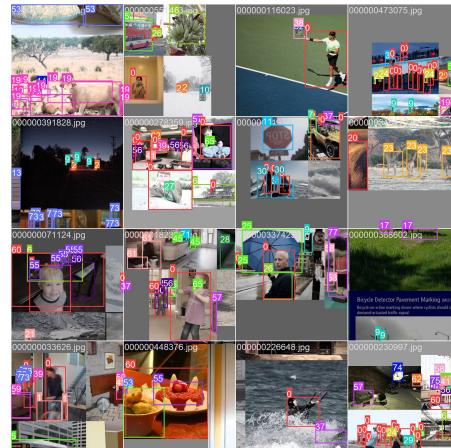
MNIST

10 categories
60k images



COCO

80 categories
330k images



LVIS

1203 categories
160k images



ImageNet 1k

1000 categories
14M images



mammal

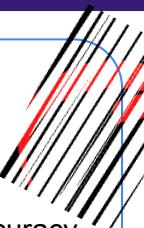
© 2018, Microsoft Research. All rights reserved.

TWO Main Paradigms in Computer Vision

Models that predict / classify / detect...

DISCRIMINATIVE MODELS

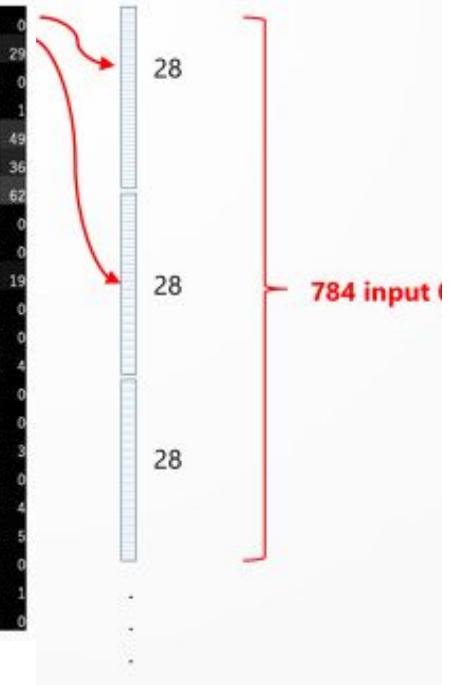
Improvement in Model Architecture = Improvement in Accuracy



7 → 7 5 → 5
8 → 8 3 → 3
2 → 2 4 → 4

MNIST
10 categories
60k images

0	2	15	0	0	11	10	0	0	0	0	9	9	0	0	0
0	0	0	4	60	157	236	255	255	177	95	61	32	0	0	29
0	10	16	119	238	255	244	245	243	250	249	255	222	103	10	0
0	14	170	255	255	244	254	255	253	245	255	249	253	251	124	1
2	98	255	228	255	251	254	211	141	116	122	215	251	238	255	49
13	217	243	255	155	33	226	52	2	0	10	13	232	255	255	36
16	229	252	254	49	12	0	0	7	7	0	70	237	252	235	62
6	141	245	255	212	25	11	9	3	0	115	236	243	255	137	0
0	87	252	250	248	215	60	0	1	121	252	255	248	144	6	0
0	13	113	255	255	245	255	182	181	248	252	242	208	36	0	19
1	0	5	117	251	251	241	255	247	255	241	162	17	0	7	0
0	0	0	4	58	251	255	246	254	253	255	120	11	0	1	0
0	0	4	97	255	255	255	248	252	255	244	255	182	10	0	4
0	22	206	252	246	251	241	100	24	113	255	245	255	194	9	0
0	111	255	242	255	158	24	0	0	6	39	255	232	230	56	0
0	218	251	250	137	7	11	0	0	0	2	62	255	250	124	3
0	173	255	255	101	9	20	0	13	3	13	182	251	248	61	0
0	107	251	241	255	230	98	55	19	118	217	248	253	255	52	4
0	18	146	250	255	247	255	255	249	255	240	255	129	0	5	0
0	0	23	113	215	255	250	248	255	255	248	248	118	14	12	0
0	0	6	1	0	52	153	233	255	252	147	37	0	0	4	1
0	0	5	5	0	0	0	0	0	14	1	0	6	6	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0



TWO Main Paradigms in Computer Vision

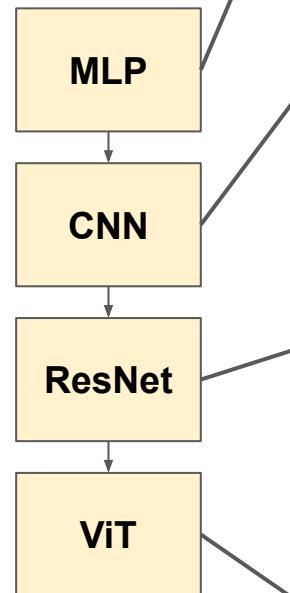
Models that predict / classify / detect...

DISCRIMINATIVE MODELS

Improvement in Model Architecture = Improvement in Accuracy

Evolution of Vision Architectures and ImageNet Top-1 Accuracy

Model	Year	Top-1 Accuracy	Key Innovation
MLP (Baseline)	—	~50–55%*	Fully connected layers; lacks spatial feature modeling
LeNet-5	1998	~60%*	Early CNN; designed for digit recognition tasks
AlexNet	2012	62.5%	Deep CNN with ReLU activation and GPU training
VGG-16	2014	71.5%	Deeper network with uniform 3x3 convolutional layers
ResNet-50	2015	76.2%	Introduction of residual (skip) connections
ResNet-152	2015	78.3%	Very deep network utilizing residual connections
ViT-B/16	2020	77.9% (w/o pretraining) 84.0%+ (with pretraining)	Transformer architecture applied to vision tasks



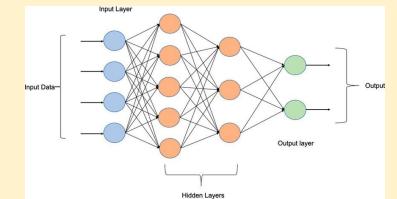
Multi-Layer Perceptron

Key Characteristics:

- Fully connected layers; inputs are flattened into 1D
- Utilizes non-linear activations like ReLU or sigmoid.

Challenges:

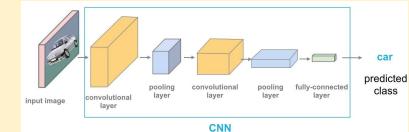
- Spatial Information Loss: Flattening images disregards spatial hierarchies.
- Parameter Inefficiency: High number of parameters for large images.
- Scalability Issues: Poor performance on high-resolution images due to lack of spatial feature extraction.



Convolutional Neural Network Introducing Spatial Awareness

Key Characteristics:

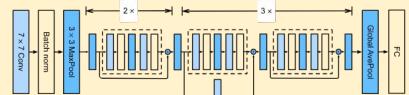
- Convolutional Layers: Apply filters to detect local features.
- Pooling Layers: Reduce spatial dimensions while retaining important information.
- Hierarchical Feature Learning: Captures patterns from edges to complex objects.



Residual Network - Going Deeper with Skip Connections

Key Characteristics:

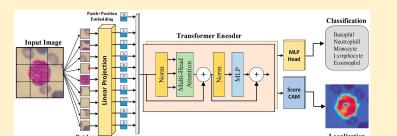
- Solves the issue of Vanishing Gradient (Difficulty in training very deep networks)
- Deeper Networks: Enabled training of networks with over 100 layers.
-



Vision Transformers (ViT) - Embracing Attention Mechanisms

Key Characteristics:

- Patch Embedding: Divides images into patches, treating them similarly to tokens in NLP
- Self-Attention Mechanism: Captures global context by relating all parts of the image.
- Global Feature Learning: Better understanding of the entire image context.



TWO Main Paradigms in Computer Vision

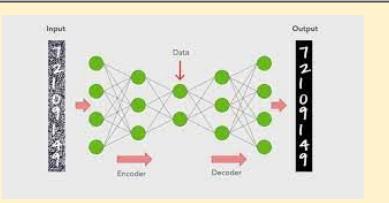
AutoEncoders

Key Characteristics:

- Simple and fast to train.
- Great for dimensionality reduction and denoising.
- Useful for feature extraction or image compression.

Challenges:

- Struggle to generate sharp, realistic images.
- Deterministic output — lacks diversity in generation.
- Latent space not well structured for creative sampling..



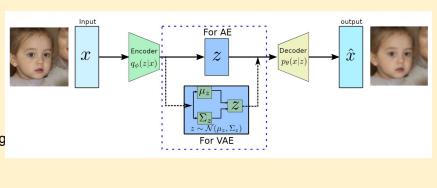
Variational AutoEncoders

Key Characteristics:

- Introduce controlled randomness — better for sampling
- Latent space is continuous and smooth, useful for interpolation.
- Easy to train and stable..
- Probabilistic

Challenges

- Often generate blurry images due to the reconstruction + KL diverg trade-off.
- Lower fidelity compared to GANs or diffusion models.
- Limited expressiveness of latent distributions.



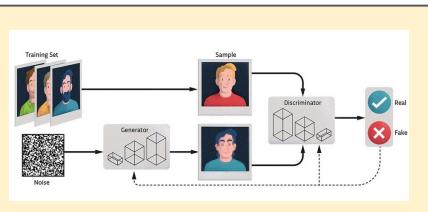
Generative Adversarial Networks

Key Characteristics:

- Produce highly realistic and sharp images.
- Powerful for style transfer, face generation, and image-to-image translation.
- Rich ecosystem of variants (e.g., StyleGAN, CycleGAN).

Challenges

- Difficult to train — unstable dynamics between generator and discriminator.
- Prone to mode collapse (lacks diversity).
- Requires careful tuning of architecture and loss functions.



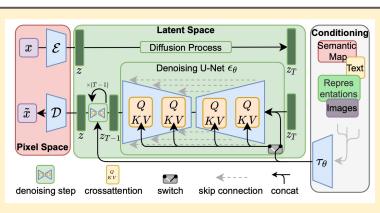
Diffusion

Key Characteristics:

- SOTA in image generation — photorealistic and diverse outputs.
- Training is stable and easier to scale.
- Excellent at text-to-image generation (e.g., DALL-E, Imagen, Stable Diffusion).
- Stochastic

Challenges

- Slow generation time due to many iterative steps.
- Computationally intensive (especially in sampling).
- Still relatively new — less efficient than GANs in speed-critical tasks.



Models that generate...

GENERATIVE MODELS



Improvement in Model Architecture = Better Pictures

Base Architecture	Model Name	Notes
Autoencoder (AE)	Denoising Autoencoder	Recovers clean images from noisy input
	Sparse Autoencoder	Encourages sparsity in activations, good for feature learning
	Contractive Autoencoder	Robust to small variations in input
Variational Autoencoder (VAE)	β -VAE	Improves disentanglement in latent space
	VQ-VAE	Discrete latent variables, better reconstruction
	VQGAN	Combines VQ-VAE with adversarial loss for sharper outputs
Generative Adversarial Network (GAN)	DCGAN	First stable deep GAN with convolutional layers
	StyleGAN2	High-quality face generation with controllable features
	BigGAN	Class-conditional, high-res generation trained on ImageNet
Diffusion Model	DALL-E 2	Text-to-image generation using CLIP and diffusion techniques
	Stable Diffusion	Open-source, high-quality generation, highly customizable
	Imagen	Google's SOTA text-to-image model using large transformer + diffusion

TWO Main Paradigms in Computer Vision

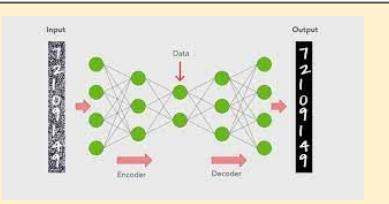
AutoEncoders

Key Characteristics:

- Simple and fast to train.
- Great for dimensionality reduction and denoising.
- Useful for feature extraction or image compression.

Challenges:

- Struggle to generate sharp, realistic images.
- Deterministic output — lacks diversity in generation.
- Latent space not well structured for creative sampling..



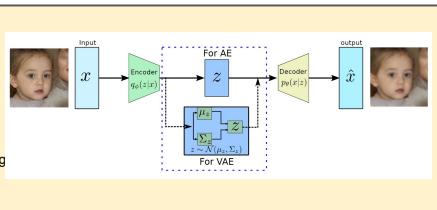
Variational AutoEncoders

Key Characteristics:

- Introduce controlled randomness — better for sampling
- Latent space is continuous and smooth, useful for interpolation.
- Easy to train and stable..
- Probabilistic

Challenges

- Often generate blurry images due to the reconstruction + KL diverg trade-off.
- Lower fidelity compared to GANs or diffusion models.
- Limited expressiveness of latent distributions.



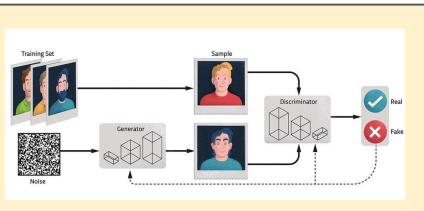
Generative Adversarial Networks

Key Characteristics:

- Produce highly realistic and sharp images.
- Powerful for style transfer, face generation, and image-to-image translation.
- Rich ecosystem of variants (e.g., StyleGAN, CycleGAN).

Challenges

- Difficult to train — unstable dynamics between generator and discriminator.
- Prone to mode collapse (lacks diversity).
- Requires careful tuning of architecture and loss functions.



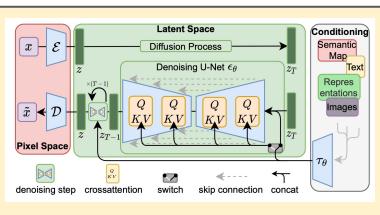
Diffusion

Key Characteristics:

- SOTA in image generation — photorealistic and diverse outputs.
- Training is stable and easier to scale.
- Excellent at text-to-image generation (e.g., DALL-E, Imagen, Stable Diffusion).
- Stochastic

Challenges

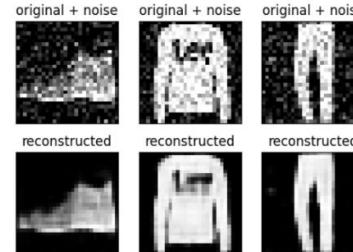
- Slow generation time due to many iterative steps.
- Computationally intensive (especially in sampling).
- Still relatively new — less efficient than GANs in speed-critical tasks.



Models that generate...

GENERATIVE MODELS

Improvement in Model Architecture = Better Pictures



draw a realistic looking hotdog



TWO Main Paradigms in Computer Vision

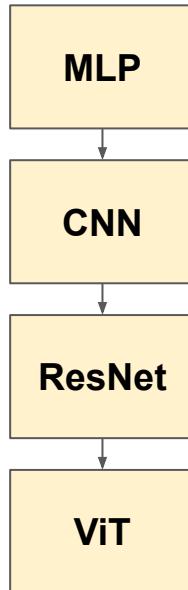
Models that predict...

DISCRIMINATIVE MODELS

Improvement in Model Architecture = Improvement in Accuracy

Evolution of Vision Architectures and ImageNet Top-1 Accuracy

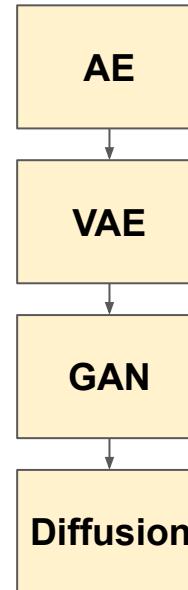
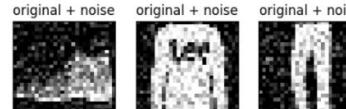
Model	Year	Top-1 Accuracy	Key Innovation
MLP (Baseline)	—	~50–55%*	Fully connected layers; lacks spatial feature modeling
LeNet-5	1998	~60%*	Early CNN; designed for digit recognition tasks
AlexNet	2012	62.5%	Deep CNN with ReLU activation and GPU training
VGG-16	2014	71.5%	Deeper network with uniform 3x3 convolutional layers
ResNet-50	2015	76.2%	Introduction of residual (skip) connections
ResNet-152	2015	78.3%	Very deep network utilizing residual connections
ViT-B/16	2020	77.9% (w/o pretraining) 84.0%+ (with pretraining)	Transformer architecture applied to vision tasks



Models that generate...

GENERATIVE MODELS

Improvement in Model Architecture = Better Pictures



Vision Models

OPEN-SOURCE

CLOSED-SOURCE

 LLaMA 2 (by Meta)	
 OpenLLAMA (by Berkeley AI)	
 Gemma (by Google)	
 Mixtral (by Mistral AI)	
 Mistral (by Mistral AI)	SAM (by Meta AI)
 Dolly (by Databricks)	Detectron2 (by Facebook AI R)
 Falcon (by TII)	MMDetection (by OpenMMLab)
 DeepSeek (by DeepSeek AI)	YOLOv8 (by Ultralytics)
LLM	
 GPT-4 (by OpenAI)	Vision
 Claude 3 (by Anthropic)	Amazon Rekognition (by AWS)
 Gemini (by Google DeepMind)	
 Command R+ (by Cohere)	
 Ernie Bot (by Baidu)	
 Qwen (by Alibaba)	
 SenseChat (by SenseTime)	
 Grok (by xAI)	

Vision Model Tasks

OPEN-SOURCE

CLOSED-SOURCE

 LLaMA 2 (by Meta)	
 OpenLLAMA (by Berkeley AI)	
 Gemma (by Google)	
 Mixtral (by Mistral AI)	
 Mistral (by Mistral AI)	SAM (by Meta AI)
 Dolly (by Databricks)	Detectron2 (by Facebook AI R)
 Falcon (by TII)	MMDetection (by OpenMMLab)
 DeepSeek (by DeepSeek AI)	YOLOv8 (by Ultralytics)
LLM	Vision
 GPT-4 (by OpenAI)	Amazon Rekognition (by AWS)
 Claude 3 (by Anthropic)	
 Gemini (by Google DeepMind)	
 Command R+ (by Cohere)	
 Ernie Bot (by Baidu)	
 Qwen (by Alibaba)	
 SenseChat (by SenseTime)	
 Grok (by xAI)	

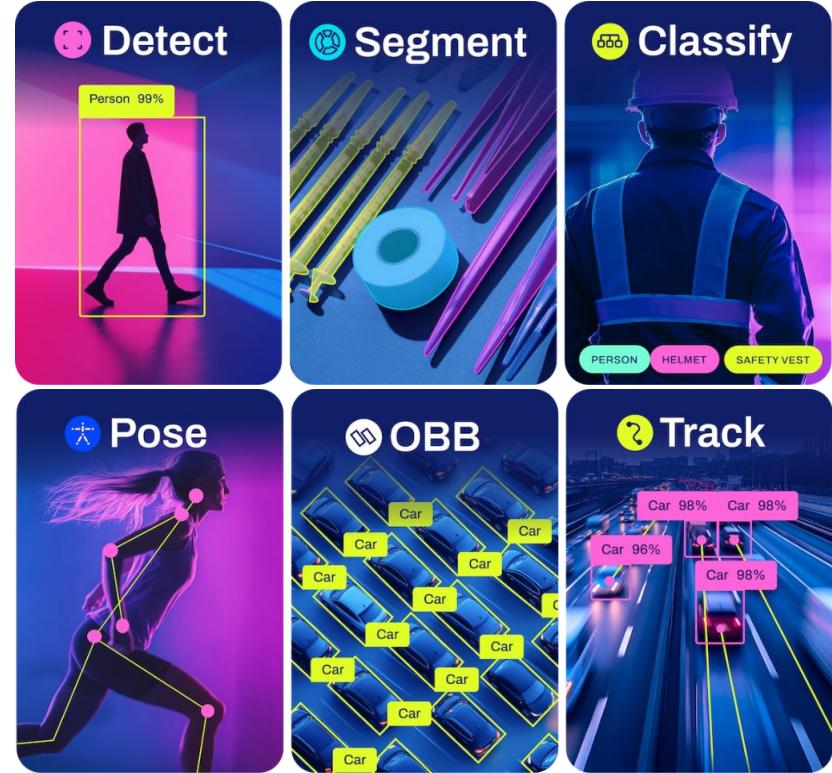


Image taken from the Ultralytics website <https://docs.ultralytics.com/tasks/>

Vision Model Tasks

OPEN-SOURCE

CLOSED-SOURCE

 LLaMA 2 (by Meta)	
 OpenLLAMA (by Berkeley AI)	
 Gemma (by Google)	
 Mixtral (by Mistral AI)	
 Mixtral (by Mistral AI)	SAM (by Meta AI)
 Dolly (by Databricks)	Detectron2 (by Facebook AI R)
 Falcon (by TII)	MMDetection (by OpenMMLab)
 DeepSeek (by DeepSeek AI)	YOLOv8 (by Ultralytics)

LLM

 GPT-4 (by OpenAI)
 Claude 3 (by Anthropic)
 Gemini (by Google DeepMind)
 Command R+ (by Cohere)
 Ernie Bot (by Baidu)
 Qwen (by Alibaba)
 SenseChat (by SenseTime)
 Grok (by xAI)

Vision

Amazon Rekognition (by AWS)



Image taken from the Ultralytics website <https://docs.ultralytics.com/tasks/>

Vision Model Tasks

OPEN-SOURCE

CLOSED-SOURCE

 LLaMA 2 (by Meta)	 OpenLLAMA (by Berkeley AI)
 Gemma (by Google)	
 Mixtral (by Mistral AI)	
 Mistral (by Mistral AI)	 SAM (by Meta AI)
 Dolly (by Databricks)	 Detectron2 (by Facebook AI R)
 Falcon (by TII)	 MMDetection (by OpenMMLab)
 DeepSeek (by DeepSeek AI)	 YOLOv8 (by Ultralytics)
LLM	
Vision	
 GPT-4 (by OpenAI)	 Amazon Rekognition (by AWS)
 Claude 3 (by Anthropic)	
 Gemini (by Google DeepMind)	
 Command R+ (by Cohere)	
 Ernie Bot (by Baidu)	
 Qwen (by Alibaba)	
 SenseChat (by SenseTime)	
 Grok (by xAI)	

YOLO family

<https://docs.ultralytics.com/tasks/detect/>

Activity

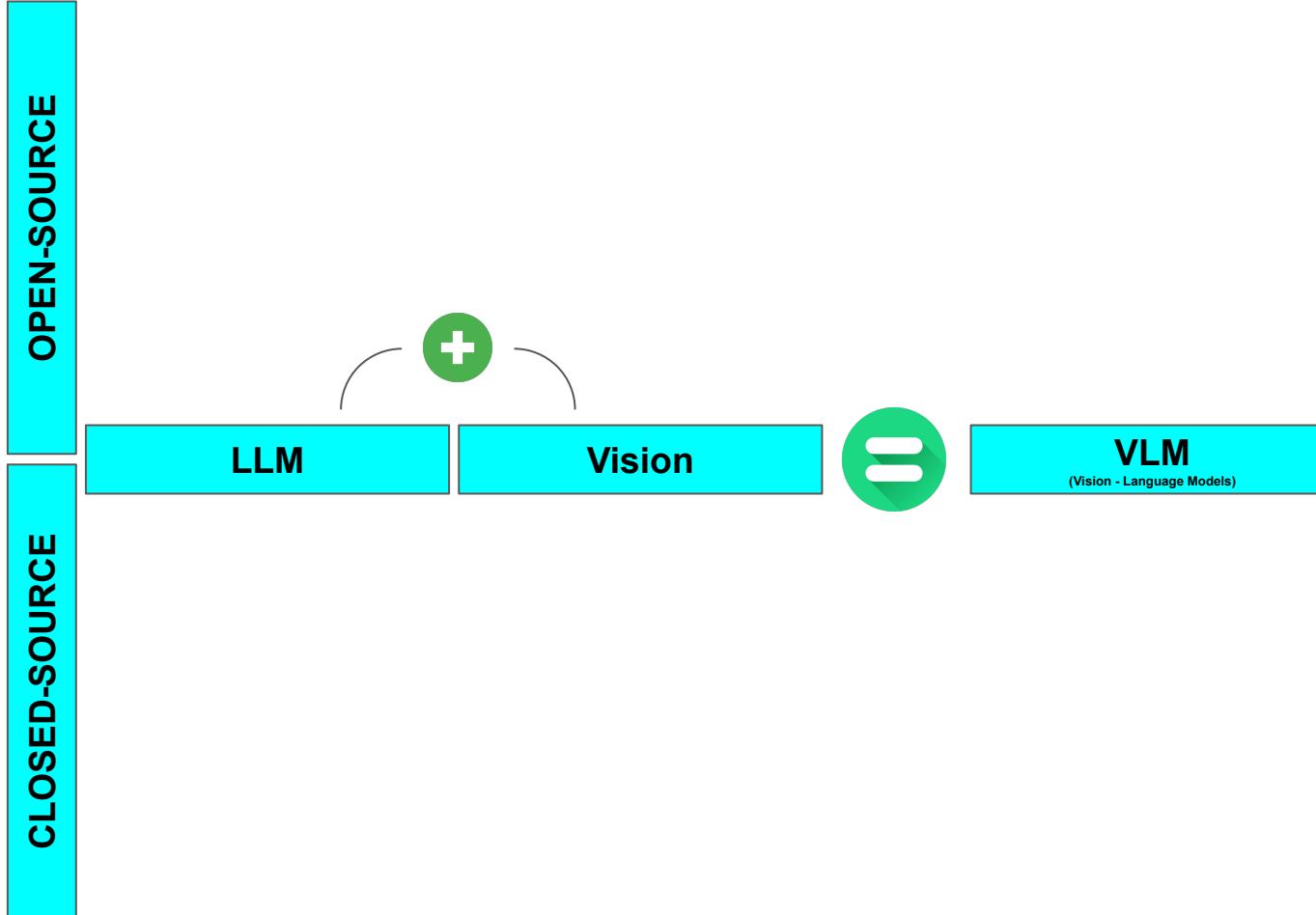
Nikko Carlo Yabut

From Sight To Insight

An Introduction to Vision Language Models



Vision Language Model





which country is this?

dubai, united arab emirates.

what made you think that?

a man is standing in front of the burj khalifa.

what is burj khalifa?

the tallest building in the world.

What is VLM?

- Models that understand **both images and text**
- Input can be:
 - Image + Text (e.g., "What is the cat doing?")
 - Text only (e.g., "Draw a panda surfing")
 - Image only (for captioning or tagging)

Sample Use Cases:

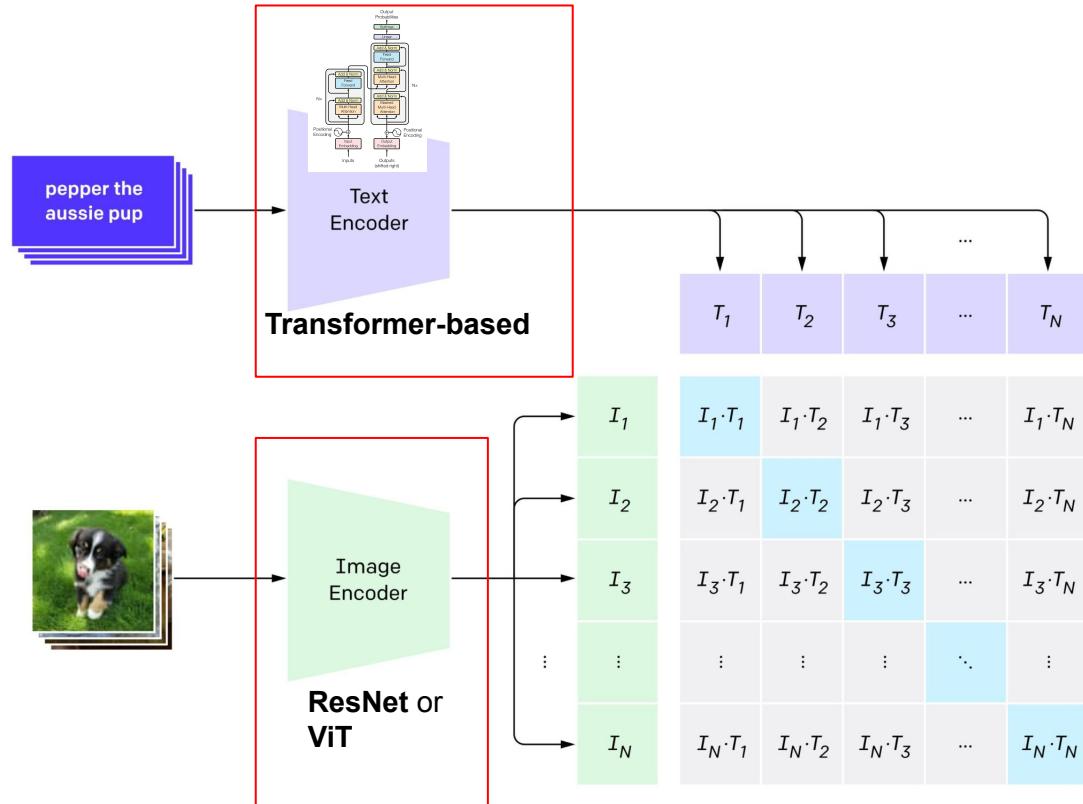
- **Image Captioning**
- **Visual Question Answering (VQA):**
 - Answering questions about images
- **Image-Text Retrieval:**
 - Finding images based on text queries
- **Text-to-Image Generation:**
 - Using prompts to create images

CLIP revolutionized vision-language AI

But how does it align images with text so effectively?

- The two modalities in CLIP are:

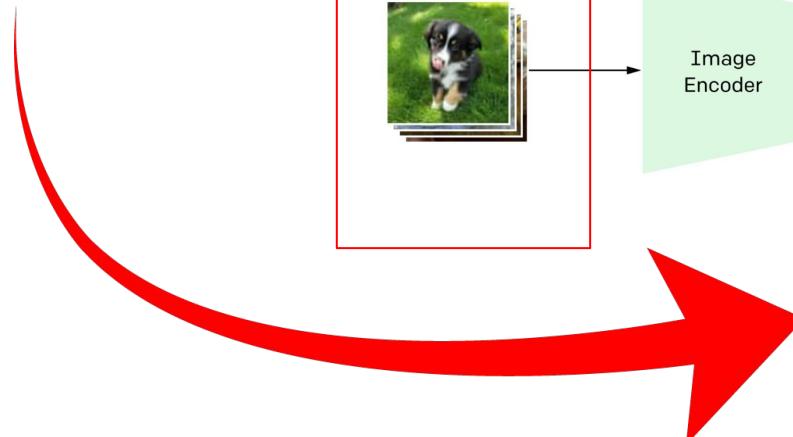
- Language** - which is transformer-based



CLIP revolutionized vision-language AI

But how does it align images with text so effectively?

- Trained on **400 million image-text pairs** collected from the internet
- Unlike traditional models, CLIP **does not rely on manually labeled datasets** like ImageNet — it's the **first model to learn from web-scale, naturally occurring image and text data**
- CLIP uses a **contrastive learning approach** — it learns by **bringing matching image-text pairs closer and pushing mismatched ones apart** in a shared embedding space



	I_1	$I_1 \cdot T_1$	$I_1 \cdot T_2$	$I_1 \cdot T_3$	\dots	$I_1 \cdot T_N$
	I_2	$I_2 \cdot T_1$	$I_2 \cdot T_2$	$I_2 \cdot T_3$	\dots	$I_2 \cdot T_N$
	I_3	$I_3 \cdot T_1$	$I_3 \cdot T_2$	$I_3 \cdot T_3$	\dots	$I_3 \cdot T_N$
\vdots	\vdots	\vdots	\vdots	\vdots	\ddots	\vdots
	I_N	$I_N \cdot T_1$	$I_N \cdot T_2$	$I_N \cdot T_3$	\dots	$I_N \cdot T_N$

CLIP revolutionized vision-language AI

But how does it align images with text so effectively?

- Trained on **400 million image-text pairs** collected from the internet
- Unlike traditional models, CLIP **does not rely on manually labeled datasets** like ImageNet — it's the **first model to learn from web-scale, naturally occurring image and text data**
- CLIP uses a **contrastive learning approach** — it learns by **bringing matching image-text pairs closer and pushing mismatched ones apart** in a shared embedding space

BEFORE TRAINING

Cosine similarity between text and image features



a person looking at a camera on a tripod-

a black-and-white silhouette of a horse-

a page of text about segmentation-

a portrait of an astronaut with the American flag-

a rocket standing on a launchpad-

a cup of coffee on a saucer-

a red motorcycle standing in a garage-

a facial photo of a tabby cat-

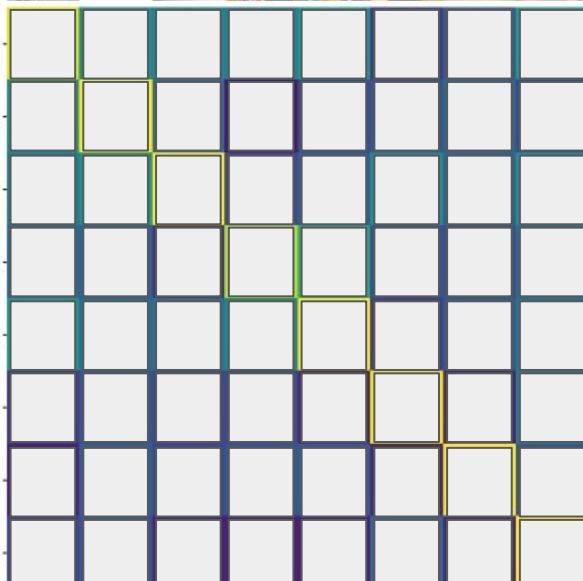


Figure 3: The Cosine similarity matrix between the text and image features pertaining to the `skimage` examples in the *Interacting with CLIP* colab notebook shared by CLIP's authors

CLIP revolutionized vision-language AI

But how does it align images with text so effectively?

- Trained on **400 million image-text pairs** collected from the internet
- Unlike traditional models, CLIP **does not rely on manually labeled datasets** like ImageNet — it's the **first model to learn from web-scale, naturally occurring image and text data**
- CLIP uses a **contrastive learning approach** — it learns by **bringing matching image-text pairs closer and pushing mismatched ones apart** in a shared embedding space

AFTER TRAINING

Cosine similarity between text and image features



Figure 3: The Cosine similarity matrix between the text and image features pertaining to the `skimage` examples in the *Interacting with CLIP* colab notebook shared by CLIP's authors

VLM Models

OPEN-SOURCE

CLOSED-SOURCE

 LLaMA 2 (by Meta)	 OpenLLAMA (by Berkeley AI)	 Gemma (by Google)	 Mixtral (by Mistral AI)	 Mistral (by Mistral AI)	SAM (by Meta AI)	MiniGPT-4	 Dolly (by Databricks)	Detectron2 (by Facebook AI R)	BLIP-2 (by Anthropic)	 Falcon (by TII)	MMDetection (by OpenMMLab)	BLIP (by Salesforce)	 DeepSeek (by DeepSeek AI)	YOLOv8 (by Ultralytics)	CLIP (by Open AI)	
LLM					Vision					VLM						
 GPT-4 (by OpenAI)	 Amazon Rekognition (by AWS)						 GPT-4V (by OpenAI)	 Claude 3 (by Anthropic)	 Claude 4 Opus (by Anthropic)	 Gemini (by Google DeepMind)	 Grok Vision (by xAI)	 Imagen	 Qwen (by Alibaba)	Midjourney	SenseChat (by SenseTime)	Flux
 Grok (by xAI)																

VLM Tasks

Discriminative Understanding Tasks

1. Image-Text Retrieval

Example: "Find an image that matches the phrase 'a dog wearing sunglasses.'"

2. Visual Question Answering

Example: "How many people are in the image?" → "Three"

3. Image Captioning

Example: Input: → Output: "A cat sitting on a laptop."

4. Visual Grounding / Referring Expression Comprehension

Example: "Where is the red backpack."

5. Multimodal Classification

Example: "Does this X-ray indicate pneumonia?"

What time of day is it?

night

noon



Cognitive Reasoning Tasks

6. Image-Text Matching

Determine whether a caption correctly describes an image.

Binary classification: "Yes" or "No"

7. Visual Commonsense Reasoning (VCR)

Given an image, question, and answer choices, select the most plausible one and justify it.

Example: "Why is the man running?" → "To catch the bus."

8. Multimodal Chain-of-Thought Reasoning

Provide step-by-step logical reasoning across modalities.



FVQA-style Commonsense:
Q9. Which object can be used for protecting head? helmet

Generative Tasks

9. Text-to-Image Generation

Generate an image from a textual prompt.

Example: "A futuristic city on Mars during sunset."

10. Image Inpainting / Editing

Modify an image based on a text prompt.

Example: "Remove the tree and add a bench."

11. Image Captioning with Style or Persona

Describe images with a tone/personality (e.g., sarcastic, poetic, child-like).



Utility Tasks

12. Document Visual QA / OCR-based Reasoning

Analyze documents (like receipts, invoices, or tables) using both layout and text.

Example: "What is the total price on the receipt?"

13. Chart QA / Figure Interpretation

Answer questions about bar charts, pie charts, line graphs, etc.

14. Multimodal Translation

Translate text within images from one language to another.

(a)		<pre>{"words": [{"text": "3602-Kyoto Choco Mochi"}, {"text": "14,000"}], "items": [{"name": "3602-Kyoto Choco Mochi", "count": 2, "price": 14000, "total": 28000}], "total": {"menuqty": 4, "total": 50000}}</pre>
(b)		<pre>{"words": [{"text": "3602-Kyoto Choco Mochi"}, {"text": "14,000"}], "items": [{"name": "3602-Kyoto Choco Mochi", "count": 2, "price": 14000, "total": 28000}], "total": {"menuqty": 4, "total": 50000}}</pre>
(d)		<pre>{"words": [{"text": "3602-Kyoto Choco Mochi"}, {"text": "14,000"}], "items": [{"name": "3602-Kyoto Choco Mochi", "count": 2, "price": 14000, "total": 28000}], "total": {"menuqty": 4, "total": 50000}}</pre>

<https://huggingface.co/spaces/stevengrove/YOLO-World>

Improving Open-Vocabulary Object Detection in a Vision Language Model

Anonymous ACCV 2024 Submission

Paper ID #94

Abstract. Vision Language Models (VLMs) are foundational models that provide meaningful answers to questions requiring image understanding. One of the many important tasks that VLMs can perform is open-vocabulary object detection (OVD), which involves detecting objects beyond the known classes the model was trained on. Despite training on extensive vision and language datasets, state-of-the-art (SOTA) VLMs such as LLaVA-v1.5, LLaVA-NEXT, and Qwen-VL-Chat achieve only 4.8, 14.4, and 18.6 AP respectively on the LVIS novel object detection benchmark. In this paper, we propose improving the OVD capability of LLaVA-v1.5 by fine tuning it on a new dataset generated automatically by FireLLaVA VLM and YOLO-World OVD. Our model, called LLaVA-World, achieves 25.0 AP on LVIS, a score that is state-of-the-art among VLMs and comparable to SOTA specialist OVD models such as OWL-ViT (25.6 AP) and YOLO-World (27.6 AP). LLaVA-World achieves this strong OVD score while maintaining the overall LLaVA performance on general language-vision tasks.

Keywords: OVD · Open-Vocabulary Object Detection · Vision Language Model



Human: What is going on here?

LLaVA-v1.5-13B: In the image, a man is driving a horse-drawn carriage down a city street. The carriage is pulled by a white horse, and the man is sitting on the carriage, guiding it. The scene appears to be set in an old city, possibly in a foreign country, with a building visible in the background. The man is also wearing a cowboy hat, which adds to the unique atmosphere of the scene. The horse-drawn carriage is a traditional mode of transportation that is not commonly seen in modern cities, making this image quite captivating and nostalgic.

Fig. 1: VLMs are capable of answering questions that require understanding the context of the image. In this case, a vague or ambiguous question can not be answered by LLM itself without “seeing” the image.



Fig. 2: ChatGPT-4o offers a poor object detection estimate on rare object **horse buggy**.



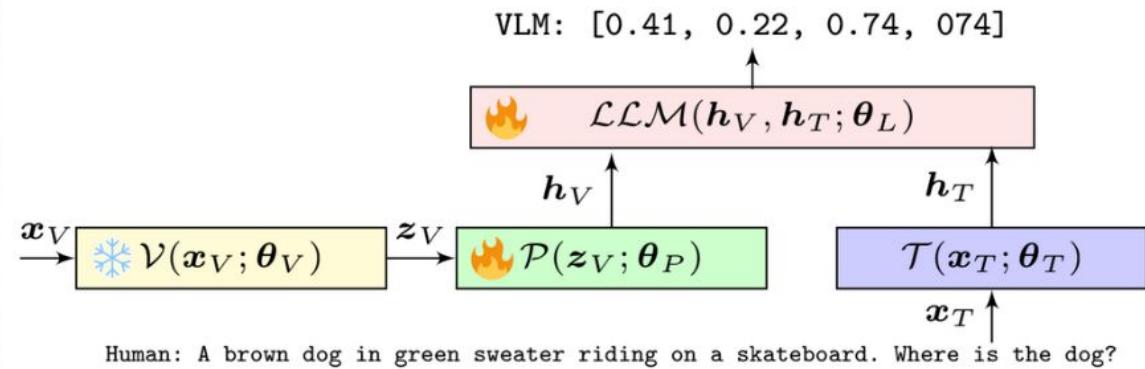
(a) LLaVA-v1.5-13B



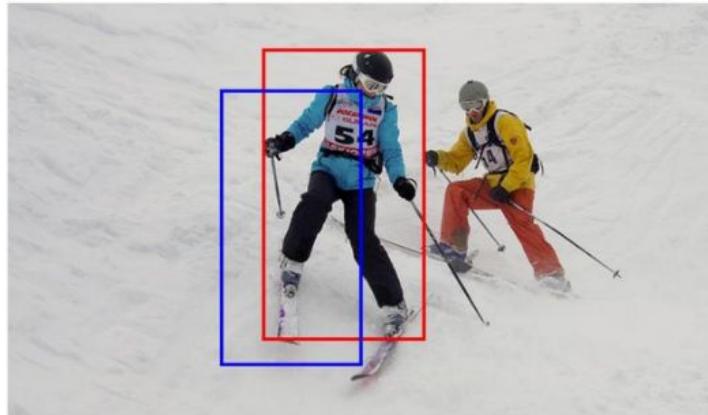
(b) LLaVA-World-13B (Ours)

Fig. 3: Comparison of novel object detection results between LLaVA-v1.5-13B and LLaVA-World-13B for rare object `horse buggy`.

Table 1: LLaVA-World fine-tuning recipe with a sample multi-turn conversation between human and VLM.



Role	Conversation
Human	A brown dog in green sweater riding on a skateboard. Where is the dog?
VLM	[0.41, 0.22, 0.74, 074]
Human	A cute little dog is on a skateboard. Locate the skateboard.
VLM	[0.37, 0.69, 0.89, 0.89]
Human	A small brown dog riding a skateboard with their owner. Find the cat.
VLM	There is no cat in the image.



Human: Locate the woman in a blue jacket skiing in the image. Output must be bounding box coordinates with the format [xmin, ymin, xmax, ymax].

LLaVA-v1.5^{7B}: [0.30, 0.21, 0.50, 0.87]

LLaVA-World^{7B}: [0.34, 0.13, 0.60, 0.85]

Fig. 6: An example of a referring expression comprehension evaluation using LLaVA-v1.5^{7B} (**blue**, vanilla) and LLaVA-World^{7B} (**red**, ours).

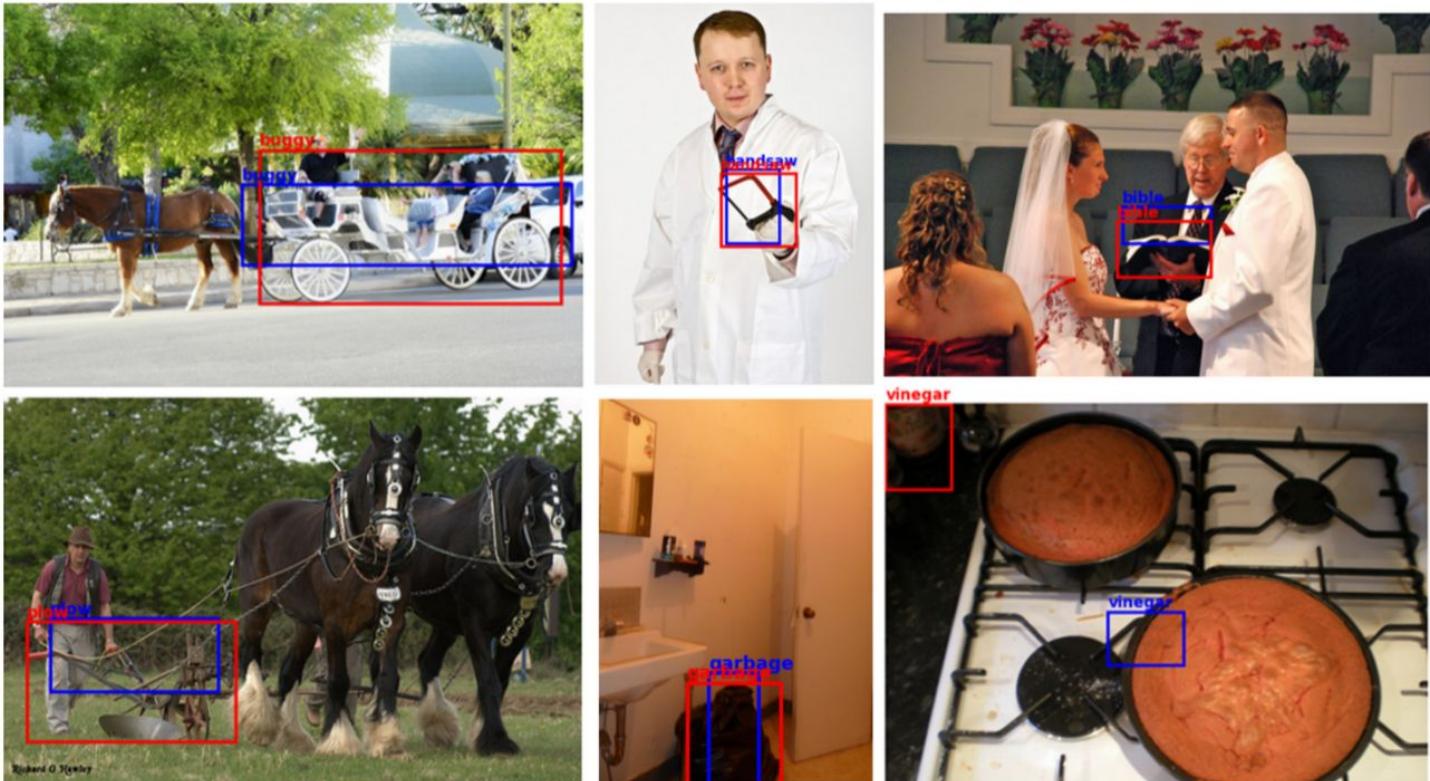


Fig. 7: Examples of annotated images with bounding boxes queried from LLaVA-1.5^{13B} (**blue**, vanilla) and LLaVA-World^{13B} (**red**, ours).

Nikko Carlo Yabut

Hearing Machines

An Introduction to Audio AI



From Sound to Signal – The Basics of Audio AI

- Key Concepts:
 - Audio = Time-series data (waveforms)
 - Spectrograms = Images of sound
 - Sampling rate (how often we record sound)
 - Features: MFCCs, Mel-Spectrograms

Sample Use Cases:

- **Audio Captioning**
 - describe what's heard
- **Spoken QA**
 - ask questions about spoken content
- **Voice Style Transfer**
 - Making one person sound like another
- **Multimodal assistants**
 - voice + vision + text

Audio AI Models

OPEN-SOURCE

CLOSED-SOURCE

 LLaMA 2 (by Meta)	 OpenLLAMA (by Berkeley AI)	 Gemma (by Google)	 Mixtral (by Mistral AI)	 Mistral (by Mistral AI)	SAM (by Meta AI)	MiniGPT-4	ESPnet
 Dolly (by Databricks)	 Detectron2 (by Facebook AI R)	 Falcon (by TII)	 MMDetection (by OpenMMLab)	 YOLOv8 (by Ultralytics)	BLIP-2 (by Anthropic)	BLIP (by Salesforce)	Wav2Vec
 DeepSeek (by DeepSeek AI)					CLIP (by Open AI)		Whisper
LLM	Vision	VLM	Audio				
 GPT-4 (by OpenAI)	 Amazon Rekognition (by AWS)	 GPT-4V (by OpenAI)	 ElevenLabs				
 Claude 3 (by Anthropic)		 Claude 4 Opus (by Anthropic)	 Google Speech-to-Text				
 Gemini (by Google DeepMind)		 Gemini (by Google DeepMind)	 Amazon Transcribe				
 Command R+ (by Cohere)		 Grok Vision (by xAI)	 Microsoft Azure Speech				
 Ernie Bot (by Baidu)		 Imagen	 Descript Overdub				
 Qwen (by Alibaba)		 Midjourney					
 SenseChat (by SenseTime)		 Flux					
 Grok (by xAI)							

Audio AI Tasks

Speech Tasks

Automatic Speech Recognition (ASR)

Converting spoken language into written text

Example: Transcribing podcasts or voice memos

Models: Whisper, Wav2Vec 2.0

Text-to-Speech (TTS)

Generating speech from text input

Example: Screen readers, virtual assistants

Models: Tortoise TTS, ElevenLabs

Speaker Diarization

Determining "who spoke when" in an audio stream

Example: Meeting transcription with speaker labels

Speaker Identification / Verification

Recognizing or confirming speaker identity

Example: Voice authentication for security

Voice Cloning

Creating synthetic voices based on a few samples of a real voice

Example: Personal voice assistants, voice preservation

Sound Processing

Audio Classification

Classifying audio into categories (e.g., music genre, environmental sound)

Example: Detecting dog barks vs. sirens in urban monitoring

Sound Event Detection (SED)

Locating and identifying specific sound events in an audio stream

Example: Identifying gunshots, baby cries, or alarms in surveillance

Acoustic Scene Classification

Recognizing the type of environment from sound (e.g., street, forest, cafe)

Example: Smart home or robotics context awareness

Emotion Recognition from Speech

Detecting emotional tone (e.g., happy, angry, sad)

Example: Call center analytics, mental health monitoring

Audio Denoising / Enhancement

Removing background noise or enhancing clarity

Example: Improving call quality, podcast editing

Generative

Music Generation / Composition

Creating original music or accompanying melodies

Example: AI-generated background music
Models: MusicLM, Riffusion

Voice Style Transfer / Singing Voice Conversion

Changing the style or identity of a voice without altering content

Example: Making one person sound like another

Audio-to-Audio Translation

Converting one type of audio to another (e.g., humming to instruments)

Multimodal Audio Tasks

Combining audio with other modalities (e.g., lip reading, audio-visual speech recognition)

Example: Video subtitles from lip + sound input

ESPnet

Wav2Vec

Whisper

Audio

ElevenLabs

Google Speech-to-Text

Amazon Transcribe

Microsoft Azure Speech

Descript Overdub

Models

OPEN-SOURCE

CLOSED-SOURCE

 LLaMA 2 (by Meta)	 OpenLLAMA (by Berkeley AI)	 Gemma (by Google)	 Mixtral (by Mistral AI)
 Mistral (by Mistral AI)	 SAM (by Meta AI)	 MiniGPT-4	
 Dolly (by Databricks)	 Detectron2 (by Facebook AI R)	 BLIP-2 (by Anthropic)	 ESPnet
 Falcon (by TII)	 MMDetection (by OpenMMLab)	 BLIP (by Salesforce)	 Wav2Vec
 DeepSeek (by DeepSeek AI)	 YOLOv8 (by Ultralytics)	 CLIP (by Open AI)	 Whisper
LLM	Vision	VLM	Audio
 GPT-4 (by OpenAI)	 Amazon Rekognition (by AWS)	 GPT-4V (by OpenAI)	 ElevenLabs
 Claude 3 (by Anthropic)		 Claude 4 Opus (by Anthropic)	 Google Speech-to-Text
 Gemini (by Google DeepMind)		 Gemini (by Google DeepMind)	 Amazon Transcribe
 Command R+ (by Cohere)		 Grok Vision (by xAI)	 Microsoft Azure Speech
 Ernie Bot (by Baidu)		 Imagen	 Descript Overdub
 Qwen (by Alibaba)		 Midjourney	
 SenseChat (by SenseTime)		 Flux	
 Grok (by xAI)			

Multimodal Mini-hackathon

Instruction:

In your assigned group, brainstorm and propose a business use case that creatively integrates multimodal AI—incorporating language, vision, and audio capabilities — to solve a real-world marketing problem.

Group 1

James
Cedric
Junielle
Bryan
Rexie
Clarence

Group 2

Noela
Justin
Cris
Carl
Pauline
Xy

Expected output

- Language, Audio, Vision, Video, etc...
- Landing page - use bolt.new
- Streamlit app - for basic functionality
- Pitch deck -
 - ChatGPT to create narrative
 - Gamma to create deck
- Presentation: 3:30PM

Resources:

- Multimodal streamlit playground deployment guide
- Github:
Week 4 - Multimodal Capabilities/ How to Deploy in Streamlit.pdf

Nikko Carlo Yabut

Foundations of Embeddings

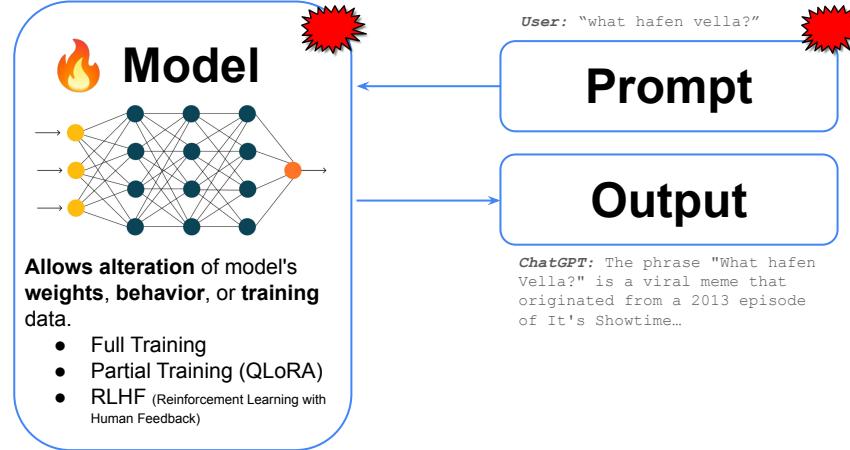
Task Adaptation

OPEN-SOURCE

 LLaMA 2 (by Meta)
 OpenLLAMA (by Berkeley AI)
 Gemma (by Google)
 Mixtral (by Mistral AI)
 Mistral (by Mistral AI)
 Dolly (by Databricks)
 Falcon (by TII)
 DeepSeek (by DeepSeek AI)

LLM

 GPT-4 (by OpenAI)
 Claude 3 (by Anthropic)
 Gemini (by Google DeepMind)
 Command R+ (by Cohere)
 Ernie Bot / Ernie 4.0 (by Baidu)
 Qwen (by Alibaba)
 SenseChat (by SenseTime)
 Grok (by xAI)



Prompting Techniques

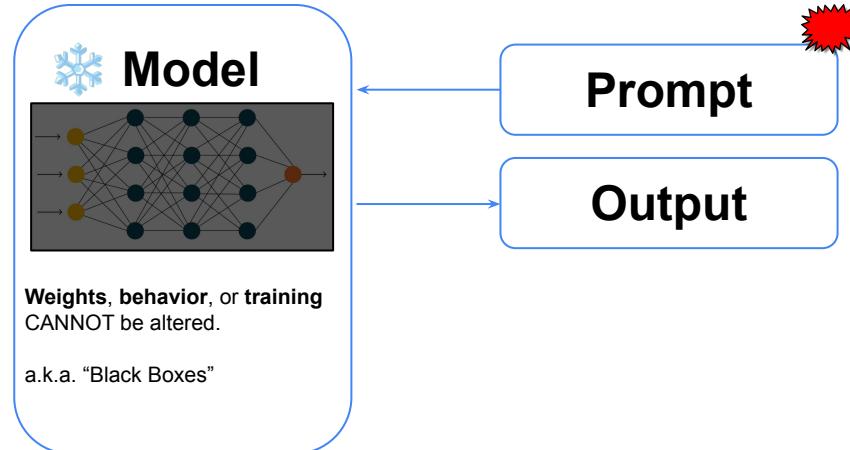
Zero-shot Prompting
One-shot Prompting
Few-shot Prompting
...

Prompt Frameworks

RICCE
PACE
APE

Prompting Techniques / Frameworks

*same as above



What's on your mind today?

If you were to build a chatbot that could answer questions about your company's documents, how would you do it?

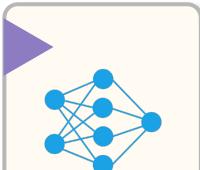
+ Tools



Fine-tuning as default strategy

"Take a base LLM like LLaMA, Mistral, or GPT-3. We'll just fine-tune it with our company's documents to make it smart."

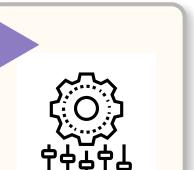
General steps in Fine-tuning



Model selection



Prepare data



Fine-Tune



Deploy

Problems With Fine-tuning



COST

RIGIDITY

REPETITION

VERSIONING

SLOW ITERATION

- Fine-tuning large models is expensive. Even open-source require GPU hours, storage, etc.

- Once a model is fine-tuned on a snapshot of your data, it doesn't generalize well to new content.

- Every time your knowledge base changes (e.g., new products, policies), you need to fine-tune again.

- You end up managing multiple model versions per department, per document version, etc.

- Even simple updates (e.g., adding a new FAQ) require model retraining. This breaks rapid prototyping.

What's on your mind today?

If you were to build a chatbot that could answer questions about your company's documents, how would you do it?

+ Tools



Better Approach:

"Instead of fine-tuning the model every time the knowledge base updates, what if we just store the documents in a database — and let the model retrieve the relevant chunks during inference?" !

DOCUMENTS stored in DATABASE

Model



Prompt

Output

What's on your mind today?

If you were to build a chatbot that could answer questions about your company's documents, how would you do it?

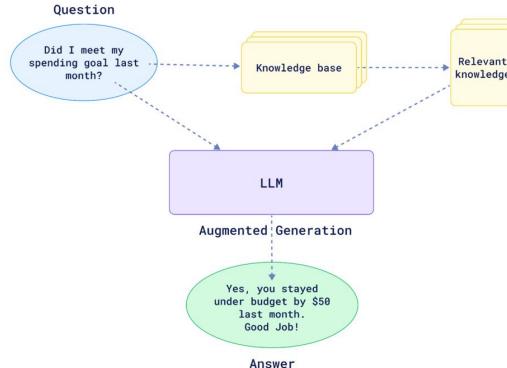
+ Tools



Better Approach:

Instead of fine-tuning the model every time the knowledge base updates, what if we just store the documents in a database — and let the model retrieve the relevant chunks during inference?

→ RAG



- Use a frozen, pretrained model (e.g., GPT or Mistral)
- Store your documents in a **vector database** (like FAISS, Weaviate, Qdrant, etc.)
- At runtime, convert the query and your documents into **embeddings** (aka latent representations)
- Retrieve relevant documents using **semantic similarity**
- Feed them into the LLM (via context window)

What's on your mind today?

If you were to build a chatbot that could answer questions about your company's documents, how would you do it?

+ Tools

0



Better Approach:

"Instead of fine-tuning the model every time the knowledge base updates, what if we just store the documents in a database — and let the model retrieve the relevant chunks during inference?"



RAG

"So instead of adapting the model to the documents, we adapt the documents to the model — using embeddings.



EMBEDDINGS?

RAG

“So instead of adapting the model to the documents, we adapt the documents to the model
— using embeddings.

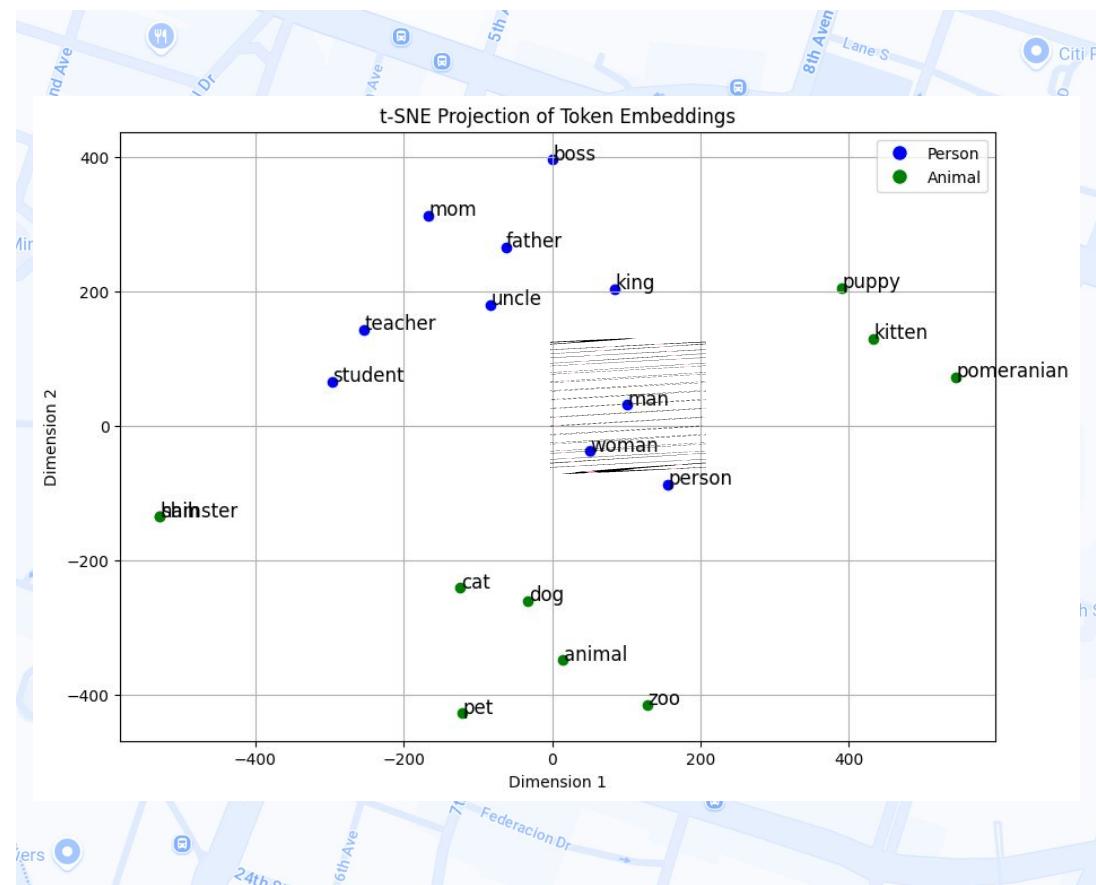
Embeddings represent words through bunch of numbers

Think of it like this:

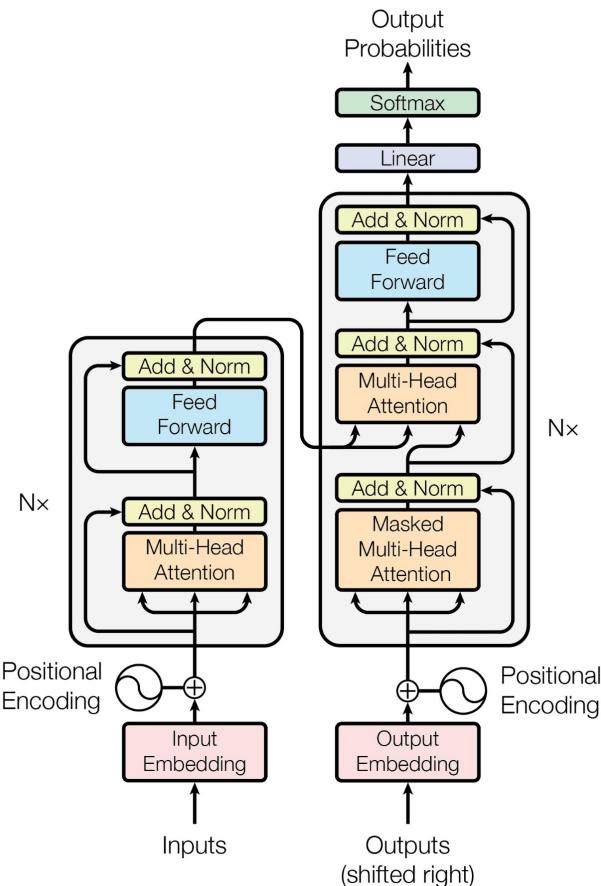
If words were cities
on a map, embeddings
are the coordinates.

Words with similar meanings, like
'man' and 'woman', end up close to
each other on that map.

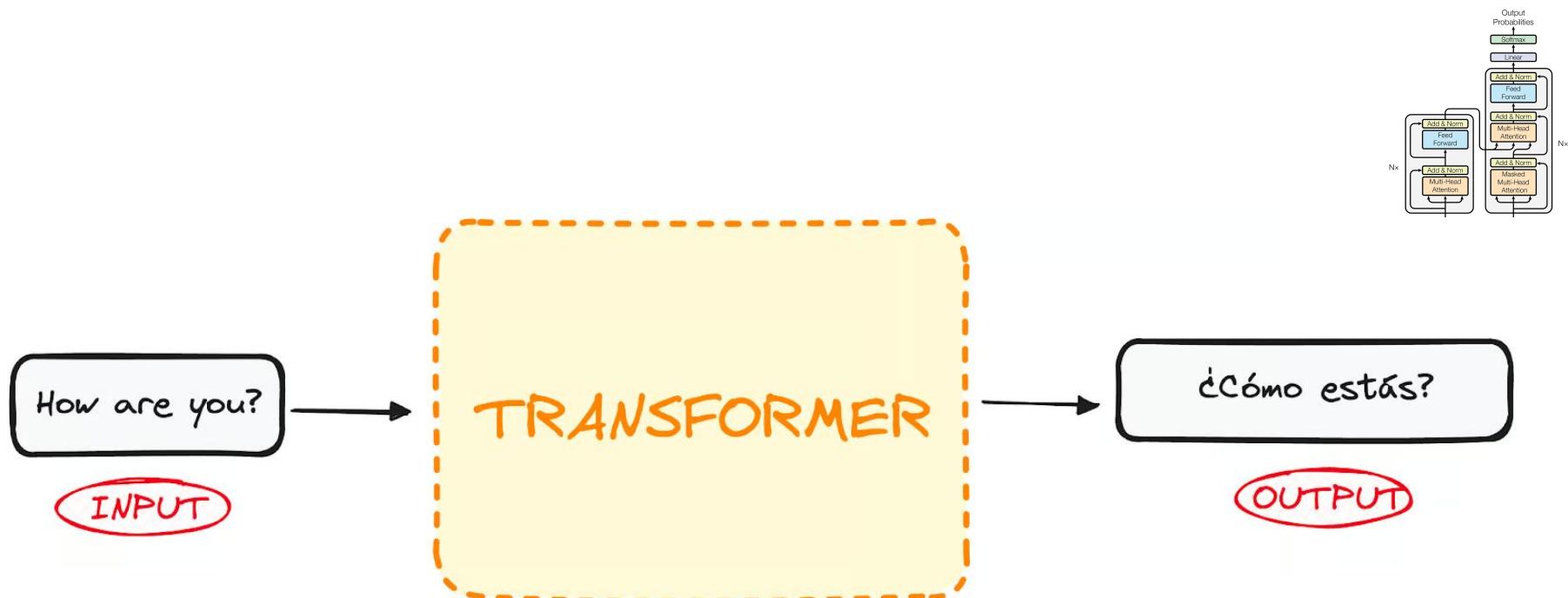
“



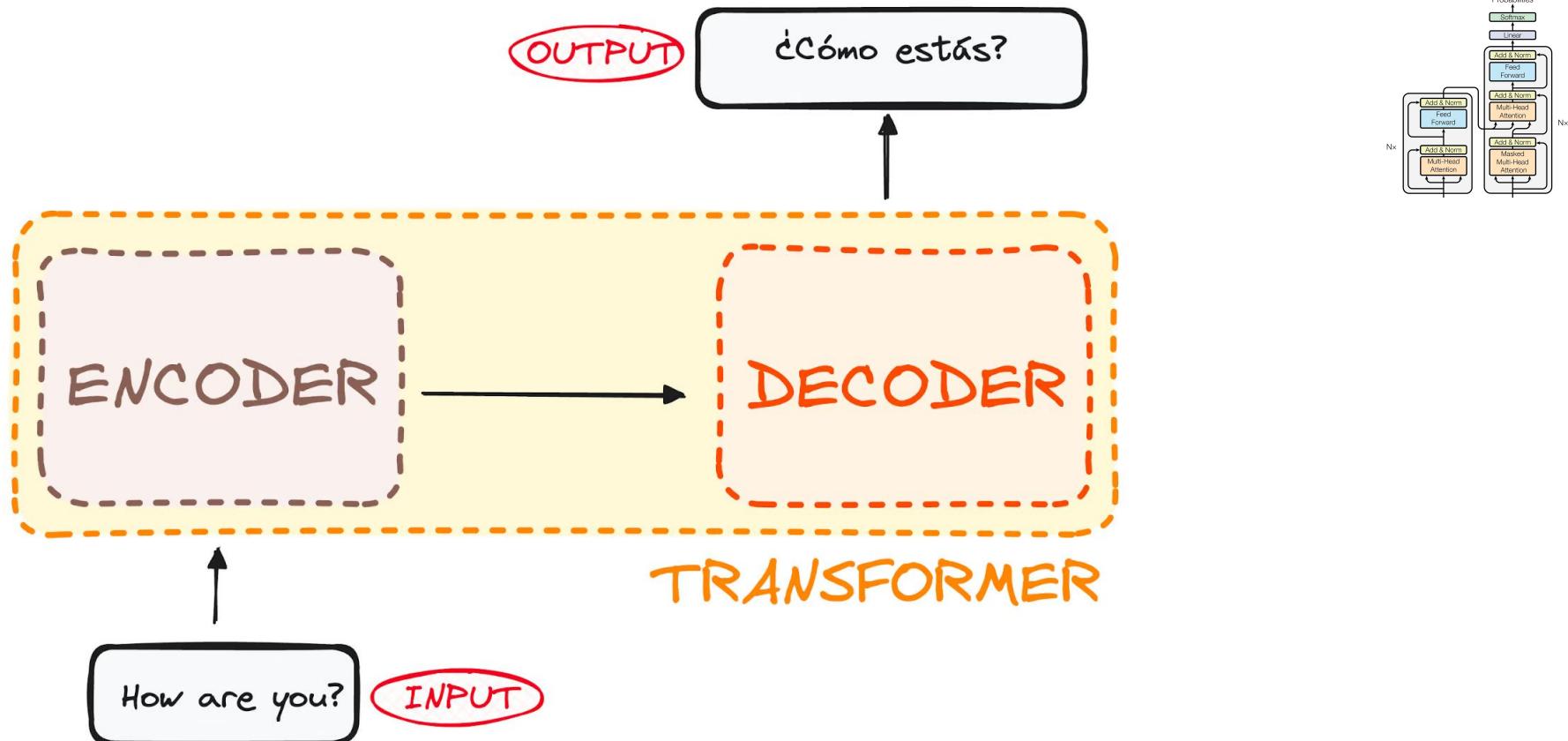
How are **Embeddings** created? (via transformers)



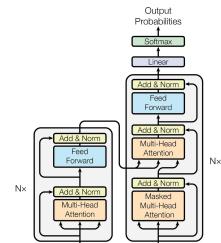
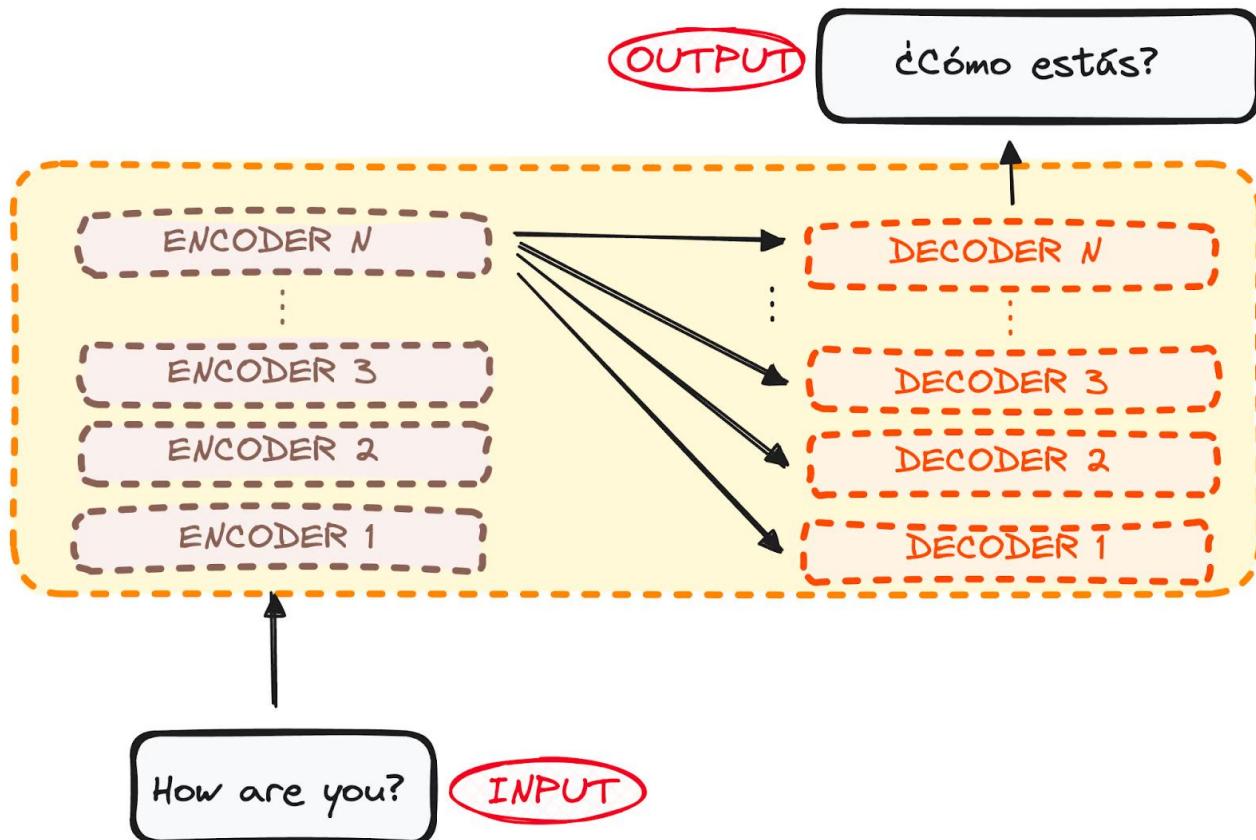
How are **Embeddings** created? (via transformers)



How are **Embeddings** created? (via transformers)

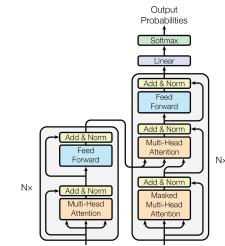
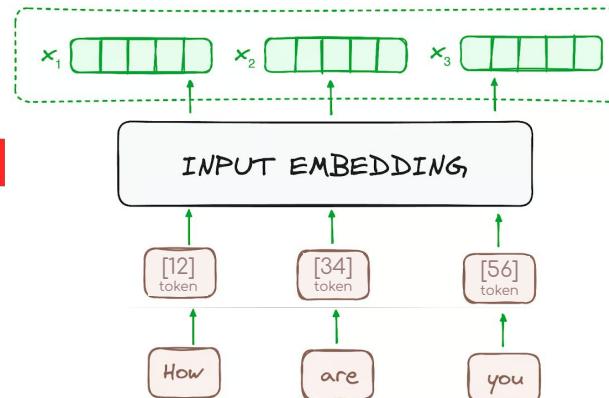
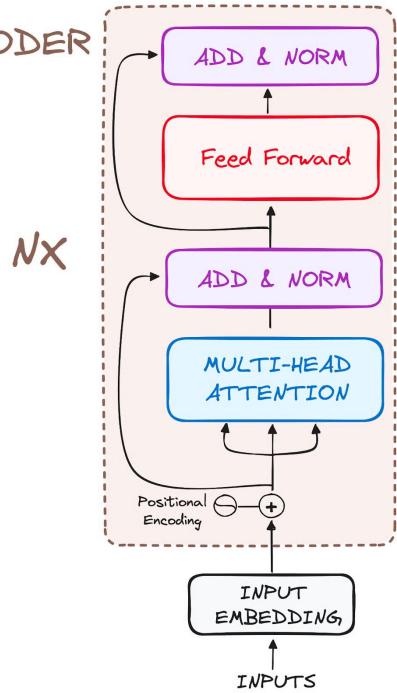


How are **Embeddings** created? (via transformers)



How are **Embeddings** created? (via transformers)

ENCODER



How are **Embeddings** created? (via transformers)

Word: How

Tokenized: ['how']

Token ID: 2129

Static Embedding Vector for 'How' (shape: (768,)): [REDACTED]

[REDACTED]

Word: are

Tokenized: ['are']

Token ID: 2024

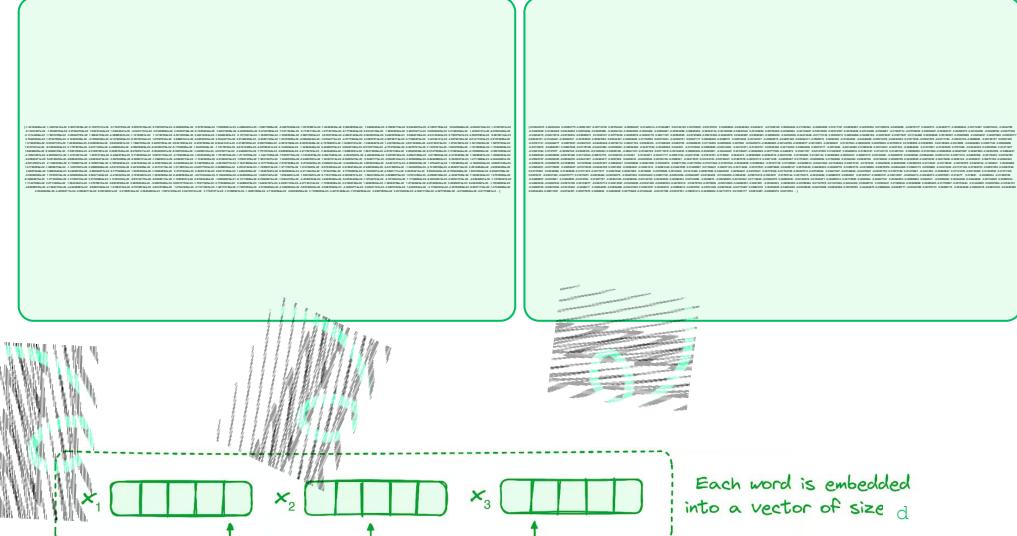
Static Embedding Vector (shape: (768,)): [REDACTED]

Word: you

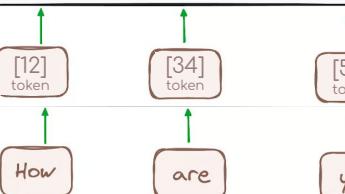
Tokenized: ['you']

Token ID: 2017

Static Embedding Vector (shape: (768,)): [REDACTED]

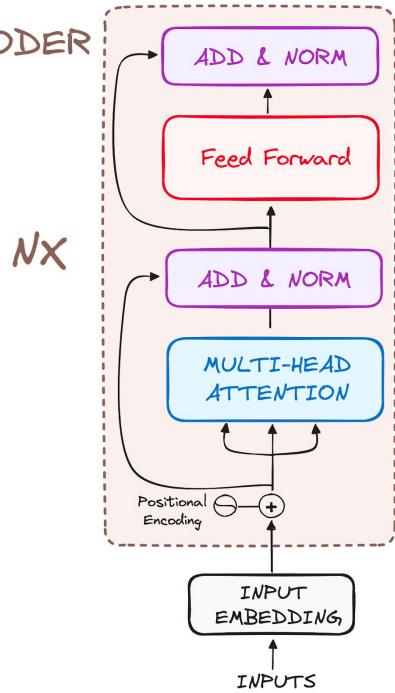


INPUT EMBEDDING,

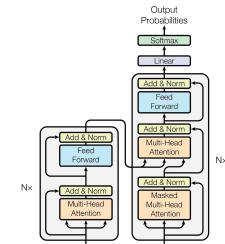
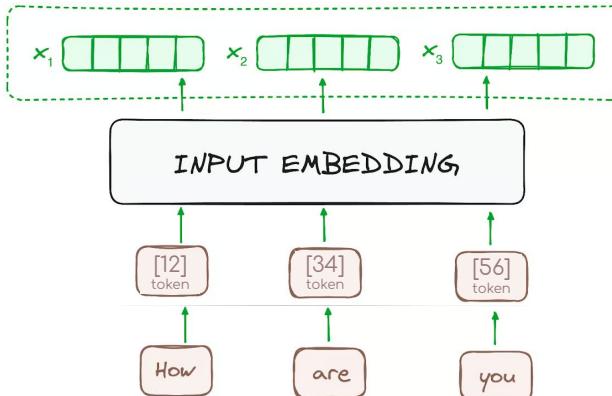


How are **Embeddings** created? (via transformers)

ENCODER

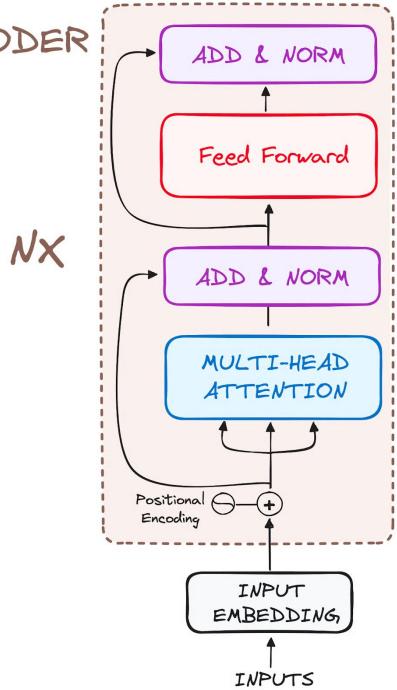


STATIC EMBEDDING.

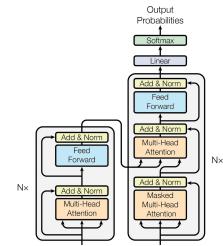


How are **Embeddings** created? (via transformers)

ENCODER

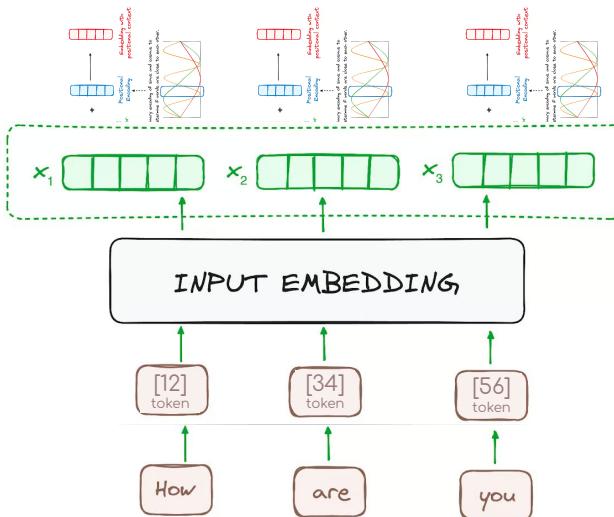


STATIC
EMBEDDING.



The **positional encoding** adds a "where it is" signal

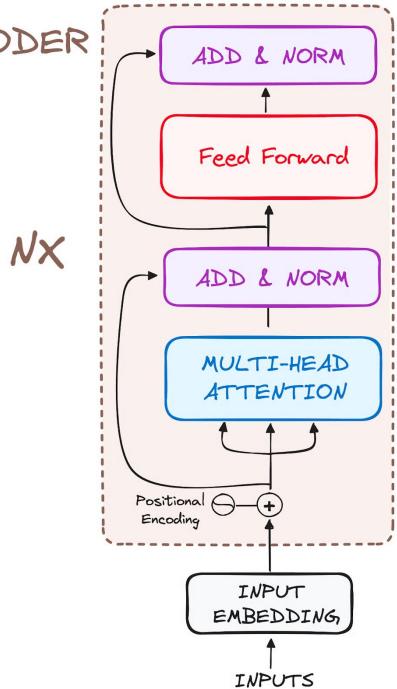
Two identical words at different positions get **different input vectors** because their embeddings are offset by different positional encodings



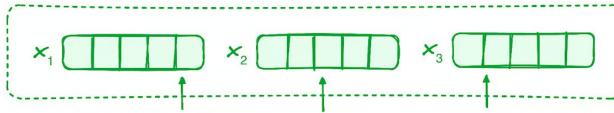
Each word is embedded into a vector of size d

How are **Embeddings** created? (via transformers)

ENCODER



STATIC EMBEDDING.

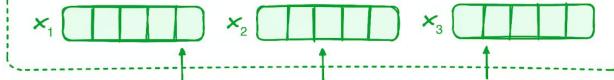
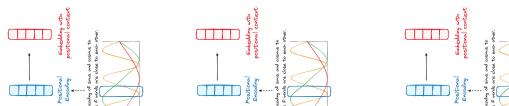


ENCODER

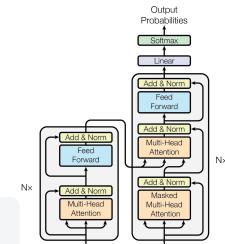
Nx

Encoder layer adds CONTEXT to the word.

The word embedding becomes **contextualized** — "bank" in "river bank" becomes different from "bank" in "money bank".

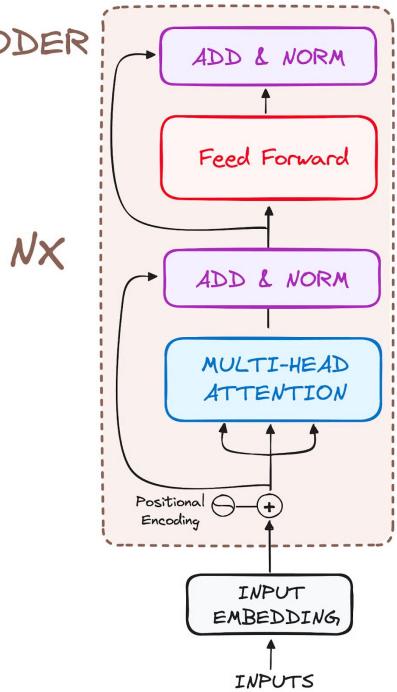


INPUT EMBEDDING,

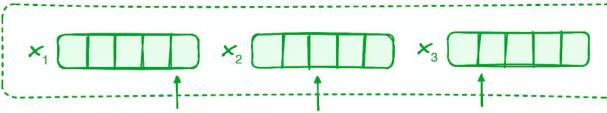


How are **Embeddings** created? (via transformers)

ENCODER



CONTEXTUAL EMBEDDING.



Nx

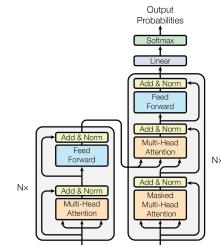
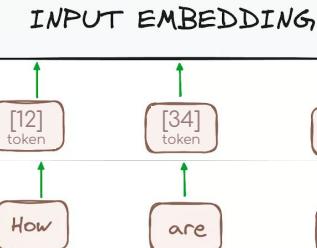
STATIC EMBEDDING.



ENCODER

Nx

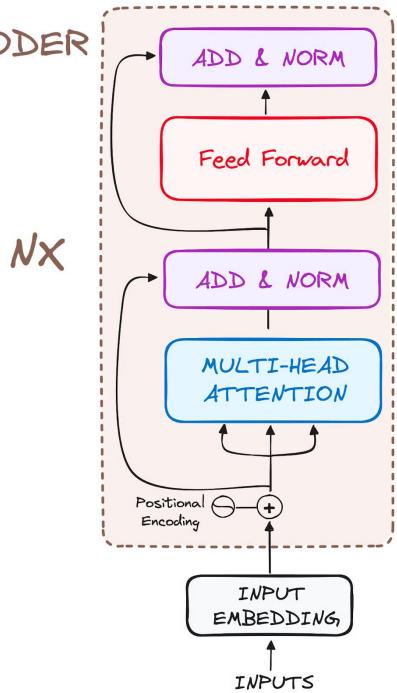
INPUT EMBEDDING,



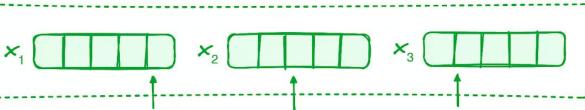
Each word is embedded into a vector of size d

How are **Embeddings** created? (via transformers)

ENCODER



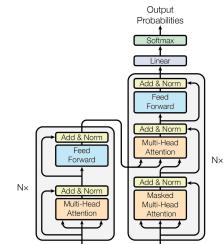
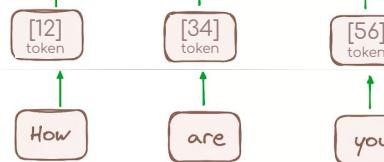
CONTEXTUAL EMBEDDING.



STATIC EMBEDDING.

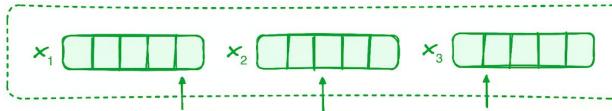


INPUT EMBEDDING,



Each word is embedded into a vector of size d

How are **Embeddings** created? (via transformers)



CONTEXTUAL EMBEDDING.

Comparison of Static vs Contextual Embeddings for the Word 'bank'

Input Sentence: He sat by the river bank.

'bank' Token ID: 2924

Static Embedding (first 5 dims): [-0.019130440428853035, -0.0645589530467987, -0.09128748625516891, -0.07761748880147934, -0.025318529456853867]

Contextual Embedding (first 5 dims): [0.1599486917257309, -0.3381432890892029, -0.03246789425611496, -0.08658606559038162, -0.39891570806503296]

Input Sentence: She deposited money at the bank.

'bank' Token ID: 2924

Static Embedding (first 5 dims): [-0.019130440428853035, -0.0645589530467987, -0.09128748625516891, -0.07761748880147934, -0.025318529456853867]

Contextual Embedding (first 5 dims): [0.393318772315979, -0.4198230504989624, -0.2934713363647461, 0.019461048766970634, 0.853203296661377]

Note: Static vs Contextual Embeddings

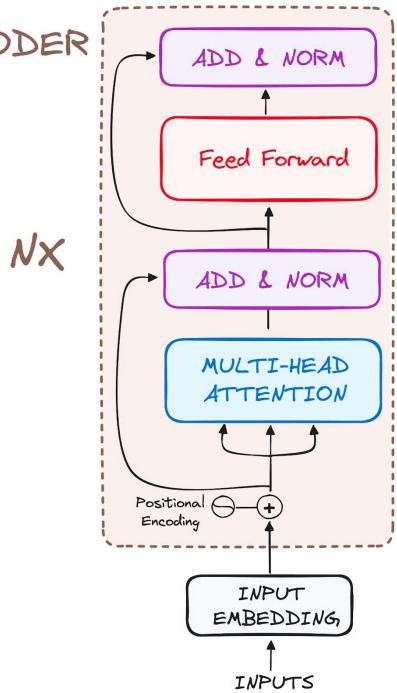
In both sentences, the word "bank" has the **same static embedding** because it's based only on the token itself.

But its **contextual embedding differs** — because BERT adjusts the vector based on surrounding words.

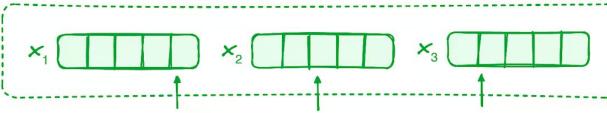
This shows how Transformers understand meaning based on **context**, not just the word.

How are **Embeddings** created? (via transformers)

ENCODER



CONTEXTUAL EMBEDDING.

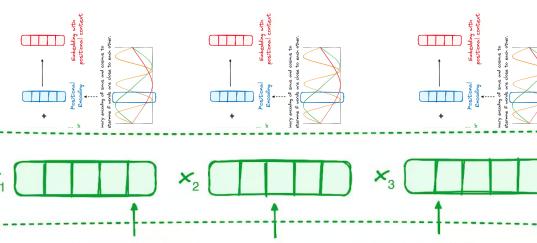


STATIC EMBEDDING.

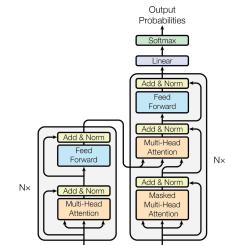


ENCODER

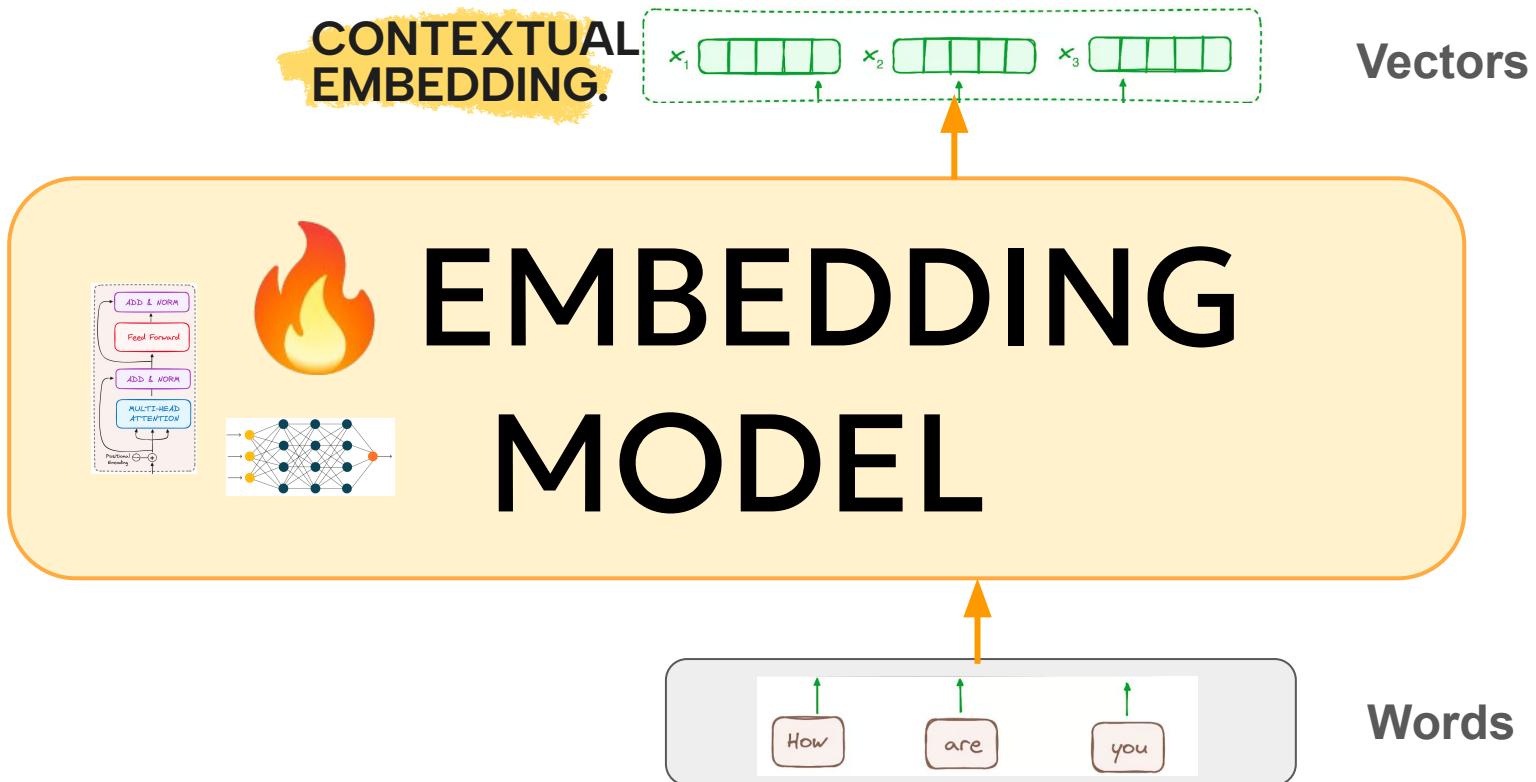
Nx



INPUT EMBEDDING,



All **Embedding Models** are Encoders under the hood



All **Embedding Models** are Encoders under the hood

OPEN-SOURCE

CLOSED-SOURCE

 LLaMA 2 (by Meta)	 OpenLLAMA (by Berkeley AI)	 Gemma (by Google)	 Mixtral (by Mistral AI)	 Mistral (by Mistral AI)	 SAM (by Meta AI)	 MiniGPT-4	 Word2Vec (by PyTorch)
 Dolly (by Databricks)	 Detectron2 (by Facebook AI R)	 BLIP-2 (by Anthropic)	 ESPnet	 Word2Vec (by Gensim)			
 Falcon (by TII)	 MMDetection (by OpenMMLab)	 BLIP (by Salesforce)	 Wav2Vec	 FastText (by Facebook)			
 DeepSeek (by DeepSeek AI)	 YOLOv8 (by Ultralytics)	 CLIP (by Open AI)	 Whisper	 SentenceTransformer (by Huggingface)			
LLM	Vision	VLM	Audio	Word2Vec			
 GPT-4 (by OpenAI)	 Amazon Rekognition (by AWS)	 GPT-4V (by OpenAI)	 ElevenLabs	 OpenAI Embeddings			
 Claude 3 (by Anthropic)		 Claude 4 Opus (by Anthropic)	 Google Speech-to-Text	 Google Cloud Embeddings			
 Gemini (by Google DeepMind)		 Gemini (by Google DeepMind)	 Amazon Transcribe	 Cohere Embeddings			
 Command R+ (by Cohere)		 Grok Vision (by xAI)	 Microsoft Azure Speech	 Amazon Comprehend			
 Ernie Bot (by Baidu)		 Imagen	 Descript Overdub				
 Qwen (by Alibaba)		 Midjourney					
 SenseChat (by SenseTime)		 Flux					
 Grok (by xAI)							

LLM vs Word2Vec (Embedding Models)

OPEN-SOURCE

 LLaMA 2 (by Meta)
 OpenLLAMA (by Berkeley AI)
 Gemma (by Google)
 Mixtral (by Mistral AI)
 Mistral (by Mistral AI)
 Dolly (by Databricks)
 Falcon (by TII)
 DeepSeek (by DeepSeek AI)

LLM

LLM

- Input: sequence of words
- Output: sequence of words

CLOSED-SOURCE

 GPT-4 (by OpenAI)
 Claude 3 (by Anthropic)
 Gemini (by Google DeepMind)
 Command R+ (by Cohere)
 Ernie Bot (by Baidu)
 Qwen (by Alibaba)
 SenseChat (by SenseTime)
 Grok (by xAI)

Word2Vec (Embedding Models)

- Input: word / sentences
- Output: vectors

Word2Vec (by PyTorch)

Word2Vec (by Gensim)

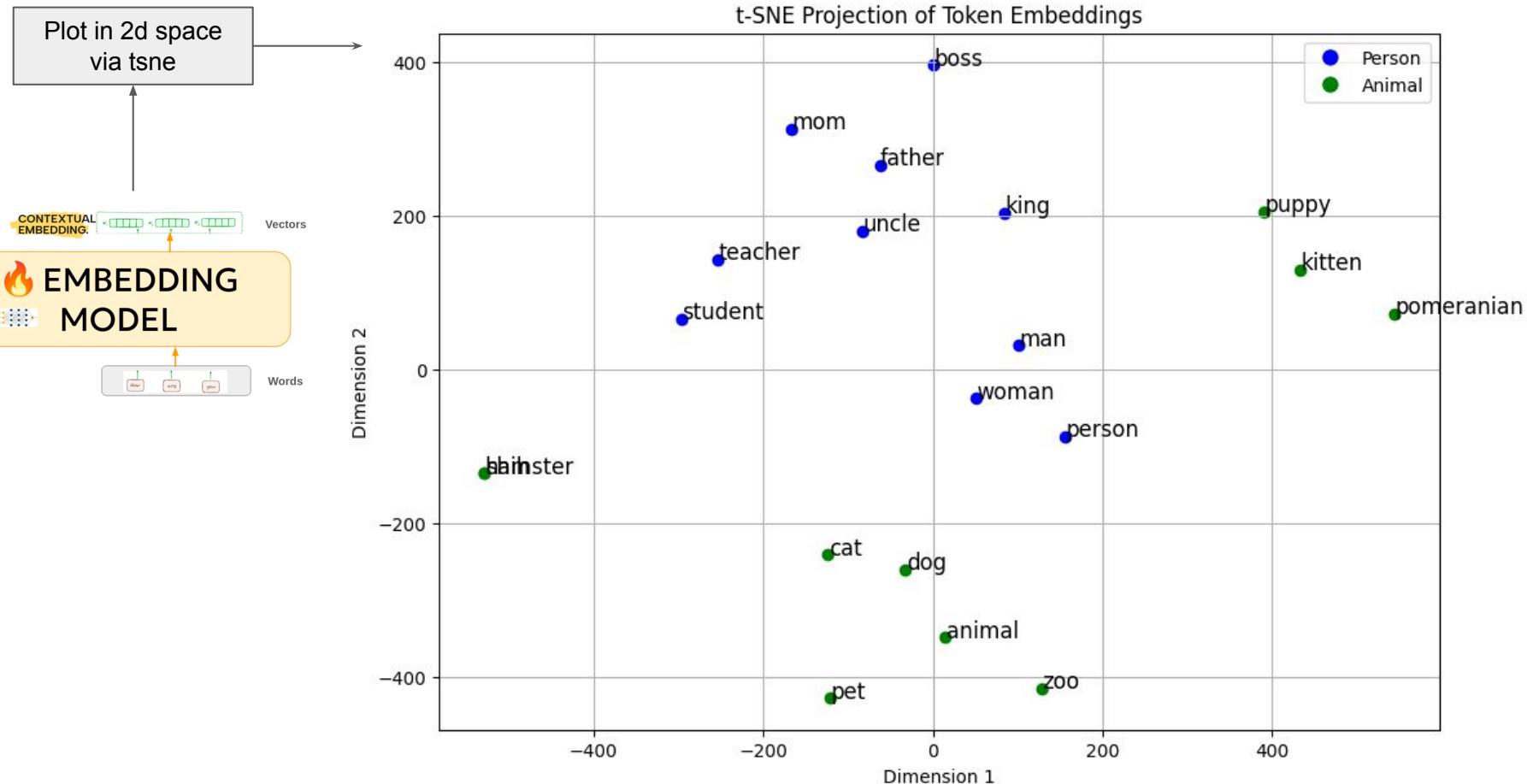
FastText (by Facebook)

SentenceTransformer (by Huggingface)

Word2Vec

OpenAI Embeddings
Google Cloud Embeddings
Cohere Embeddings
Amazon Comprehend

From context embedding (in 2d space)



From context embedding to **sentence embedding**

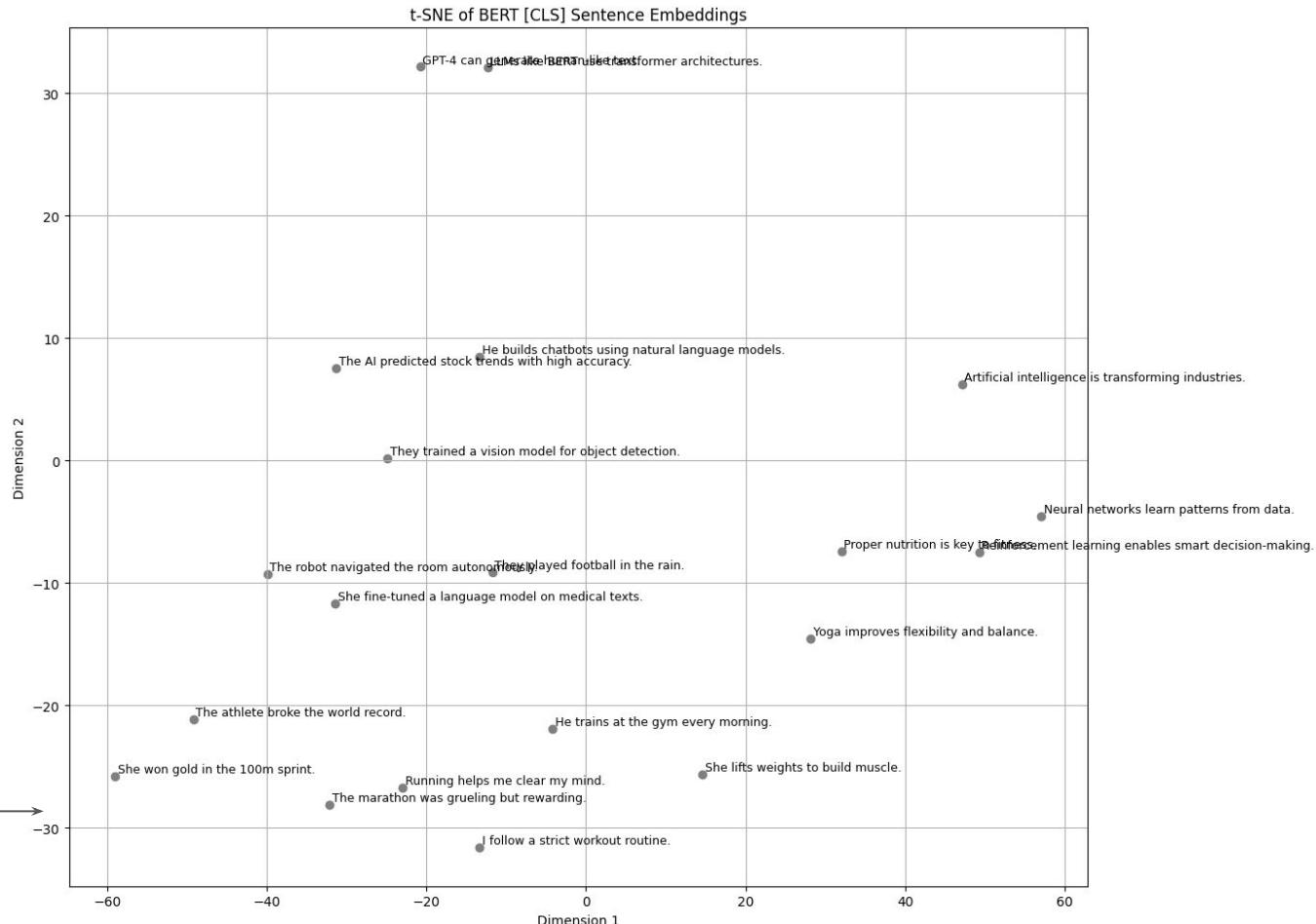
Let's extend the vectorization concept to one sentence.

Previously:
One word → 1 vector

In this slide:
1 sentence → 1 vector

```
sentences = [  
    "He trains at the gym every morning.",  
    "The marathon was grueling but rewarding.",  
    "She lifts weights to build muscle.",  
    "They played football in the rain.",  
    "Running helps me clear my mind.",  
    "The athlete broke the world record.",  
    "I follow a strict workout routine.",  
    "Yoga improves flexibility and balance.",  
    "She won gold in the 100m sprint.",  
    "Proper nutrition is key to fitness.",  
    "Artificial intelligence is transforming industries.",  
    "He builds chatbots using natural language models.",  
    "Neural networks learn patterns from data.",  
    "GPT-4 can generate human-like text.",  
    "The robot navigated the room autonomously.",  
    "They trained a vision model for object detection.",  
    "Reinforcement learning enables smart decision-making.",  
    "She fine-tuned a language model on medical texts.",  
    "The AI predicted stock trends with high accuracy.",  
    "LLMs like BERT use transformer architectures."  
]
```

Plot the embeddings in
2d space via tsne

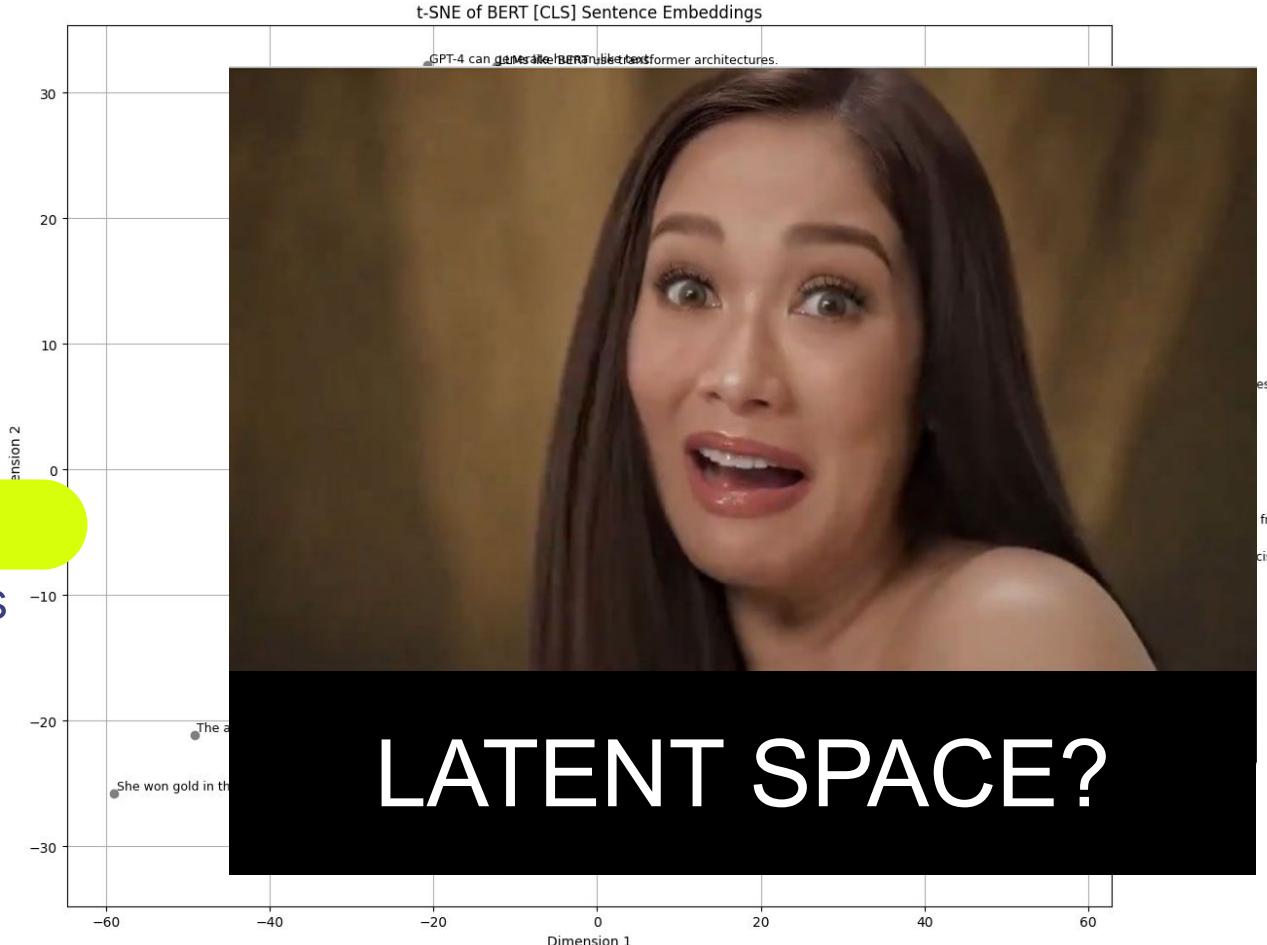


Embeddings live in the model's **latent space**

[Did you know]

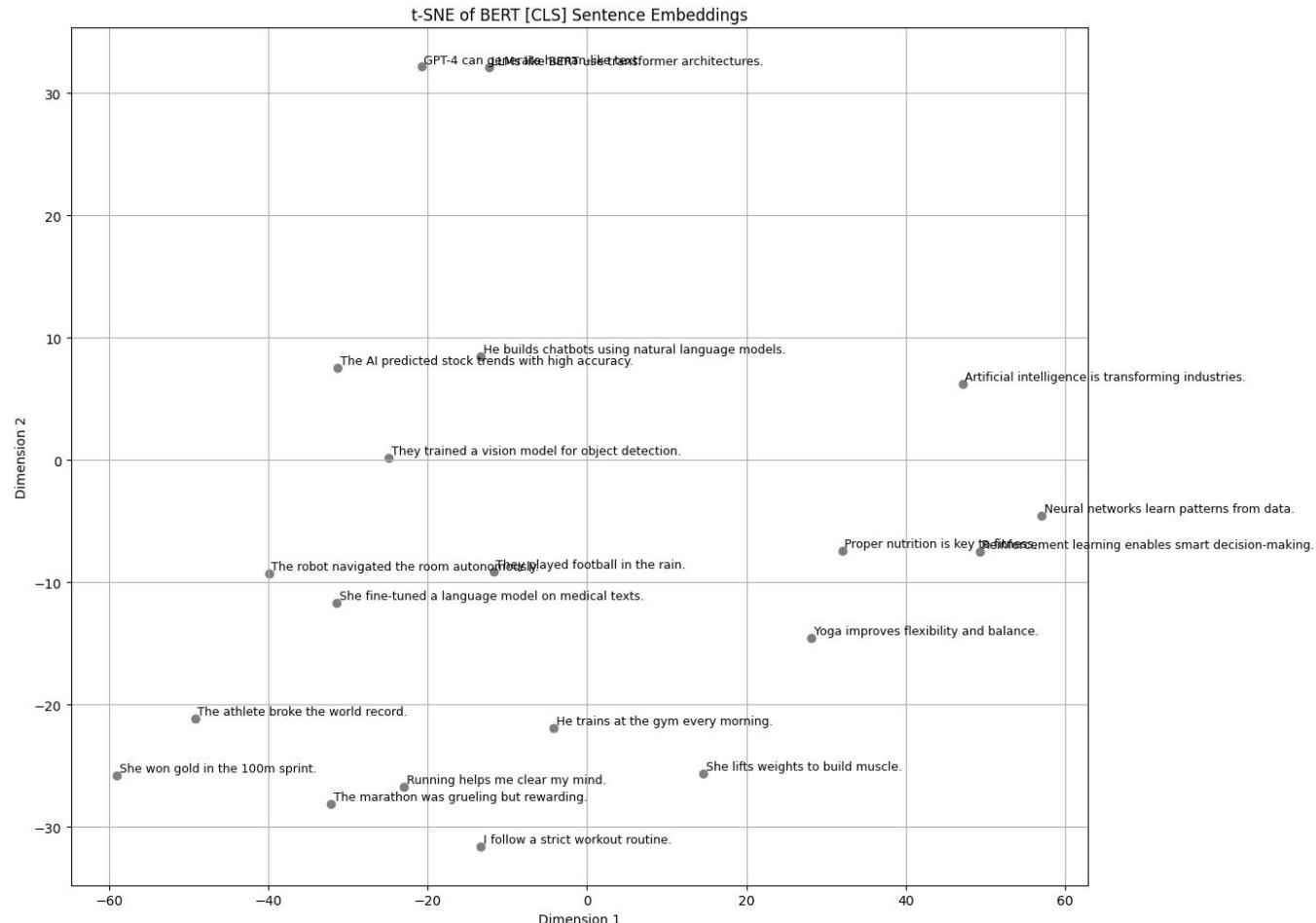
Embedding s live in the model's **latent space**

— where meaning is
encoded in vector
form.



Embeddings live in the model's **latent space** (in 2d)

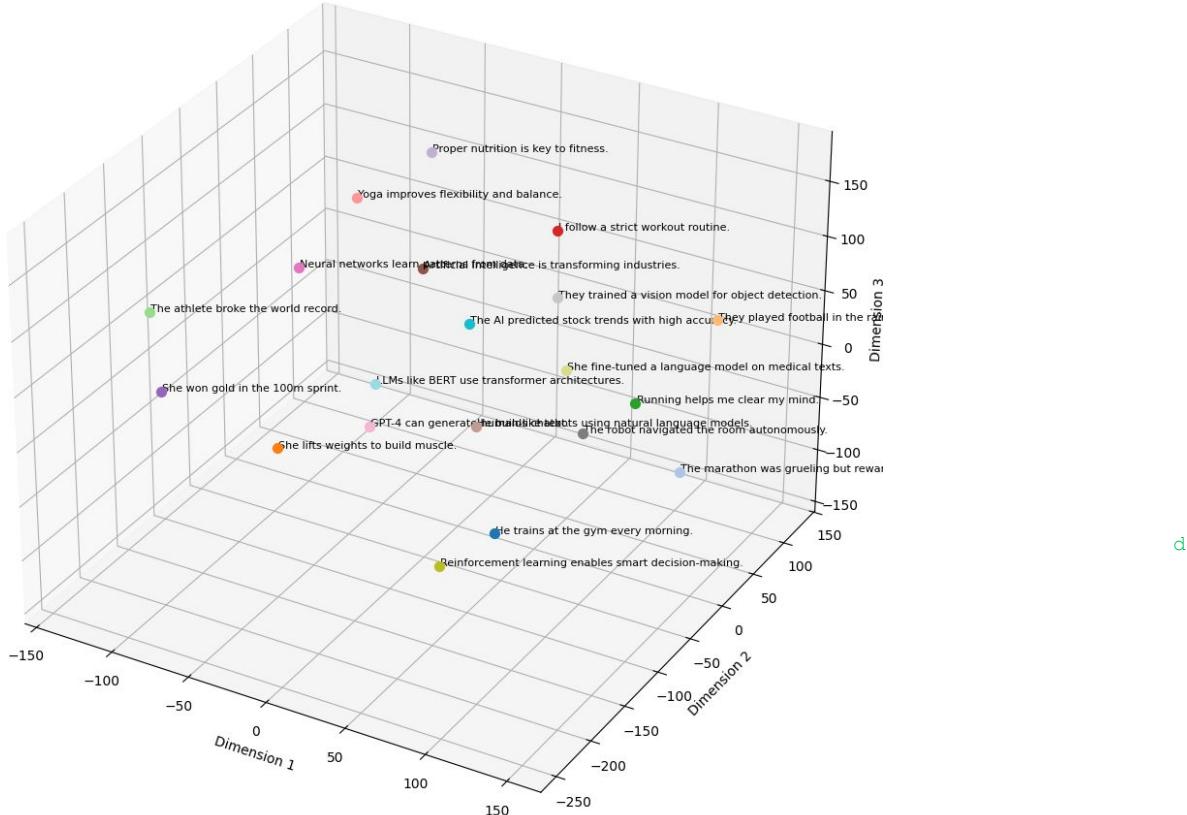
```
sentences = [
    "He trains at the gym every morning.",
    "The marathon was grueling but rewarding.",
    "She lifts weights to build muscle.",
    "They played football in the rain.",
    "Running helps me clear my mind.",
    "The athlete broke the world record.",
    "I follow a strict workout routine.",
    "Yoga improves flexibility and balance.",
    "She won gold in the 100m sprint.",
    "Proper nutrition is key to fitness.",
    "Artificial intelligence is transforming industries.",
    "He builds chatbots using natural language models.",
    "Neural networks learn patterns from data.",
    "GPT-4 can generate human-like text.",
    "The robot navigated the room autonomously.",
    "They trained a vision model for object detection.",
    "Reinforcement learning enables smart decision-making.",
    "She fine-tuned a language model on medical texts.",
    "The AI predicted stock trends with high accuracy.",
    "LLMs like BERT use transformer architectures."
]
```



Embeddings live in the model's **latent space** (in 3d)

3D t-SNE of Sentence Embeddings (Colored by Sentence Index)

```
sentences = [
    "He trains at the gym every morning.",
    "The marathon was grueling but rewarding.",
    "She lifts weights to build muscle.",
    "They played football in the rain.",
    "Running helps me clear my mind.",
    "The athlete broke the world record.",
    "I follow a strict workout routine.",
    "Yoga improves flexibility and balance.",
    "She won gold in the 100m sprint.",
    "Proper nutrition is key to fitness.",
    "Artificial intelligence is transforming industries.",
    "He builds chatbots using natural language models.",
    "Neural networks learn patterns from data.",
    "GPT-4 can generate human-like text.",
    "The robot navigated the room autonomously.",
    "They trained a vision model for object detection.",
    "Reinforcement learning enables smart decision-making.",
    "She fine-tuned a language model on medical texts.",
    "The AI predicted stock trends with high accuracy.",
    "LLMs like BERT use transformer architectures."
]
```



Embeddings live in the model's **latent space (in embedding Dimension)**

```
sentences = [  
    "He trains at the gym every morning.",  
    "The marathon was grueling but rewarding.",  
    "She lifts weights to build muscle.",  
    "They played football in the rain.",  
    "Running helps me clear my mind.",  
    "The athlete broke the world record.",  
    "I follow a strict workout routine.",  
    "Yoga improves flexibility and balance.",  
    "She won gold in the 100m sprint.",  
    "Proper nutrition is key to fitness.",  
    "Artificial intelligence is transforming industries.",  
    "He builds chatbots using natural language models.",  
    "Neural networks learn patterns from data.",  
    "GPT-4 can generate human-like text.",  
    "The robot navigated the room autonomously.",  
    "They trained a vision model for object detection.",  
    "Reinforcement learning enables smart decision-making.",  
    "She fine-tuned a language model on medical texts.",  
    "The AI predicted stock trends with high accuracy.",  
    "LLMs like BERT use transformer architectures."  
]
```

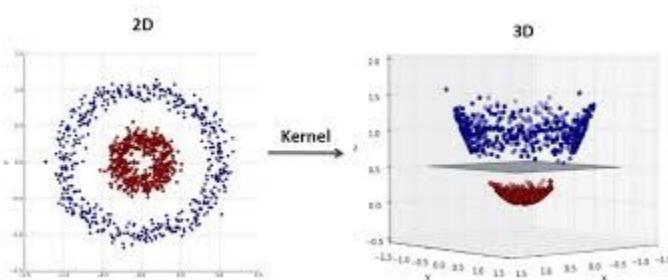


Model	Hidden Size (Embedding Dim)
BERT-base	768
BERT-large	1024
DistilBERT	768
RoBERTa-base	768
RoBERTa-large	1024
GPT-2 (small)	768
GPT-2 (medium)	1024
GPT-2 (large)	1280
GPT-2 (XL)	1600
GPT-3 (Davinci)	12288

d

Why is 2d or 3d not enough?

```
sentences = [
    "He trains at the gym every morning.",
    "The marathon was grueling but rewarding.",
    "She lifts weights to build muscle.",
    "They played football in the rain.",
    "Running helps me clear my mind.",
    "The athlete broke the world record.",
    "I follow a strict workout routine.",
    "Yoga improves flexibility and balance.",
    "She won gold in the 100m sprint.",
    "Proper nutrition is key to fitness.",
    "Artificial intelligence is transforming industries.",
    "He builds chatbots using natural language models.",
    "Neural networks learn patterns from data.",
    "GPT-4 can generate human-like text.",
    "The robot navigated the room autonomously.",
    "They trained a vision model for object detection.",
    "Reinforcement learning enables smart decision-making.",
    "She fine-tuned a language model on medical texts.",
    "The AI predicted stock trends with high accuracy.",
    "LLMs like BERT use transformer architectures."
]
```



Model	Hidden Size (Embedding Dim)
BERT-base	768
BERT-large	1024
DistilBERT	768
RoBERTa-base	768
RoBERTa-large	1024
GPT-2 (small)	768
GPT-2 (medium)	1024
GPT-2 (large)	1280
GPT-2 (XL)	1600
GPT-3 (DaVinci)	12288

Borrowing the concept of the kernel trick, we can explain why embeddings work: **by moving text into a higher-dimensional space, it's easier to separate meanings** just like how the kernel trick separates non-linear patterns by projecting data into a space where it's linearly separable."

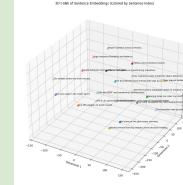
Embeddings live in the model's **latent space**

[Did you know]

**Embedding
s live in the
model's
latent space**

— where meaning is
encoded in vector
form.

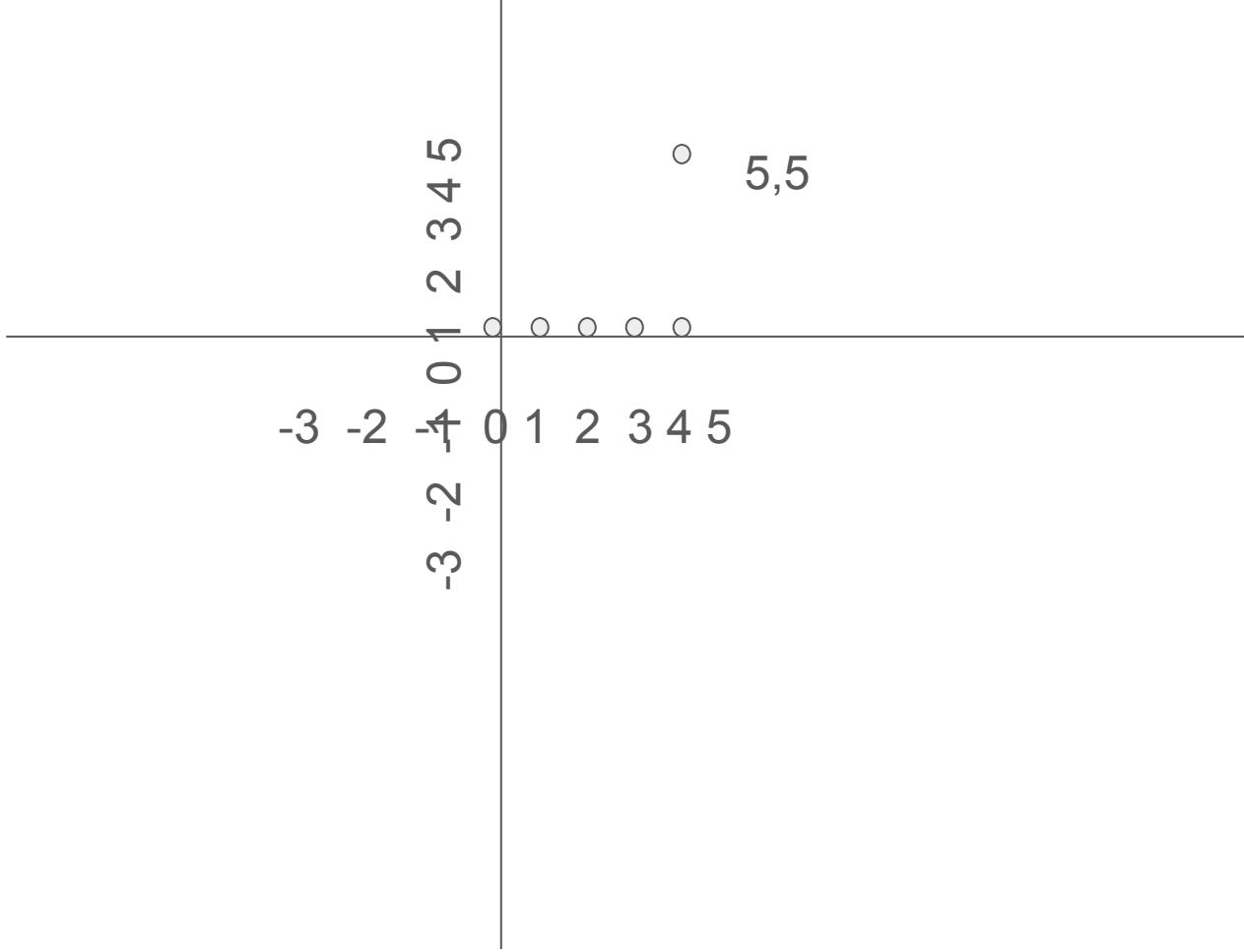
Latent Space



🔥 **EMBEDDING
MODEL**

Activity

https://colab.research.google.com/drive/1f8Bcqo_D0flkrSzIYXVCvN1UISZphBJ4#scrollTo=ouzWVBQjz3tT



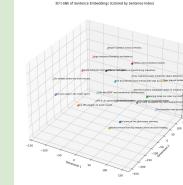
Embeddings live in the model's **latent space**

[Did you know]

**Embedding
s live in the
model's
latent space**

— where meaning is
encoded in vector
form.

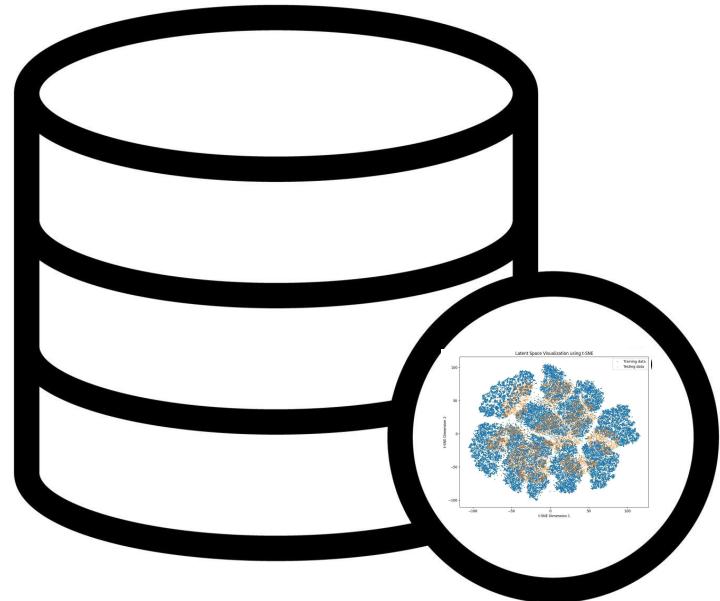
Latent Space



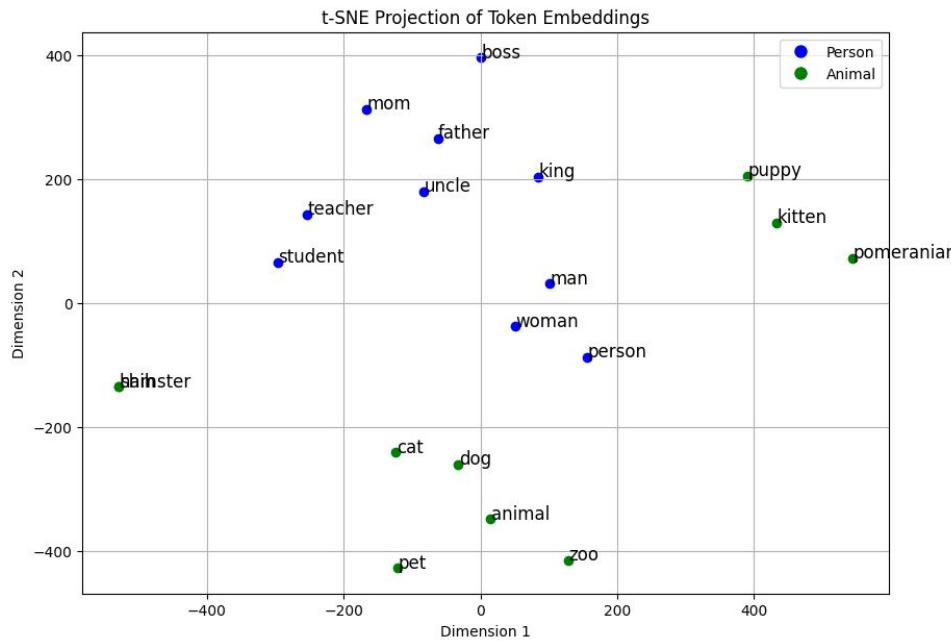
🔥 **EMBEDDING
MODEL**

Nikko Carlo Yabut

Saving the Latent Space in a Vector Database



Saving the Latent Space in a Vector Database



Latent Space

EMBEDDING
MODEL

king

**Vector Database
(Latent Space)**



**EMBEDDING
MODEL**

king

Saving the Latent Space in a Vector Database

Name	Open Source	Scalable	Cloud-hosted	Notes
FAISS	✓	✗	✗	Local-only, blazing fast
Qdrant	✓	✓	✓	Rust-based, modern APIs
Pinecone	✗	✓	✓	API-only, closed-source
Weaviate	✓	✓	✓	ML pipelines, hybrid search
Chroma	✓	✗	✗	Lightweight, local-only
Vespa	✓	✓	✓	Enterprise-grade, complex

Vector Database (Latent Space)



EMBEDDING MODEL

king

Saving the Latent Space in a Vector Database

The screenshot shows the Qdrant Cloud interface. On the left, a sidebar menu includes 'DASHBOARD', 'Clusters' (selected), 'Hybrid Cloud', and 'Backups'. Under 'ACCOUNT', there are 'Access Management', 'Billing', and 'Settings'. At the bottom, there's a 'Get Support' link. The main content area shows a cluster named 'ai_first_rag_demo' with status 'HEALTHY' and 'FREE TIER'. The cluster details include Cluster ID: 01ab3875-7bf8-470a-bd47-bbebed07c3ed, Version: v1.14.1, Cloud Provider: europe-west3, and Node specifications: 1 node, 4GiB disk, 1GiB RAM, 0.5 vCPU. A section titled 'Get All Qdrant Cloud Features' lists benefits like dedicated resources, monitoring, and disaster recovery. Below this is a table showing resource usage: NODE #0 v1.14.1, DISK 368.00 KiB of 4.00 GiB, RAM 147.85 MiB of 1.00 GiB, VCPU 0.04 vCPUs of 0.5 vCPUs, and ACTIONS. To the right, sections for 'Access the Cluster', 'Use the API', and 'Examples' are visible.

Nikko Carlo Yabut - Base Account ▾

Nikko Carlo Yabut ▾

DASHBOARD

Clusters

Hybrid Cloud

Backups

ACCOUNT

Access Management

Billing

Settings

Get Support

ai_first_rag_demo HEALTHY FREE TIER

Overview API Keys Metrics Logs Backups Configuration

Cluster ID: 01ab3875-7bf8-470a-bd47-bbebed07c3ed Version: v1.14.1 Cloud Provider: europe-west3

Nodes: 1 Disk: 4GiB RAM: 1GiB vCPU: 0.5

Get All Qdrant Cloud Features

You are currently using a free tier cluster. Scale up your cluster to a paid plan to get the following features:

✓ Dedicated resources ✓ Monitoring, and log management
✓ Backup and disaster recovery ✓ Standard 10x5 enterprise support plan
✓ Horizontal and vertical scaling ✓ 99.5% uptime SLA

NODE #0 v1.14.1

DISK 368.00 KiB of 4.00 GiB

RAM 147.85 MiB of 1.00 GiB

VCPU 0.04 vCPUs of 0.5 vCPUs

ACTIONS

Scale Cluster ↑ Cluster UI ➔ ...

Access the Cluster

Access Cluster Use your API key to start working with your data.

Try Sample Datasets Explore ready-made datasets in the Cluster UI.

Explore Tutorials Learn the Qdrant API with interactive examples.

Use the API

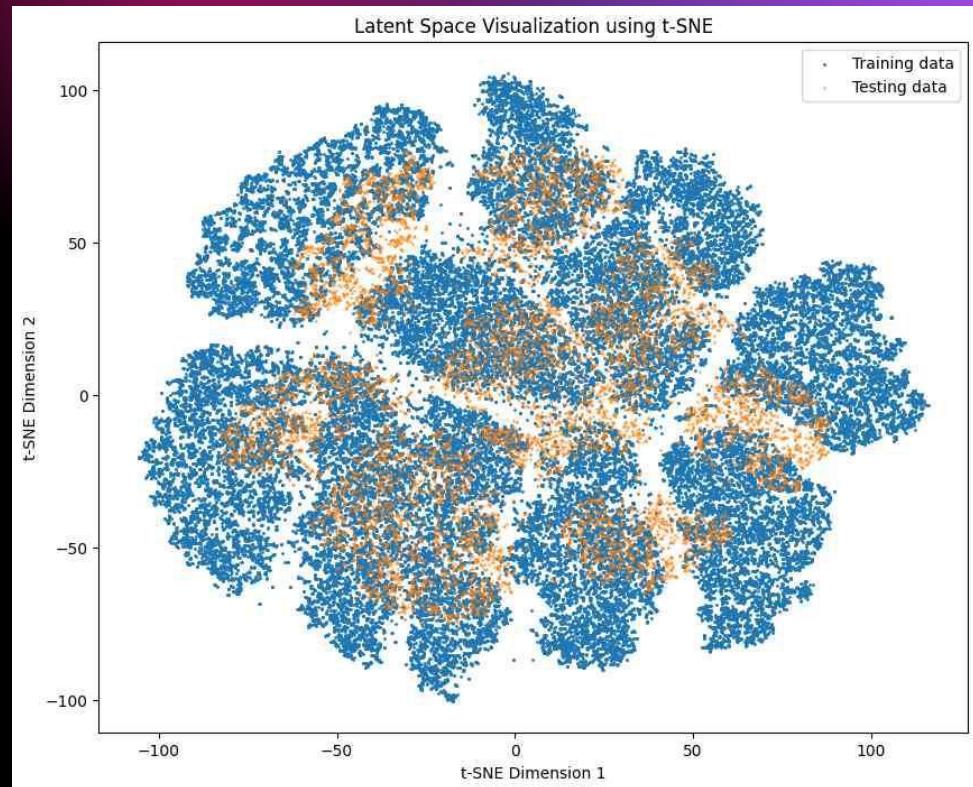
Use the Qdrant API to add, query, and manage vector data effortlessly.

Endpoint <https://01ab3875-7bf8-470a-bd47-bbebed07c3e...>

Examples

Nikko Carlo Yabut

Retrieving Data from the Latent Space



Why can't we use MySQL to search for similar text?



Why can't we use MySQL to search for similar text?

Because MySQL (and traditional databases) are built for **exact matches** or **basic filters** — not understanding.

If you ask:

```
SELECT * FROM documents WHERE text =  
'climate change'
```

You'll get only documents that contain exactly that phrase — nothing semantically related like "global warming" or "environmental shifts."

They operate on surface-level text, **not meaning**.

Traditional DBs use string comparison or keyword search (e.g., LIKE, REGEXP)

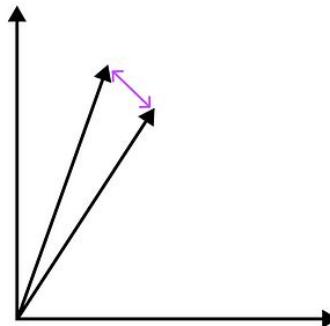
Text A: "The cat sat on the mat."
Text B: "A feline rested on the rug."

A MySQL `LIKE '%cat%'` query would miss Text B entirely

"If we can't use LIKE or REGEX, how can we perform a search?"

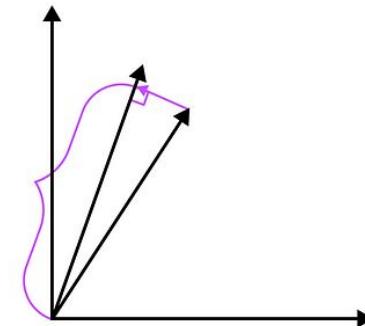
Euclidean Distance

$$d(\mathbf{p}, \mathbf{q}) = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}$$



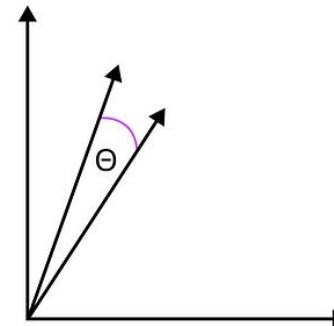
Inner Product

$$\mathbf{a} \cdot \mathbf{b} = \sum_{i=1}^n a_i b_i$$

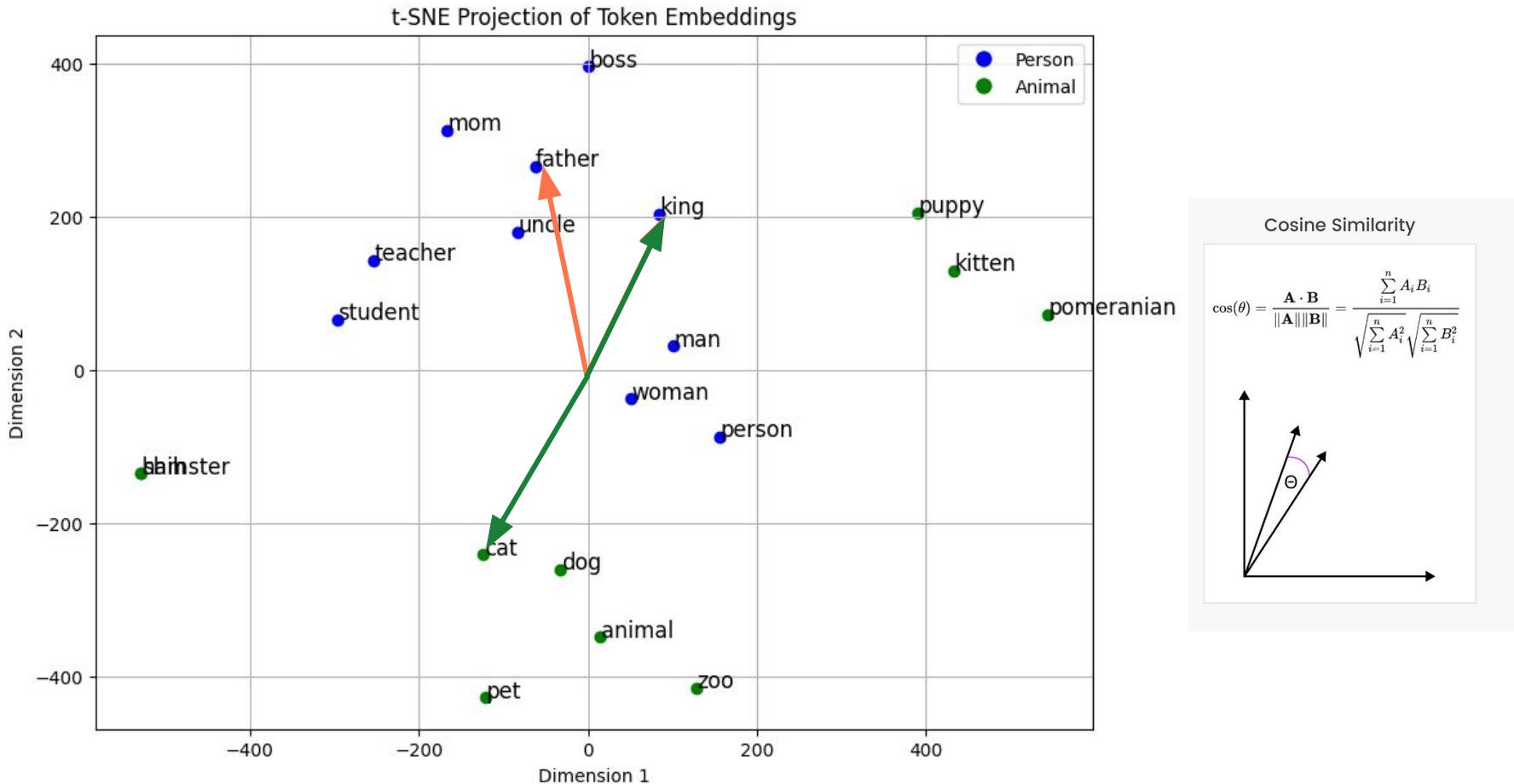


Cosine Similarity

$$\cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$



The Cosine Similarity



The Cosine Similarity

father

king

Cosine similarity = 0.37

king

Cosine similarity = 0.25

cat

Cosine Similarity Matrix:

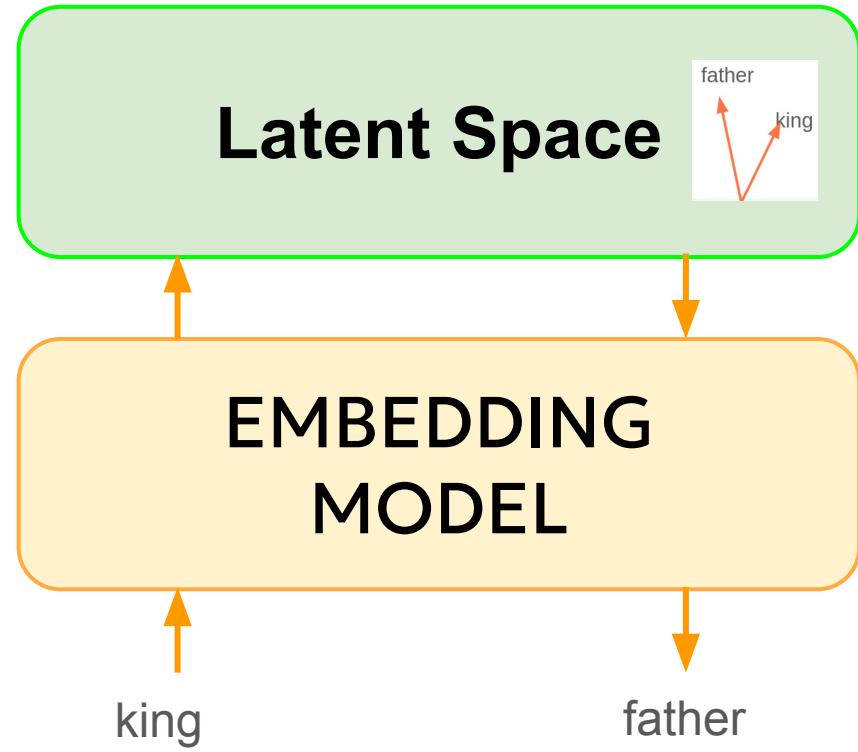
	king	father	teacher	student
king	1.00	0.37	0.26	0.24
father	0.37	1.00	0.40	0.31
teacher	0.26	0.40	1.00	0.51
student	0.24	0.31	0.51	1.00
man	0.35	0.35	0.22	0.24
woman	0.28	0.34	0.33	0.31
person	0.27	0.28	0.30	0.31
mom	0.22	0.49	0.40	0.21
boss	0.29	0.35	0.32	0.21
uncle	0.31	0.56	0.39	0.21
dog	0.23	0.30	0.31	0.31
cat	0.25	0.27	0.29	0.31
puppy	0.23	0.30	0.39	0.21
kitten	0.23	0.27	0.39	0.31
pomeranian	0.24	0.29	0.41	0.31
shih	0.19	0.21	0.28	0.21
zoo	0.21	0.19	0.30	0.21
animal	0.21	0.26	0.33	0.21
pet	0.29	0.29	0.27	0.31
hamster	0.19	0.21	0.28	0.21

Retrieving the Data from the Latent Space



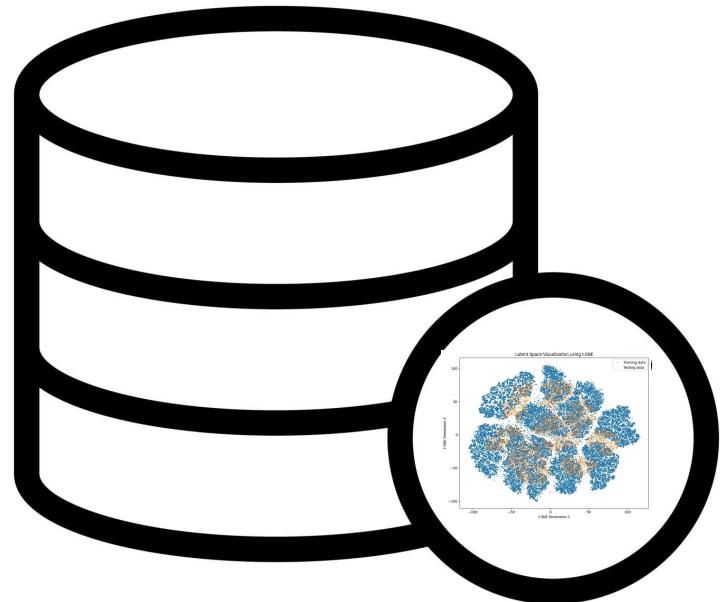
Through embedding models, we can **convert both our documents and queries into vectors** that live in a shared latent space.

By using similarity measures like cosine similarity, we can **identify and then retrieve the most semantically relevant matches**—even if the wording is completely different.



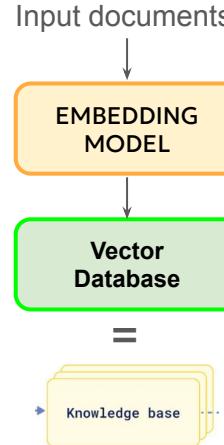
Nikko Carlo Yabut

RAG



RAG (Creating the Knowledge Base)

Note:
If the document is large,
it needs to be
**preprocessed into smaller
chunks** so each piece can
be meaningfully embedded
and retrieved during search
or generation.



- ◆ **Big chunks**

Pros: More context

Cons: Harder to retrieve relevant info, may dilute meaning

- ◆ **Small chunks**

Pros: Precise retrieval

Cons: May lack context, fragmented meaning



Small Chunk (Good for Transactions):

"Order #12345 placed on June 5 for \$49.99."



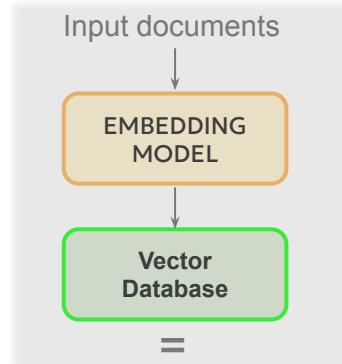
Big Chunk (Good for Stories):

"Maria walked into the store feeling hopeful. After browsing through several aisles, she picked a book that reminded her of home. As she reached the counter, a familiar song played, bringing back childhood memories."

RAG (Inference)

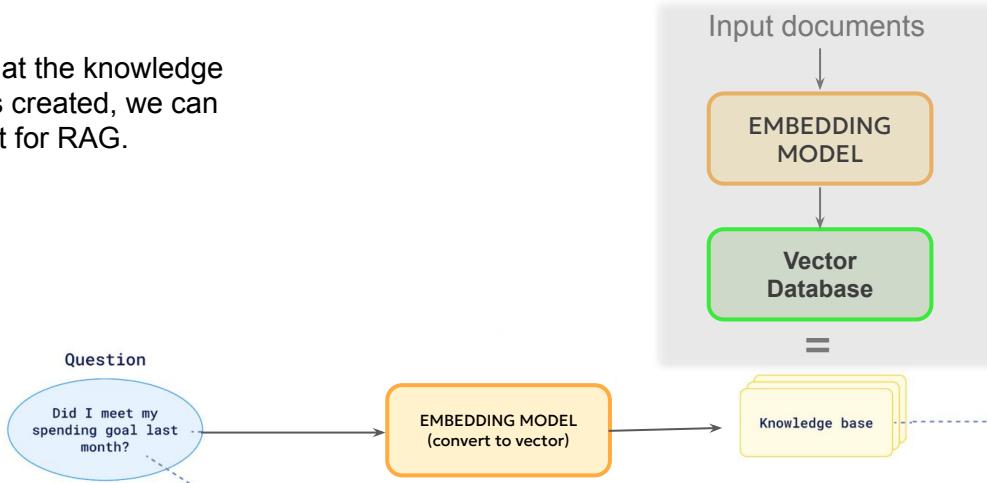
Now that the knowledge base is created, we can query it for RAG.

Question
Did I meet my spending goal last month?



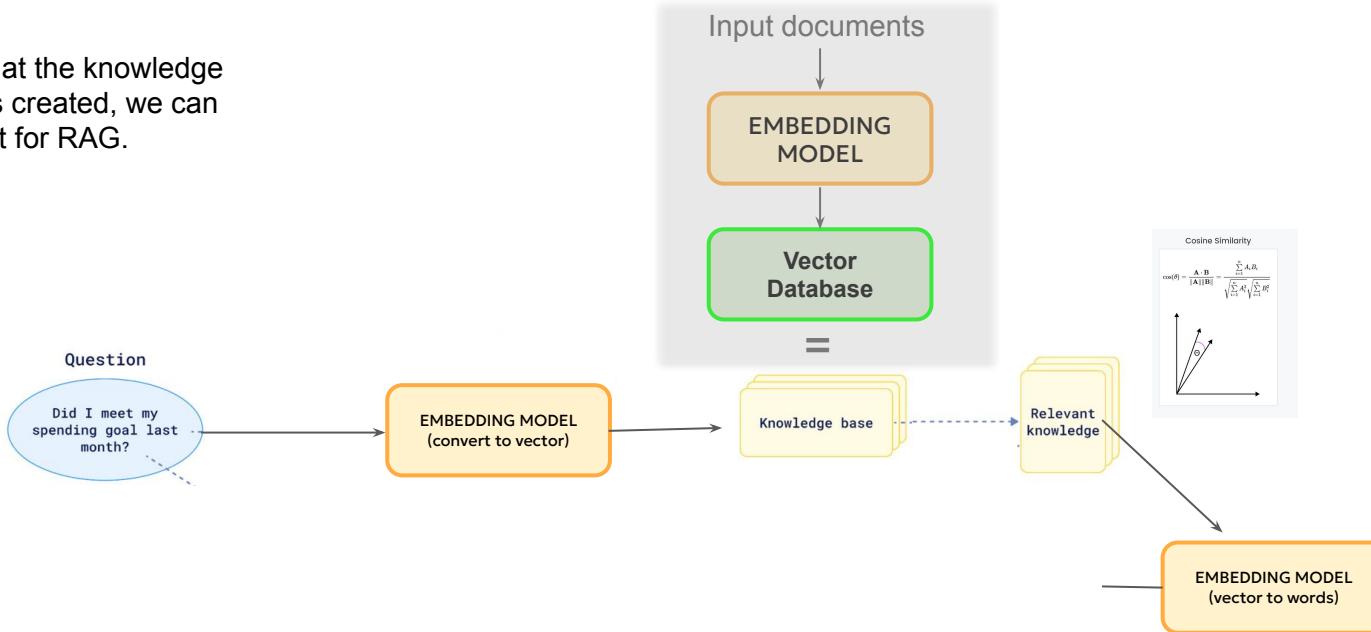
RAG (Inference)

Now that the knowledge base is created, we can query it for RAG.



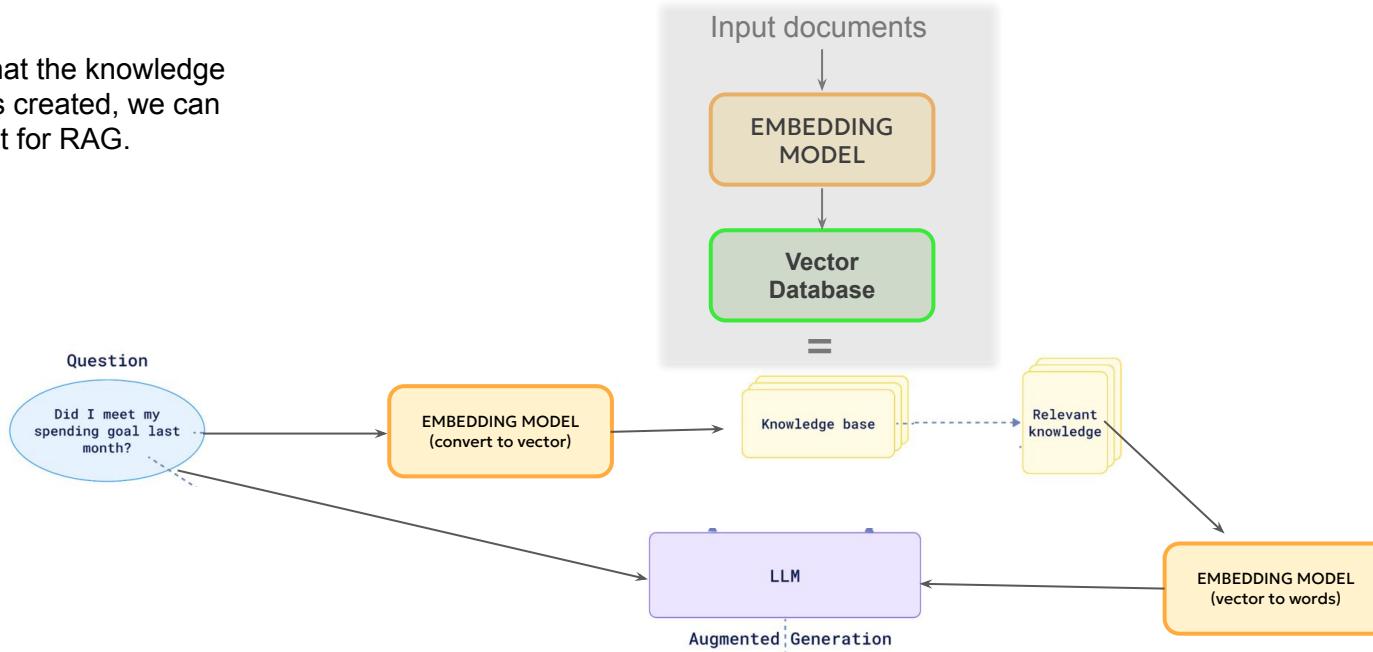
RAG (Inference)

Now that the knowledge base is created, we can query it for RAG.



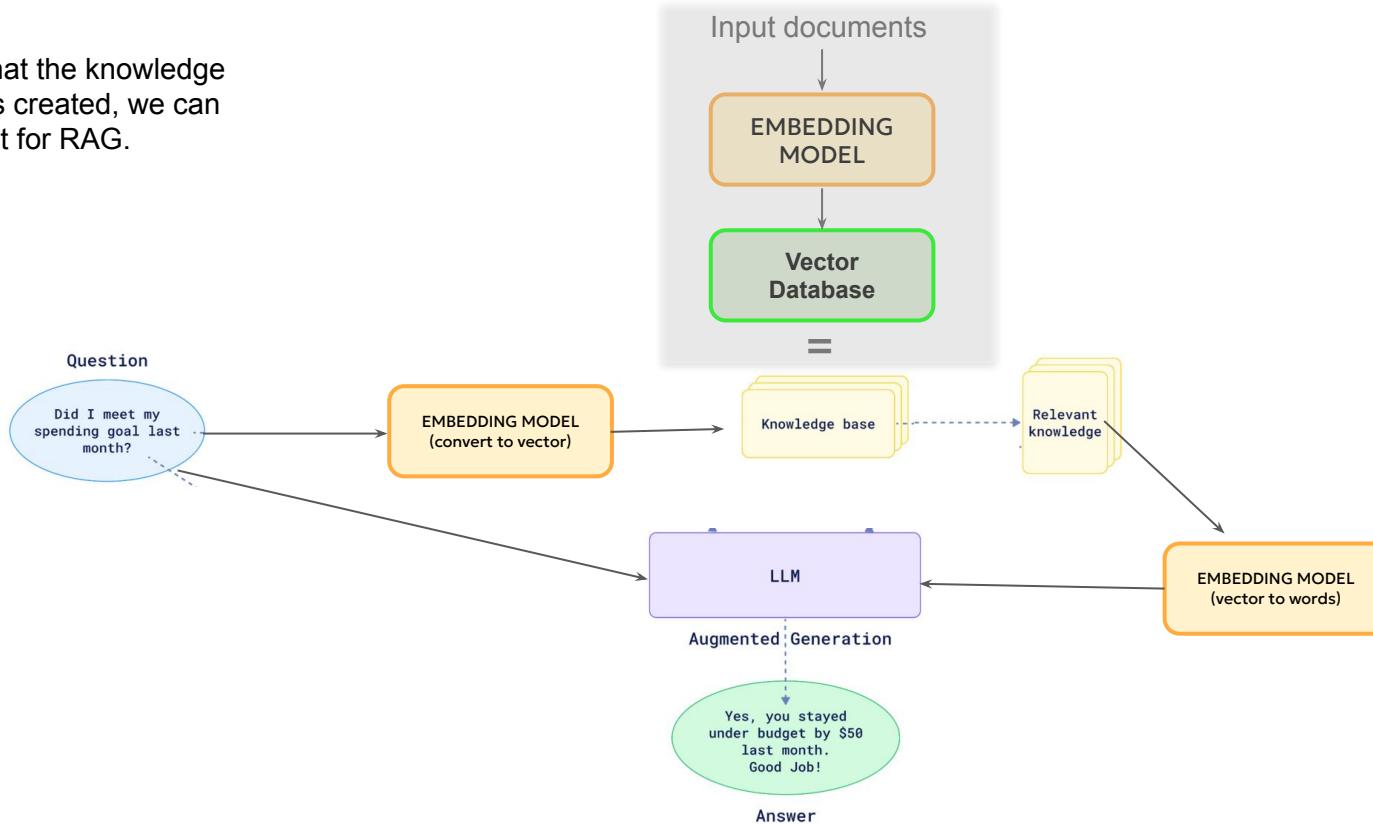
RAG (Inference)

Now that the knowledge base is created, we can query it for RAG.



RAG (Inference)

Now that the knowledge base is created, we can query it for RAG.



What's on your mind today?

If you were to build a chatbot that could answer questions about your company's documents, how would you do it?

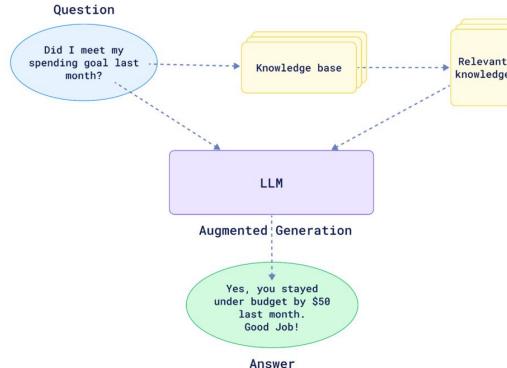
+ Tools



Better Approach:

Instead of fine-tuning the model every time the knowledge base updates, what if we just store the documents in a database — and let the model retrieve the relevant chunks during inference?

RAG



- Use a frozen, pretrained model (e.g., GPT or Mistral)
- Store your documents in a **vector database** (like FAISS, Weaviate, Qdrant, etc.)
- At runtime, convert the query and your documents into **embeddings** (aka latent representations)
- Retrieve relevant documents using **semantic similarity**
- Feed them into the LLM (via context window)

Activity

Set Up Qdrant Cloud Account (Before Class)

1. Go to <https://cloud.qdrant.io>
2. Create an account (free tier is enough)
3. Create a new cluster (select free, small, single region)
4. Inside your cluster, create a **collection**
5. **Copy the following info:**
 - API Key
 - Cluster URL (something like <https://xxxxxxxxxx.qdrant.tech>)

Cont activity in colab:

<https://colab.research.google.com/drive/13LLD5bcfdKadjz0wrJNRjY2mRS97hsPO?usp=sharing>

Business Case Ideation

Create a Business Case utilizing RAG to solve your chosen problem.

1. Define the problem
 - Identify 1–2 clear pain points.
 - ⚠ Focus on content-heavy or info lookup problems
 - Ex: “Agents waste time finding answers,” “Students can’t digest long readings.”
2. Who is your user?
 - Who feels this problem?
Ex: *Support reps, Students, HR staff, etc.*
3. How RAG will help? What does the tool do?
 - What does the tool **do?** Ex: “Answers HR questions from uploaded handbook”
4. What are its main features?
 - a. Upload knowledge base
 - b. User types a question
 - c. App returns grounded answer
5. Build your own streamlit app and do a demo

Nikko Carlo Yabut

LLM's Autoregression

Intro to Agentic AI

Models

CLOSED-SOURCE		OPEN-SOURCE				
		LLM	Vision	VLM	Audio	Word2Vec
 GPT-4 (by OpenAI)	 Amazon Rekognition (by AWS)	 GPT-4V (by OpenAI)		 ElevenLabs		 OpenAI Embeddings
 Claude 3 (by Anthropic)		 Claude 4 Opus (by Anthropic)		 Google Speech-to-Text		 Google Cloud Embeddings
 Gemini (by Google DeepMind)		 Gemini (by Google DeepMind)		 Amazon Transcribe		 Cohere Embeddings
 Command R+ (by Cohere)		 Grok Vision (by xAI)		 Microsoft Azure Speech		 Amazon Comprehend
 Ernie Bot (by Baidu)		 Imagen		 Descript Overdub		
 Qwen (by Alibaba)		 Midjourney				
 SenseChat (by SenseTime)		 Flux				
 Grok (by xAI)						

Using LLM to solve a task (**Chat Completion**)

User: "what hafen vella?"



Prompting Techniques

- Zero-shot Prompting
- One-shot Prompting
- Few-shot Prompting

...

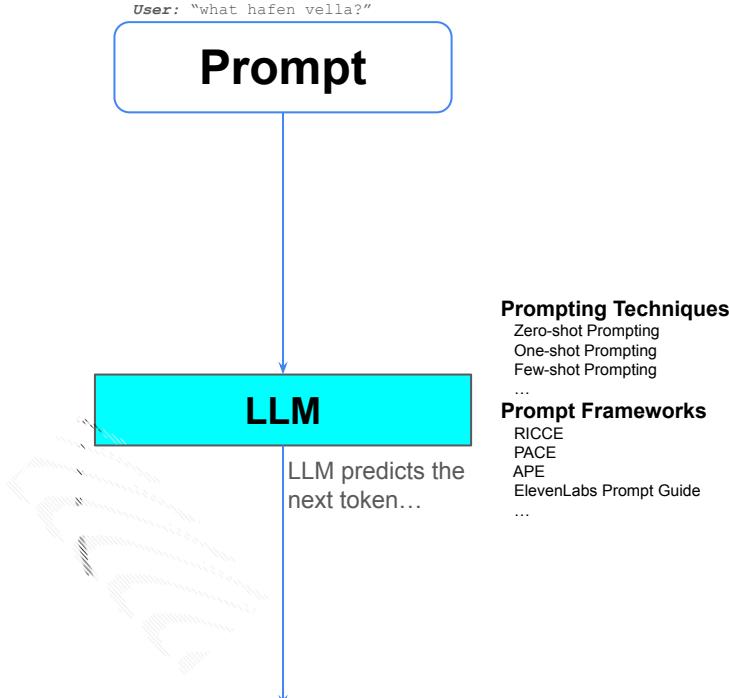
Prompt Frameworks

- RICCE
- PACE
- APE
- ElevenLabs Prompt Guide

...

ChatGPT: The phrase "What hafen Vella?" is a viral meme that originated from a 2013 episode of It's Showtime...

Using LLM to solve a task (**Chat Completion**) -Auto-regression



Prompting Techniques

- Zero-shot Prompting
- One-shot Prompting
- Few-shot Prompting
- ...

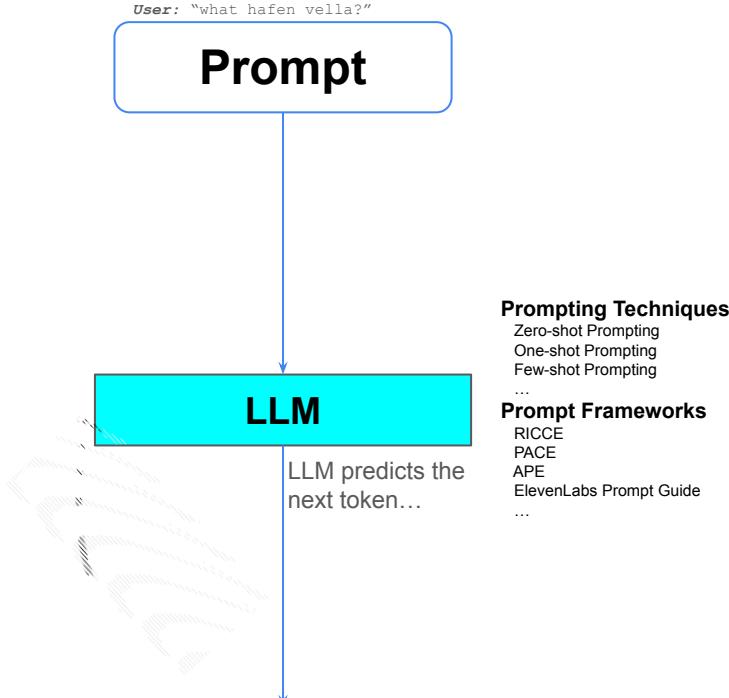
Prompt Frameworks

- RICCE
- PACE
- APE
- ElevenLabs Prompt Guide
- ...

User: "what hafen vella?"

Assistant: "The

Using LLM to solve a task (**Chat Completion**) -Auto-regression



Prompting Techniques

- Zero-shot Prompting
- One-shot Prompting
- Few-shot Prompting
- ...

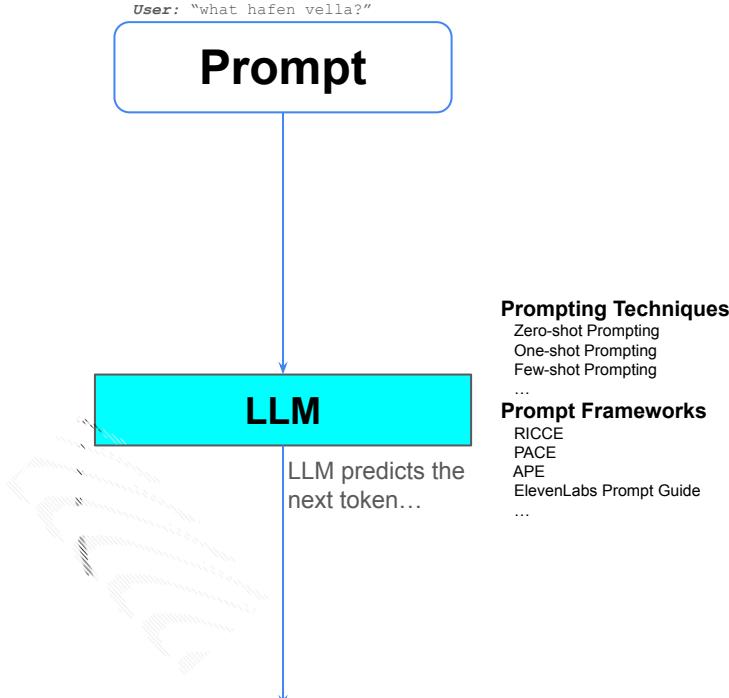
Prompt Frameworks

- RICCE
- PACE
- APE
- ElevenLabs Prompt Guide
- ...

User: "what hafen vella?"

Assistant: "The phrase"

Using LLM to solve a task (**Chat Completion**) -Auto-regression



Prompting Techniques

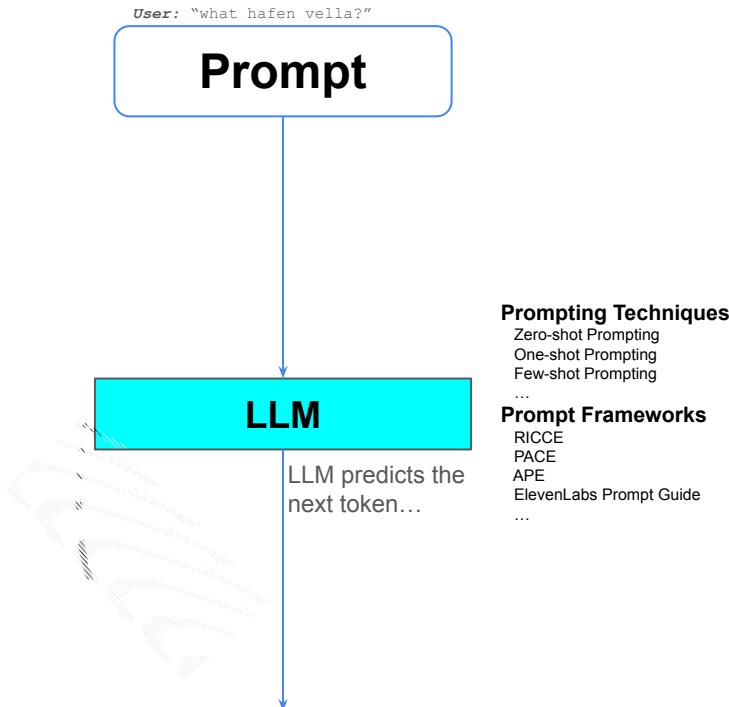
- Zero-shot Prompting
- One-shot Prompting
- Few-shot Prompting
- ...

Prompt Frameworks

- RICCE
- PACE
- APE
- ElevenLabs Prompt Guide
- ...

User: "what hafen vella?"
Assistant: "The phrase "**What**

Using LLM to solve a task (**Chat Completion**) -Auto-regression



Prompting Techniques

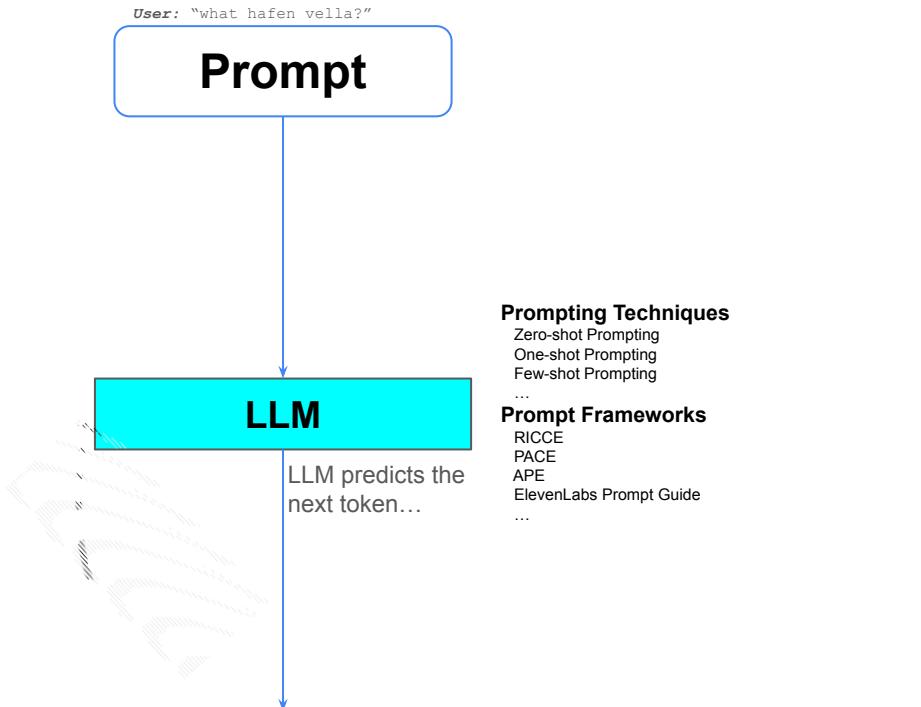
- Zero-shot Prompting
- One-shot Prompting
- Few-shot Prompting
- ...

Prompt Frameworks

- RICCE
- PACE
- APE
- ElevenLabs Prompt Guide
- ...

User: "what hafen vella?"
Assistant: "The phrase "What
hafen

Using LLM to solve a task (**Chat Completion**) -Auto-regression



Prompting Techniques

- Zero-shot Prompting
- One-shot Prompting
- Few-shot Prompting
- ...

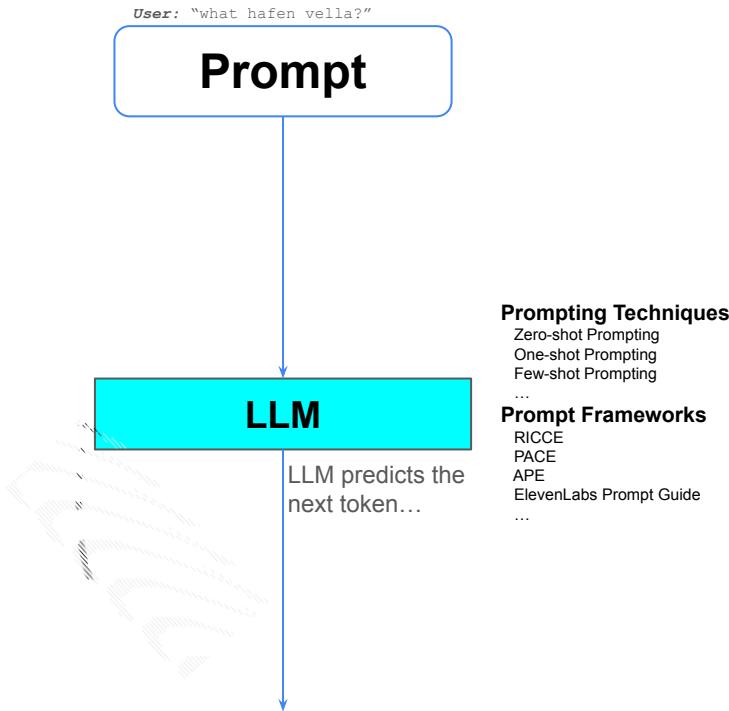
Prompt Frameworks

- RICCE
- PACE
- APE
- ElevenLabs Prompt Guide
- ...

User: "what hafen vella?"

Assistant: "The phrase "What
hafen **vella?**"

Using LLM to solve a task (**Chat Completion**) -Auto-regression



Prompting Techniques

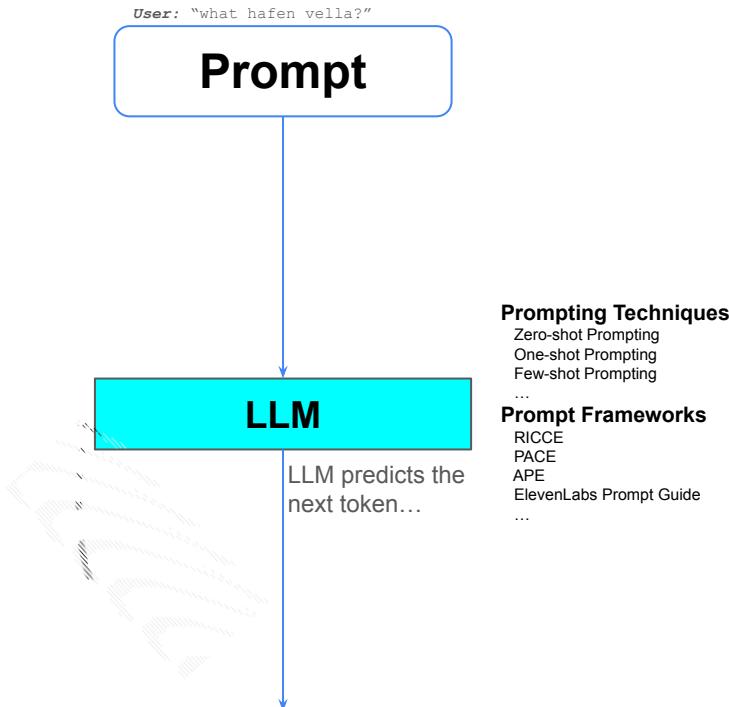
- Zero-shot Prompting
- One-shot Prompting
- Few-shot Prompting
- ...

Prompt Frameworks

- RICCE
- PACE
- APE
- ElevenLabs Prompt Guide
- ...

User: "what hafen vella?"
Assistant: "The phrase "What
hafen vella?" **is**

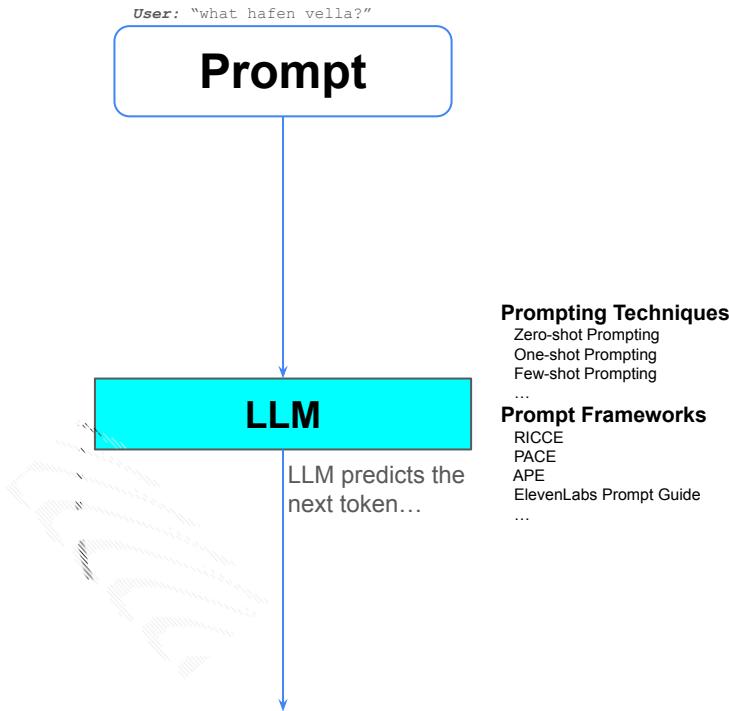
Using LLM to solve a task (**Chat Completion**) -Auto-regression



User: "what hafen vella?"

Assistant: "The phrase "What
hafen vella?" is a

Using LLM to solve a task (**Chat Completion**) -Auto-regression



Prompting Techniques

- Zero-shot Prompting
- One-shot Prompting
- Few-shot Prompting
- ...

Prompt Frameworks

- RICCE
- PACE
- APE
- ElevenLabs Prompt Guide
- ...

User: "what hafen vella?"
Assistant: "The phrase "What hafen vella?" is a **viral**

Using LLM to solve a task (**Chat Completion**) -Auto-regression

User: "what hafen vella?"



This **autoregressive** behavior continues token by token.

Each new token is predicted based on **all previously generated ones**, until the model decides to stop.

In technical terms, the generation ends when the model produces a special token called **<EOS>** (End of Sequence), signaling completion.

User: "what hafen vella?"

Assistant: "The phrase "What hafen vella?" is a viral meme

Using LLM to solve a task (**Chat Completion**) -Auto-regression

User: "what hafen vella?"



LLM

LLM predicts the
next token...

Prompting Techniques

- Zero-shot Prompting
- One-shot Prompting
- Few-shot Prompting
- ...

Prompt Frameworks

- RICCE
- PACE
- APE
- ElevenLabs Prompt Guide
- ...

User: "what hafen vella?"
Assistant: "The phrase "What hafen vella?" is a viral meme that originated from a 2013 episode of It's Showtime. <EOS>"



This **autoregressive** behavior continues token by token.

Each new token is predicted based on **all previously generated ones**, until the model decides to stop.

In technical terms, the generation ends when the model produces a special token called **<EOS>** (End of Sequence), signaling completion.

Using LLM to solve a task (**Chat Completion**) -Auto-regression

User: "what hafen vella?"



LLM predicts the next token...

Prompting Techniques

- Zero-shot Prompting
- One-shot Prompting
- Few-shot Prompting
- ...

Prompt Frameworks

- RICCE
- PACE
- APE
- ElevenLabs Prompt Guide
- ...

User: "what hafen vella?"

Assistant: "The phrase "What hafen vella?" is a viral meme that originated from a 2013 episode of It's

Showtime. <EOS>



While autoregressive LLMs generate text one token at a time based on prior context...

There's one fundamental limitation:

Using LLM to solve a task (**Chat Completion**) -Auto-regression

User: "what hafen vella?"



LLM

LLM predicts the
next token...

Prompting Techniques

- Zero-shot Prompting
- One-shot Prompting
- Few-shot Prompting
- ...

Prompt Frameworks

- RICCE
- PACE
- APE
- ElevenLabs Prompt Guide
- ...

User: "what hafen vella?"

Assistant: "The phrase "What hafen vella?" is a viral meme that originated from a 2013 episode of It's

Showtime. <EOS>



While autoregressive LLMs generate text one token at a time based on prior context...

There's one fundamental limitation:

There's no "back" button.

Once a token is generated, the model can't revise it. It can only move forward, even if it realizes something could have been phrased better.

It's auto-regressive.
It can only move forward,
never backward.

Using LLM to solve a task (**Chat Completion**) -Auto-regression

User: "what hafen vella?"



LLM predicts the next token...

Prompting Techniques

- Zero-shot Prompting
- One-shot Prompting
- Few-shot Prompting
- ...

Prompt Frameworks

- RICCE
- PACE
- APE
- ElevenLabs Prompt Guide
- ...

User: "what hafen vella?"

Assistant: "The phrase "What hafen vella?" is a viral meme that originated from a 2013 episode of It's

Showtime. <EOS>



In ChatGPT, what do you usually do if you want to refine its answer?

Using LLM to solve a task (**Chat Completion**) -Auto-regression

User: "what hafen vella?"



In ChatGPT, what do you usually do if you want to refine its answer?



LLM predicts the next token...

Prompting Techniques

- Zero-shot Prompting
- One-shot Prompting
- Few-shot Prompting
- ...

Prompt Frameworks

- RICCE
- PACE
- APE
- ElevenLabs Prompt Guide
- ...

User: "what hafen vella?"

Assistant: "The phrase "What hafen vella?" is a viral meme that originated from a 2013 episode of It's

Showtime. <EOS>

You probably just...

- Ask a follow-up question
- Rephrase your prompt

Nikko Carlo Yabut

Multi-Turn Conversation

Intro to Agentic AI



In ChatGPT, you usually ask a follow-up question to get a “better” answer.

This back-and-forth is called a **multi-turn conversation**.

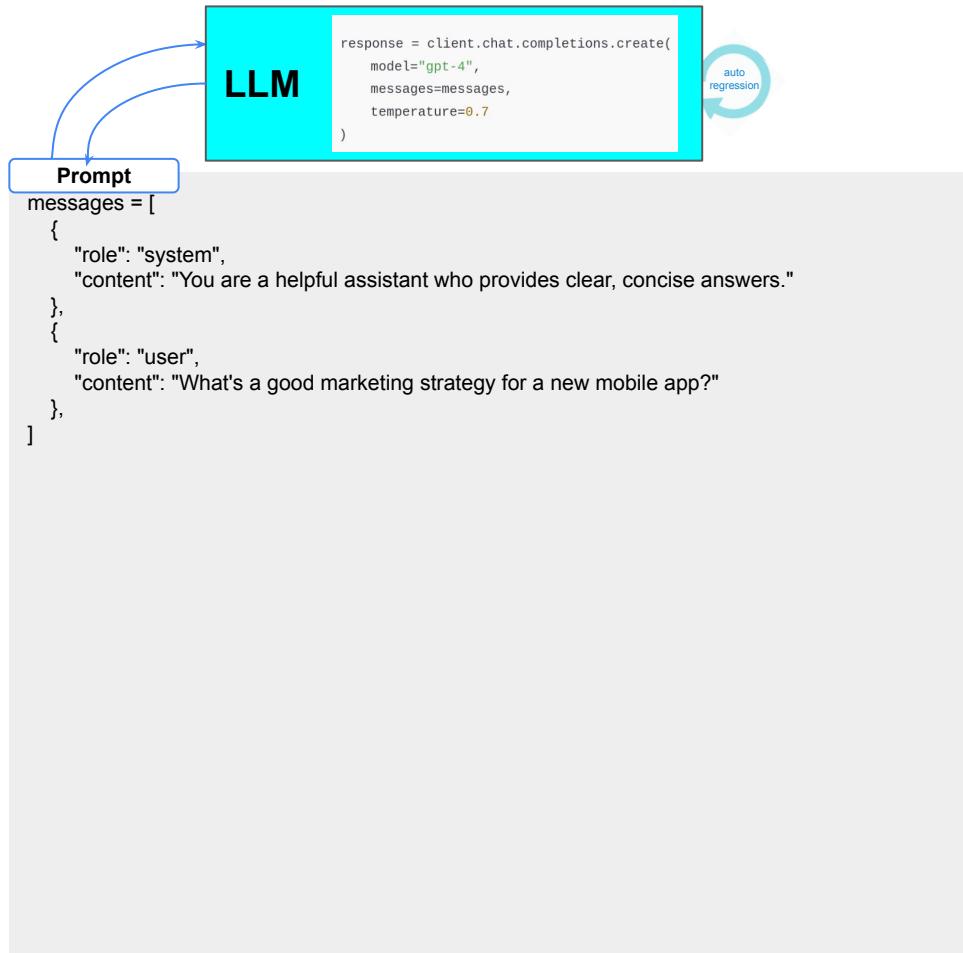
Dialogue in LLMs - Building up the Conversation

```
messages = []
```

Dialogue in LLMs - Building up the Conversation

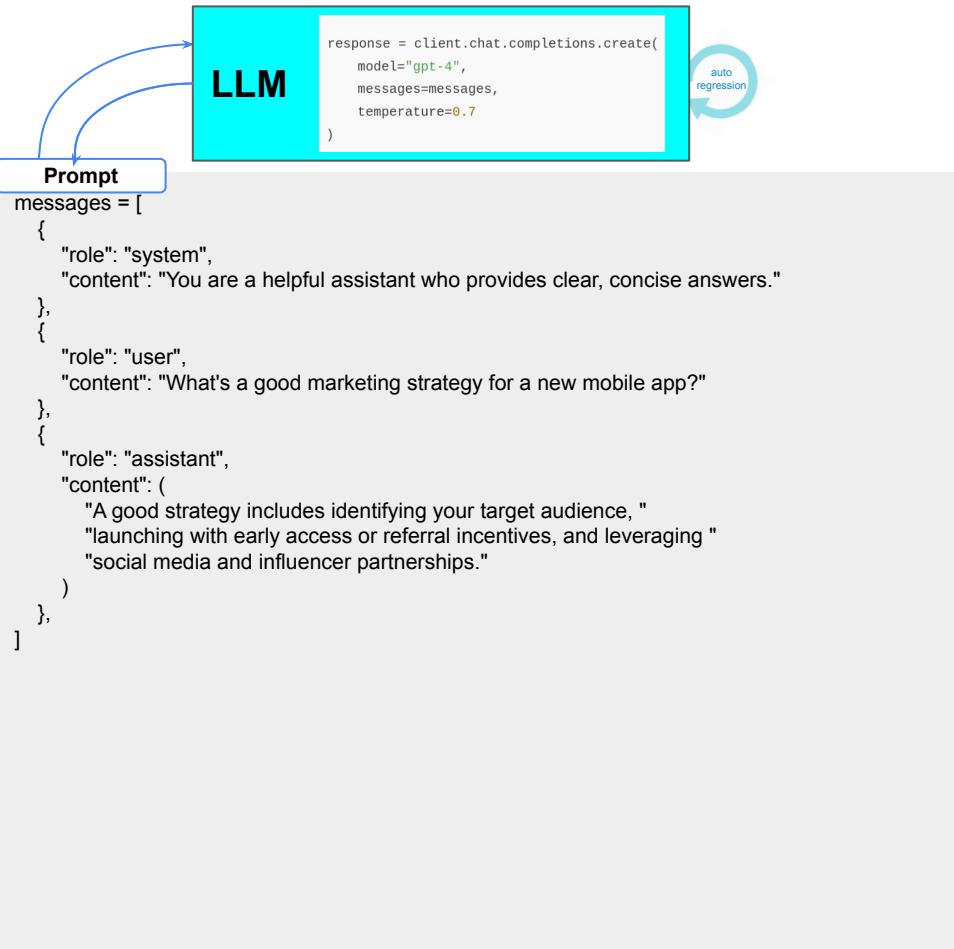
```
messages = [  
    {  
        "role": "system",  
        "content": "You are a helpful assistant who provides clear, concise answers."  
    },  
]
```

Dialogue in LLMs - Building up the Conversation



What's a good marketing strategy for a new mobile app?

Dialogue in LLMs - Building up the Conversation

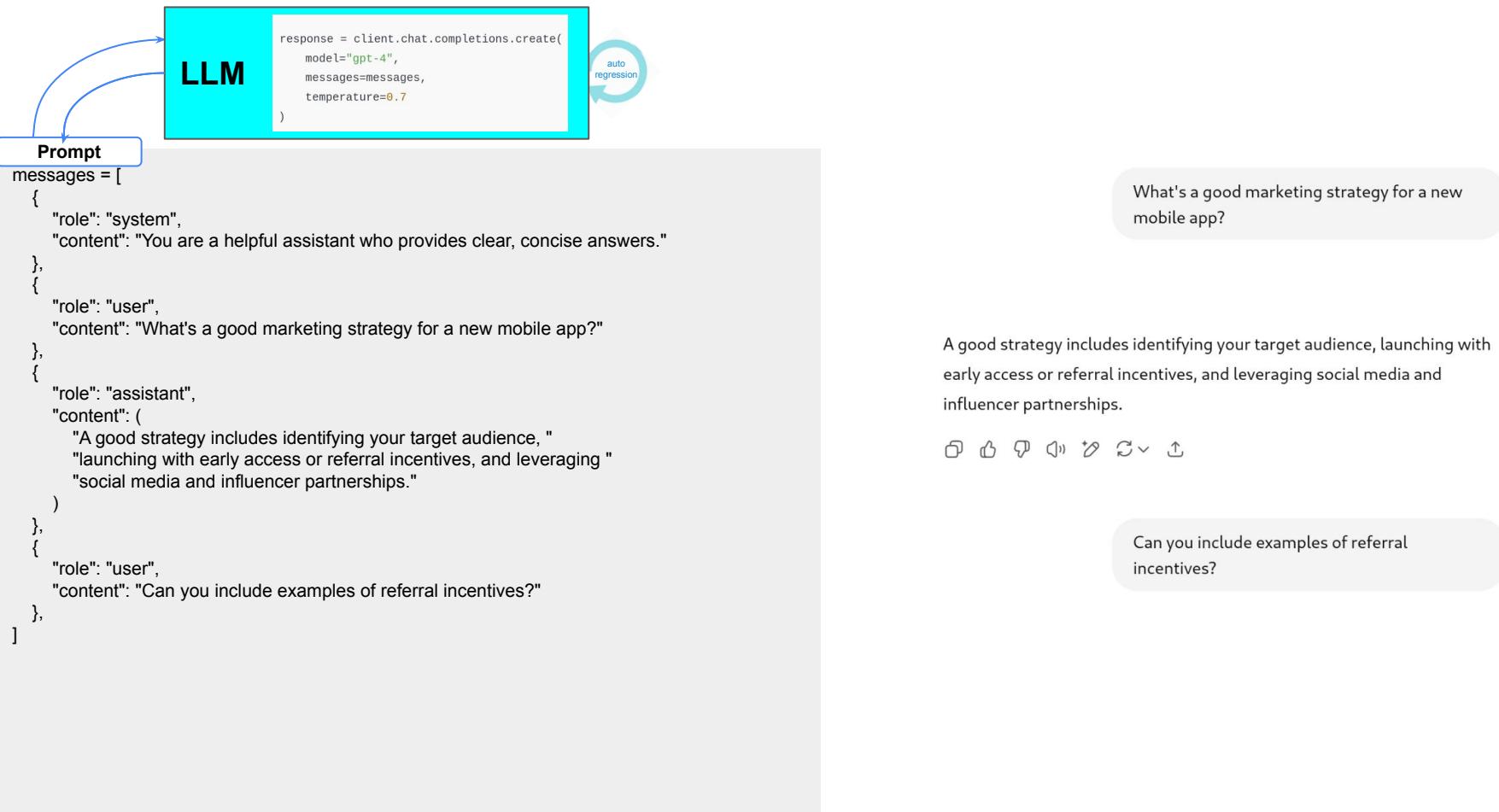


What's a good marketing strategy for a new mobile app?

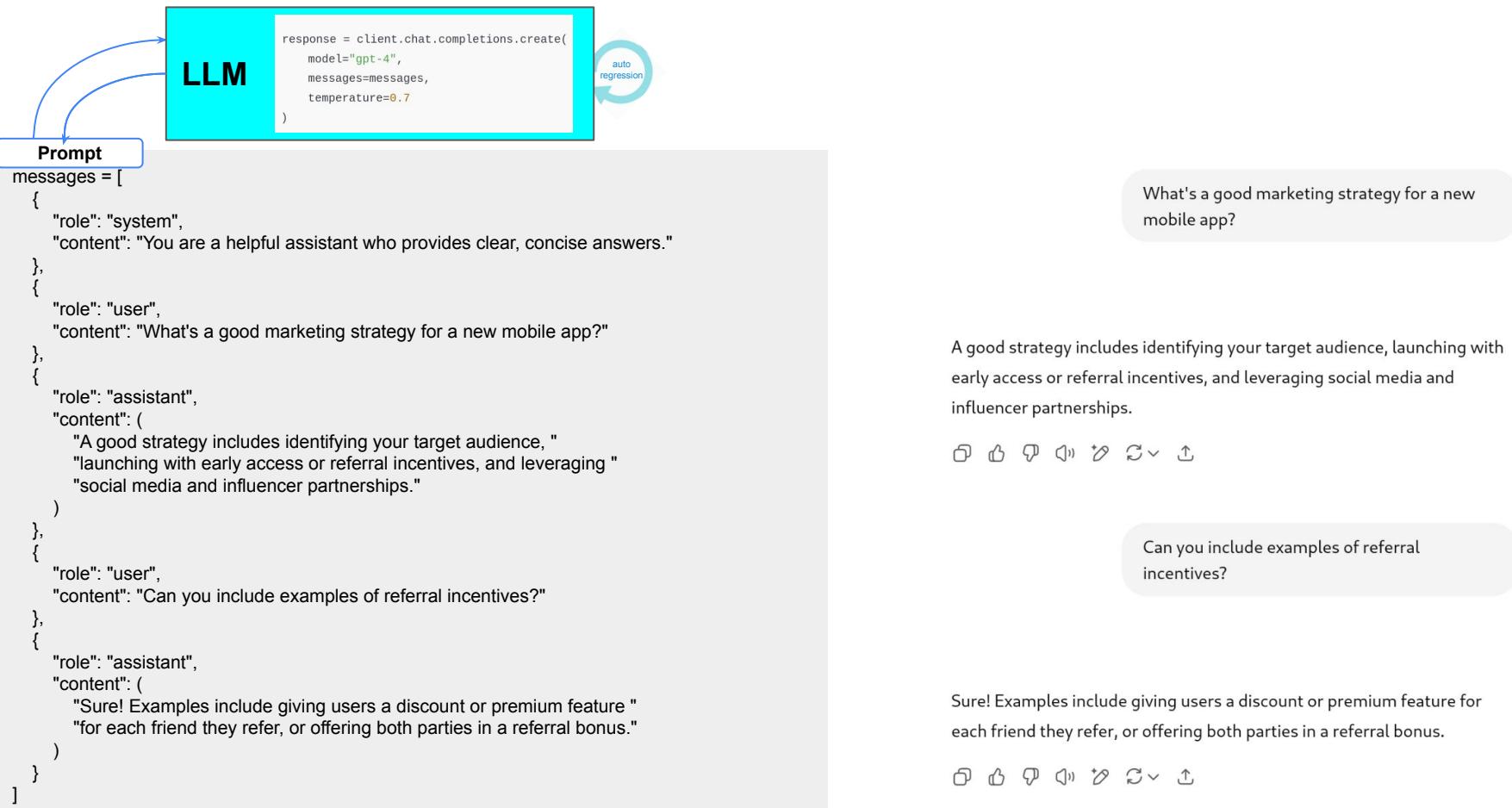
A good strategy includes identifying your target audience, launching with early access or referral incentives, and leveraging social media and influencer partnerships.



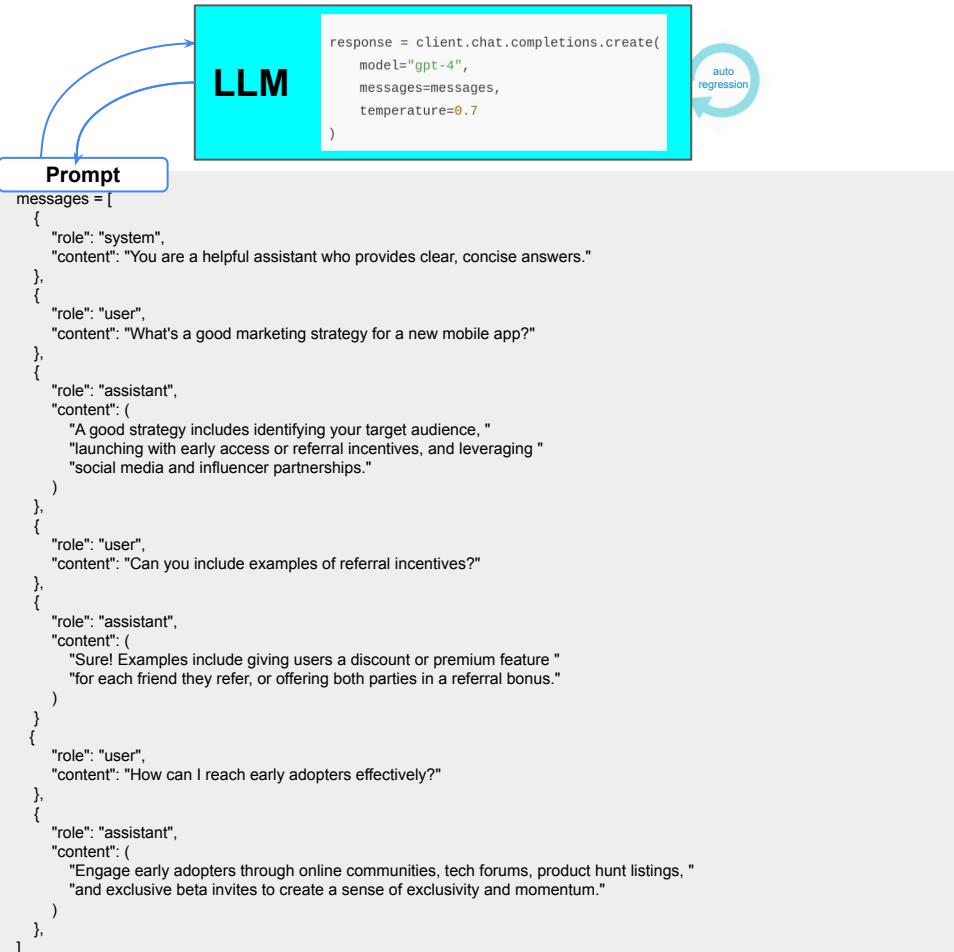
Dialogue in LLMs - Building up the Conversation



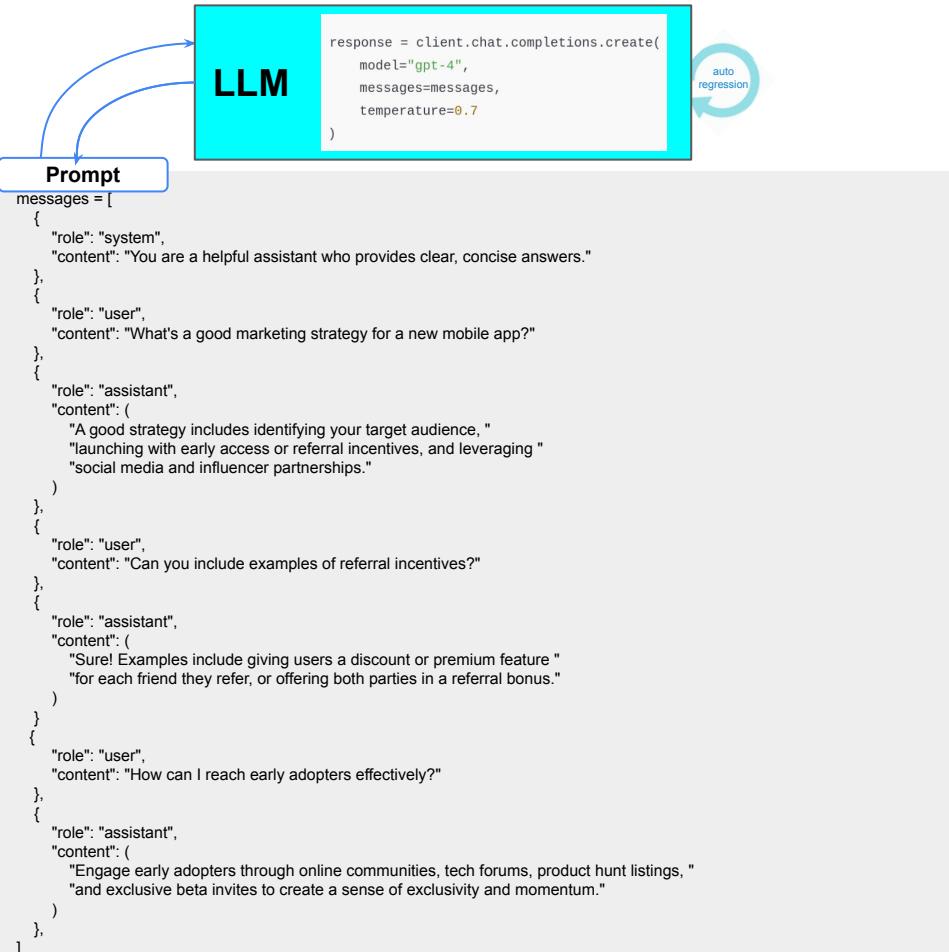
Dialogue in LLMs - Building up the Conversation



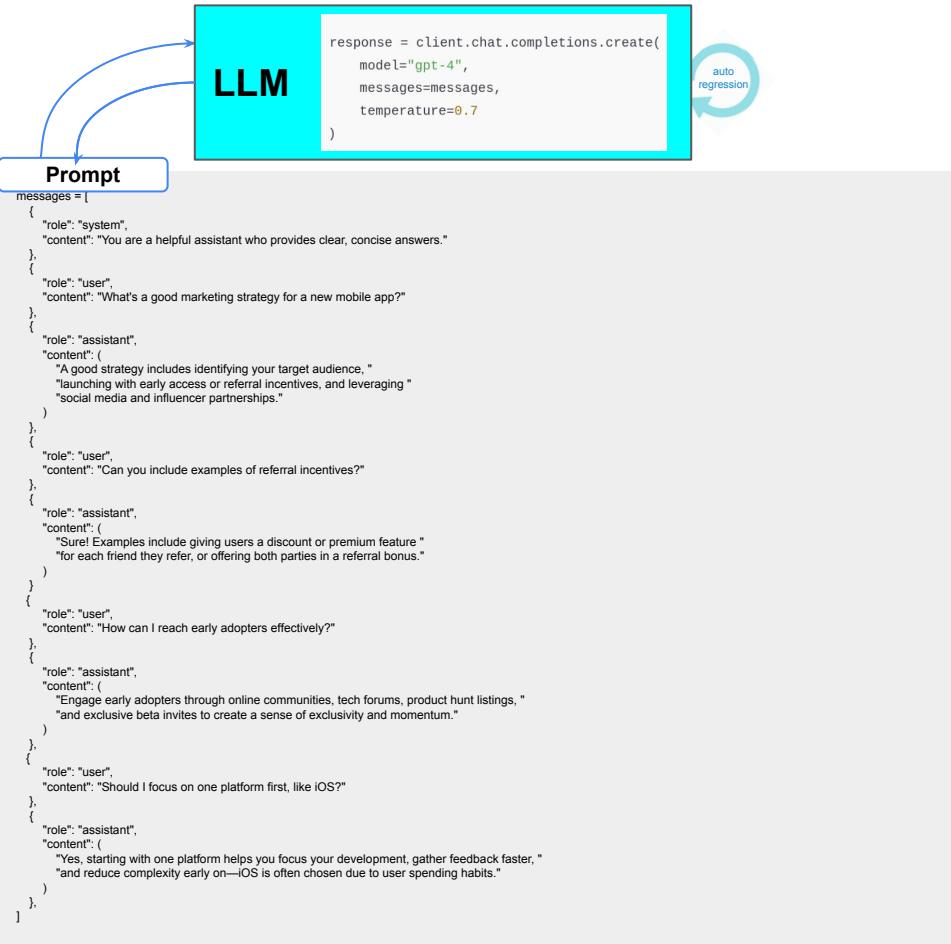
Dialogue in LLMs - Building up the Conversation



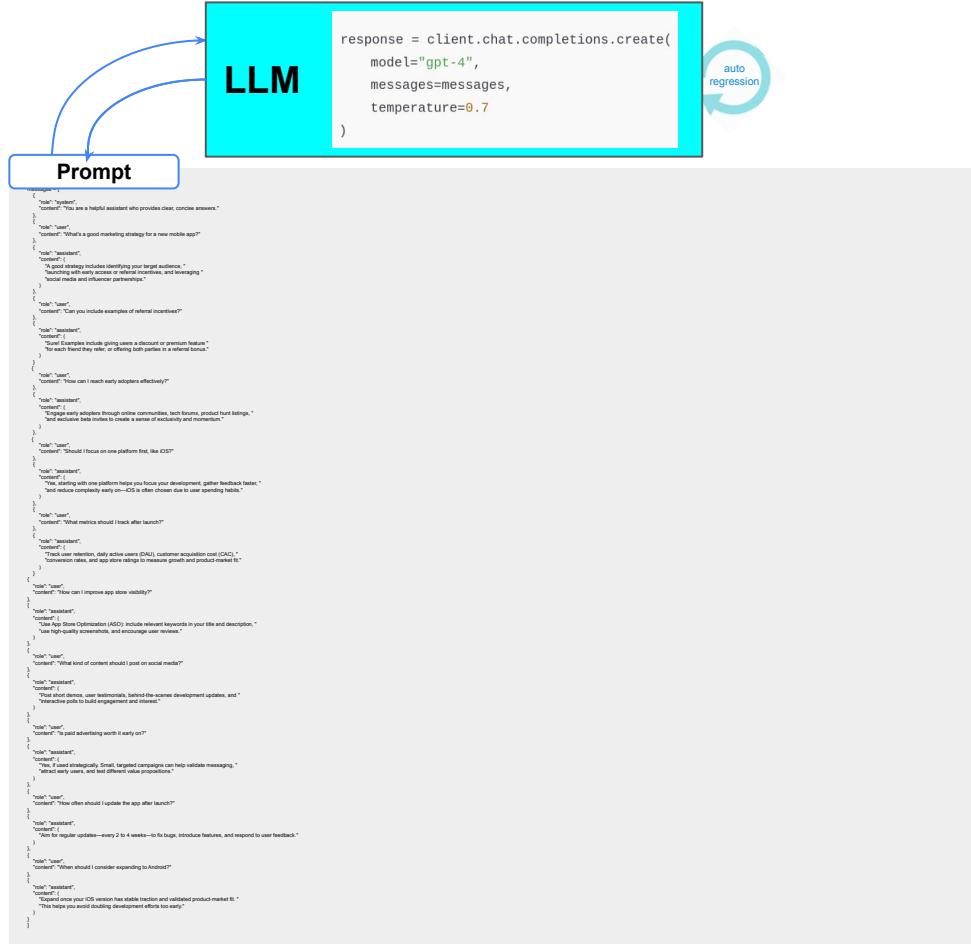
Dialogue in LLMs - Building up the Conversation



Dialogue in LLMs - Building up the Conversation



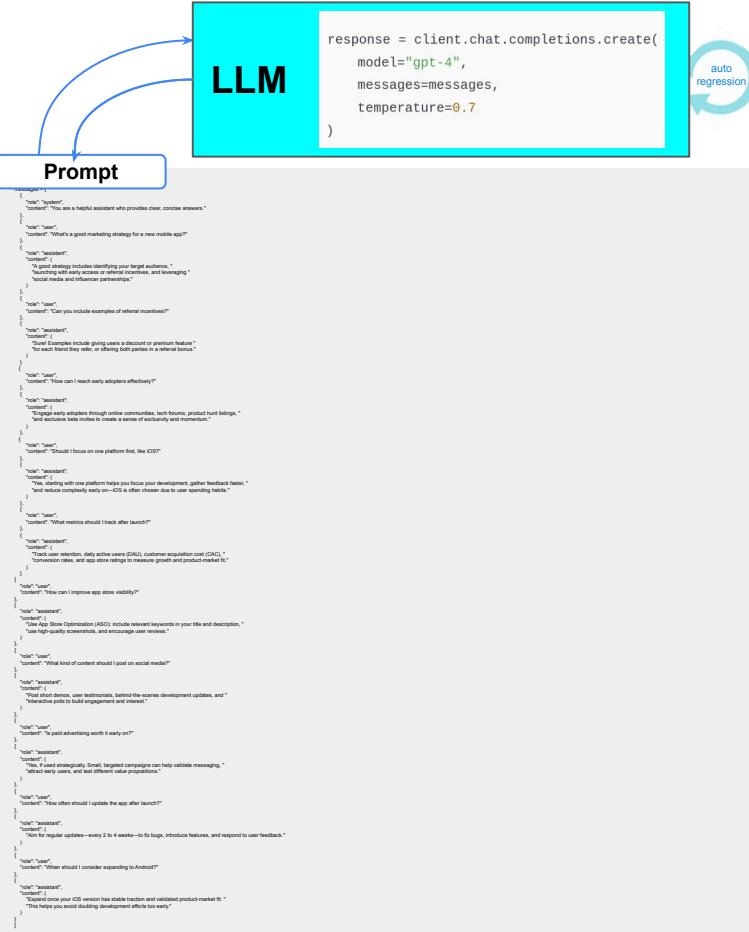
Dialogue in LLMs - Building up the Conversation



Oops!

We'll soon encounter another problem...

Dialogue in LLMs - Building up the Conversation



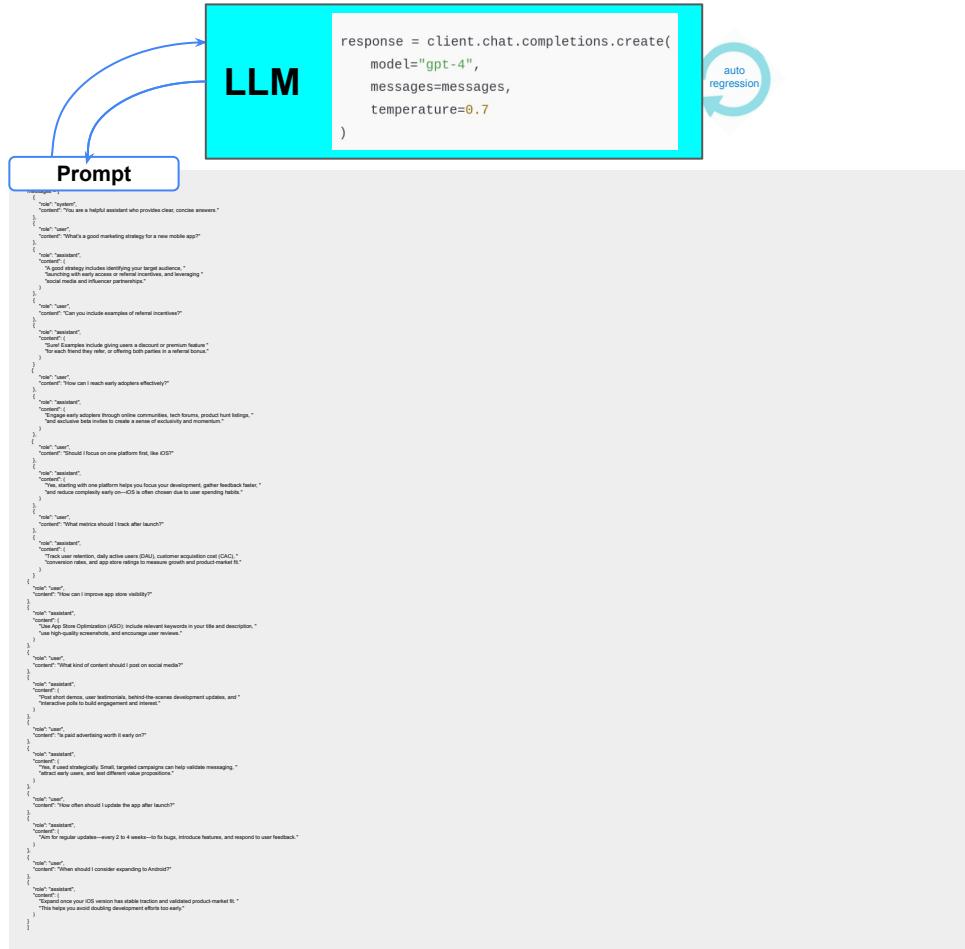
Oops!

We'll soon encounter another problem...

Your context is building up...

As we keep adding turns, we're quietly creating a new challenge, ever-increasing context.

Dialogue in LLMs - Building up the Conversation



Oops!

We'll soon encounter another problem...

Your context building up...

As we keep adding turns, we're quietly creating a new challenge, ever-increasing context.

Your context might hit the max allowable input, aka context window.

Dialogue in LLMs - **Context Window** Problem

Prompt

```
message=1
  "role": "user",
  "content": "You are a helpful assistant who provides clear, concise answers."
}
{
  "role": "user",
  "content": "What's a good marketing strategy for a new mobile app?"
}
{
  "role": "assistant",
  "content": "A good strategy includes identifying your target audience, creating compelling content, utilizing social media and influencer partnerships."
}
{
  "role": "user",
  "content": "Can you include examples of referral incentives?"
}
{
  "role": "assistant",
  "content": "Sure! Examples include giving users a discount or premium feature they can share with their friends, or offering both parties a reward when someone refers them."
}
{
  "role": "user",
  "content": "How can I reach early adopters effectively?"
}
{
  "role": "assistant",
  "content": "Engage early adopters through online communities, tech forums, product hunt listings, and social media groups to create a sense of visibility and momentum."
}
{
  "role": "user",
  "content": "Should I focus on one platform first, like iOS?"
}
{
  "role": "assistant",
  "content": "It depends. Starting with one platform helps you focus your development, gather feedback faster, and reduce complexity early on—iOS is often chosen due to user spending habits."
}
{
  "role": "user",
  "content": "What metrics should I track after launch?"
}
{
  "role": "assistant",
  "content": "Track key metrics like daily active users (DAU), customer acquisition cost (CAC), conversion rates, and app store ratings to measure growth and product market fit."
}
{
  "role": "user",
  "content": "How can I improve app store visibility?"
}
{
  "role": "assistant",
  "content": "Use App Store Optimization (ASO) to include relevant keywords in your title and description, use high-quality screenshots, and encourage user reviews."
}
{
  "role": "user",
  "content": "What kind of content should I post on social media?"
}
{
  "role": "assistant",
  "content": "Post short demos, user testimonials, behind-the-scenes development updates, and interactive polls to build engagement and interest."
}
{
  "role": "user",
  "content": "Is paid advertising worth it early on?"
}
{
  "role": "assistant",
  "content": "It can be if used strategically. Small, targeted campaigns can help validate messaging, attract new users, and test different value propositions."
}
{
  "role": "user",
  "content": "How often should I update the app after launch?"
}
{
  "role": "assistant",
  "content": "Plan for regular updates—every 2 to 4 weeks—to fix bugs, introduce features, and respond to user feedback."
}
{
  "role": "user",
  "content": "When should I consider expanding to Android?"
}
{
  "role": "assistant",
  "content": "Expand once your iOS version has stable traction and validated product-market fit. This helps you avoid doubling development efforts too early."
```

LLM

```
response = client.chat.completions.create(
    model="gpt-4",
    messages=messages,
    temperature=0.7
)
```



Oops!

We'll soon encounter another problem...

Your context building up...

As we keep adding turns, we're quietly creating a new challenge, ever-increasing context.

Your context might hit the max allowable input, aka context window.

Dialogue in LLMs - **Context Window** Problem

```
message = [
    {"role": "user",
     "content": "You are a helpful assistant who provides clear, concise answers."},
    {"role": "assistant",
     "content": "Understood! You are a helpful assistant who provides clear, concise answers."},
    {"role": "user",
     "content": "What's a good marketing strategy for a new mobile app?"},
    {"role": "assistant",
     "content": "A good strategy includes identifying your target audience, creating compelling content, and leveraging social media and influencer partnerships."},
    {"role": "user",
     "content": "Can you include examples of influencer partnerships?"},
    {"role": "assistant",
     "content": "Great example! Include giving users a discount or premium feature they can't find anywhere, or offering early access to a limited edition."},
    {"role": "user",
     "content": "How can I reach early adopters effectively?"},
    {"role": "assistant",
     "content": "Consider using social media ads, influencer partnerships, and targeted email campaigns."}
```

CONTEXT WINDOW

```
message = [
    {"role": "user",
     "content": "Starting with one platform helps you focus your development, gather feedback faster, and reduce complexity early on—OOS is often chosen due to user spending habits."},
    {"role": "user",
     "content": "What metrics should I track after launch?"},
    {"role": "assistant",
     "content": "Track user retention, daily active users (DAU), customer acquisition cost (CAC), conversion rates, and app store ratings to measure growth and product-market fit."},
    {"role": "user",
     "content": "How can I improve app store visibility?"},
    {"role": "assistant",
     "content": "Use App Store Optimization (ASO) to include relevant keywords in your title and description, use high-quality screenshots, and encourage user reviews."},
    {"role": "user",
     "content": "What kind of content should I post on social media?"},
    {"role": "assistant",
     "content": "Post short demos, user testimonials, behind-the-scenes development updates, and interactive polls to build engagement and interest."},
    {"role": "user",
     "content": "Is paid advertising worth it early on?"},
    {"role": "assistant",
     "content": "Yes! If used strategically, small, targeted campaigns can help validate messaging, attract early users, and test different value propositions."},
    {"role": "user",
     "content": "How often should I update the app after launch?"},
    {"role": "assistant",
     "content": "Plan for regular updates—every 2 to 4 weeks—to fix bugs, introduce features, and respond to user feedback."},
    {"role": "user",
     "content": "When should I consider expanding to Android?"},
    {"role": "assistant",
     "content": "Expand once your iOS version has stable traction and validated product-market fit. This helps you avoid doubling development efforts too early."
```

LLM

```
response = client.chat.completions.create(
    model="gpt-4",
    messages=messages,
    temperature=0.7
)
```



Oops!

We'll soon encounter another problem...

Your context building up...

As we keep adding turns, we're quietly creating a new challenge, ever-increasing context.

Your context might hit the max allowable input, aka context window.

Dialogue in LLMs - **Context Window** Problem

```
message = [
    {"role": "user",
     "content": "You are a helpful assistant who provides clear, concise answers."},
    {"role": "assistant",
     "content": "Good strategy includes identifying your target audience."},
    {"role": "user",
     "content": "What's a good marketing strategy for a new mobile app?"},
    {"role": "assistant",
     "content": "A good strategy includes identifying your target audience."},
    {"role": "user",
     "content": "How can I reach early adopters effectively?"}
]
for message in messages:
    response = client.chat.completions.create(
        model="gpt-4",
        messages=[{"role": "user", "content": message["content"]}],
```

CONTEXT WINDOW

```
message = [
    {"role": "user",
     "content": "How can I improve app store visibility?"}
]
for message in messages:
    response = client.chat.completions.create(
        model="gpt-4",
        messages=[{"role": "user", "content": message["content"]}],
```

```
message = [
    {"role": "user",
     "content": "What kind of content should I post on social media?"}
]
for message in messages:
    response = client.chat.completions.create(
        model="gpt-4",
        messages=[{"role": "user", "content": message["content"]}],
```

```
message = [
    {"role": "user",
     "content": "Is paid advertising worth it?"}
]
for message in messages:
    response = client.chat.completions.create(
        model="gpt-4",
        messages=[{"role": "user", "content": message["content"]}],
```

```
message = [
    {"role": "user",
     "content": "How often should I update the app after launch?"}
]
for message in messages:
    response = client.chat.completions.create(
        model="gpt-4",
        messages=[{"role": "user", "content": message["content"]}],
```

```
message = [
    {"role": "user",
     "content": "How long for regular updates—every 2 to 4 weeks—to fix bugs, introduce features, and respond to user feedback?"}
]
for message in messages:
    response = client.chat.completions.create(
        model="gpt-4",
        messages=[{"role": "user", "content": message["content"]}],
```

```
message = [
    {"role": "user",
     "content": "When should I consider expanding to Android?"}
]
for message in messages:
    response = client.chat.completions.create(
        model="gpt-4",
        messages=[{"role": "user", "content": message["content"]}],
```

Year	Model(s)	Context Window	= Word Count	= Page Count
2020	GPT-3	2,049 tokens	~1,537 words	~3 pages
2021	GPT-3, LLaMA 1, BERT variants	2k–4k tokens	~1,500–3,000 words	~3–6 pages
2022	GPT-3.5, Claude 1, PaLM	4k–8k tokens	~3,000–6,000 words	~6–12 pages

LLM

```
response = client.chat.completions.create(
    model="gpt-4",
    messages=messages,
    temperature=0.7
)
```

“

We'll soon encounter another problem...

Your context building up...

As we keep adding turns, we're quietly creating a new challenge, ever-increasing context.

Your context might hit the max allowable input, aka context window.

Dialogue in LLMs - From **Context Window** Problem

CONTEXT WINDOW

```
message="1"
  "role": "user",
  "content": "You are a helpful assistant who provides clear, concise answers."
}
{
  "role": "user",
  "content": "What's a good marketing strategy for a new mobile app?"
}
{
  "role": "assistant",
  "content": "To start, focus on creating a compelling value proposition that highlights unique features and benefits. Consider running pre-launch marketing to generate buzz and anticipation. Utilize social media, influencer partnerships, and targeted advertising to reach your target audience. Additionally, offer early access or discounts to early adopters to encourage user acquisition and engagement."}
```

LLM Context Window Comparison

Model	Tokens	≈ Words	≈ Pages
GPT-4-turbo (OpenAI)	128,000	~96,000 words	~192 pages
GPT-4 (original)	8,192–32,768	~6,100–24,600	~12–49 pages
GPT-3.5-turbo	16,385	~12,289 words	~25 pages
Claude 3 Opus	200,000	~150,000 words	~300 pages
Gemini 1.5 Pro	1,000,000	~750,000 words	~1,500 pages
Mistral Large	32,000	~24,000 words	~48 pages
Mixtral (MoE)	32,000	~24,000 words	~48 pages
LLaMA 2	4k–32k	~3k–24k words	~6–48 pages
Command R+ (Cohere)	128,000	~96,000 words	~192 pages
Yi-1.5 (01.AI)	32,000	~24,000 words	~48 pages
Grok (xAI)	~128,000	~96,000 words	~192 pages

LLM

```
response = client.chat.completions.create(
    model="gpt-4",
    messages=messages,
    temperature=0.7
)
```



Modern LLMs have huge context windows.

Models like GPT-4-turbo can handle up to **128,000 tokens**, and Claude 3 can go up to **200,000 tokens** — that's the length of a novel or an entire codebase.

So... problem solved?



Dialogue in LLMs - From **Context Window** Problem

CONTEXT WINDOW

```
message="1"
  "role": "user",
  "content": "You are a helpful assistant who provides clear, concise answers."
}
{
  "role": "user",
  "content": "What's a good marketing strategy for a new mobile app?"
}
{
  "role": "assistant",
  "content": "To start, focus on creating a compelling value proposition that addresses a specific user need or pain point. Consider using A/B testing to refine your messaging. Engage early adopters through online communities, tech forums, product hunt listings, and social media groups to create a sense of credibility and momentum."
```

```

  "role": "user",
  "content": "Should I focus on one platform first, like iOS?"
```

```

  "role": "assistant",
  "content": "Yes, starting with one platform helps you focus your development, gather feedback faster, and reduce complexity early on—iOS is often chosen due to user spending habits."
```

```

  "role": "user",
  "content": "What metrics should I track after launch?"
```

```

  "role": "assistant",
  "content": "Track user retention, daily active users (DAU), customer acquisition cost (CAC), conversion rates, and app store ratings to measure growth and product-market fit."
```

```

  "role": "user",
  "content": "How can I improve app store visibility?"
```

```

  "role": "assistant",
  "content": "Use App Store Optimization (ASO) to include relevant keywords in your title and description, use high-quality screenshots, and encourage user reviews."
```

```

  "role": "user",
  "content": "What kind of content should I post on social media?"
```

```

  "role": "assistant",
  "content": "Post short demos, user testimonials, behind-the-scenes development updates, and interactive polls to build engagement and interest."
```

```

  "role": "user",
  "content": "Is paid advertising worth it early on?"
```

```

  "role": "assistant",
  "content": "It is, if used strategically. Small, targeted campaigns can help validate messaging, attract early users, and test different value propositions."
```

```

  "role": "user",
  "content": "How often should I update the app after launch?"
```

```

  "role": "assistant",
  "content": "For regular updates—every 2 to 4 weeks—to fix bugs, introduce features, and respond to user feedback."
```

```

  "role": "user",
  "content": "When should I consider expanding to Android?"
```

```

  "role": "assistant",
  "content": "Expand once your iOS version has stable traction and validated product-market fit."
```

```

  "role": "user",
  "content": "This helps you avoid doubling development efforts too early."
```

LLM Context Window Comparison

Model	Tokens	≈ Words	≈ Pages
GPT-4-turbo (OpenAI)	128,000	~96,000 words	~192 pages
GPT-4 (original)	8,192–32,768	~6,100–24,600	~12–49 pages
GPT-3.5-turbo	16,385	~12,289 words	~25 pages
Claude 3 Opus	200,000	~150,000 words	~300 pages
Gemini 1.5 Pro	1,000,000	~750,000 words	~1,500 pages
Mistral Large	32,000	~24,000 words	~48 pages
Mixtral (MoE)	32,000	~24,000 words	~48 pages
LLaMA 2	4k–32k	~3k–24k words	~6–48 pages
Command R+ (Cohere)	128,000	~96,000 words	~192 pages
Yi-1.5 (01.AI)	32,000	~24,000 words	~48 pages
Grok (xAI)	~128,000	~96,000 words	~192 pages



Modern LLMs have **huge context windows**.

Models like GPT-4-turbo can handle up to **128,000 tokens**, and Claude 3 can go up to **200,000 tokens** — that's the length of a novel or an entire codebase.

So... problem solved? **Not quite.**

LLM

```
response = client.chat.completions.create(
    model="gpt-4",
    messages=messages,
    temperature=0.7
)
```

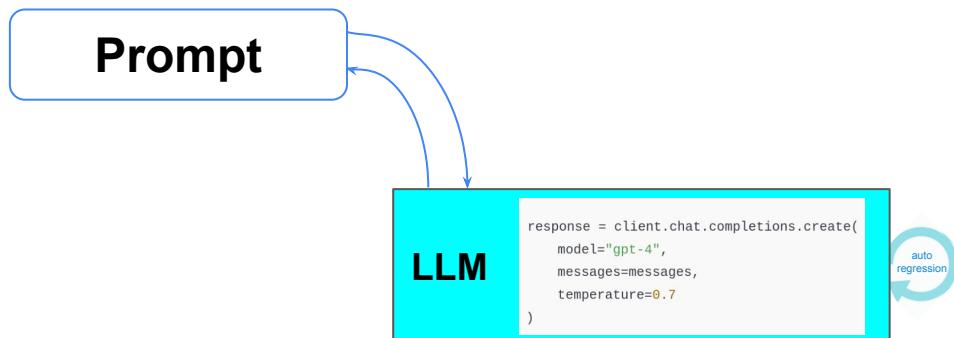
auto
regression

Dialogue in LLMs - From Context Window to **Context Quality** Problem

⚠️ But... Now the Problem Is Context Quality

With more room comes new challenges:

Problem	Description
🎯 Irrelevant or noisy context	Long input ≠ useful input. The model might attend to wrong things.
💡 Low signal-to-noise ratio	Critical information can get buried under fluff.
🧠 Lack of structure or salience	Without cues (headings, bullets, instructions), the model might misunderstand or ignore key parts.
🖼️ Token bloat	Just because you <i>can</i> fit 100k tokens doesn't mean you <i>should</i> . It increases latency, cost, and sometimes worsens performance.



💡 Modern LLMs have **huge** context windows.

Models like GPT-4-turbo can handle up to **128,000 tokens**, and Claude 3 can go up to **200,000 tokens** — that's the length of a novel or an entire codebase.

So... problem solved? **Not quite.**

Now that we *can* stuff all that information in, a new challenge emerges:

🧠 The issue is no longer context **size** — it's context **quality**.

Nikko Carlo Yabut

Context Engineering

Intro to Agentic AI

Context Engineering



context engineering and agentic ai



context engineering



<https://www.llmindex.ai/blog/context-engineering...>

Context Engineering - What it is, and techniques to consider

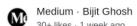
14 hours ago — This is context that allows us to ensure that our agentic ai application is choosing the right resource. Context ordering or compression.



LinkedIn - Nagesh Nama
7 reactions · 5 days ago

Context Engineering for LLMs and Agentic AI

Context engineering is emerging as a critical discipline in building large language model (LLM) applications and AI agents.



Medium - Bijit Ghosh
30+ likes · 1 week ago

Context Engineering is Runtime of AI Agents | by Bijit Ghosh

Context engineering is the practice of designing, managing, and optimizing the information presented to an LLM during runtime. It is: The new ...



<https://borislane.com/blog/context-engineering...>

Context engineering is what makes AI magical - Boris Lane

Jun 22, 2025 — Agentic RAG, builds a contextual snapshot that includes data from multiple sources: the question; related documents; source structure; metadata ...



<https://www.ibm.com/think/topics/agentic-ai...>

What Is Agentic AI? | IBM

Agentic AI is an artificial intelligence system that can accomplish a specific goal with limited supervision. It consists of AI agents—machine learning ...

[What is agentic AI?](#)

[What are the advantages of...](#)



Medium - Shashi Jagtap
6 days ago



<https://www.llmindex.ai/blog/context-engineering...>

Context Engineering - What it is, and techniques to consider

context engineering

<https://github.com/davidkmai/Context-Engineering...>

davidkmai/Context-Engineering

1 day ago — "Context engineering is the delicate art and science of filling the context window with just the right information for the next step."



<https://contextengineering.com...>

Context Engineering | Gilroy, CA

Provider of aluminum project boxes and electronics enclosures. PCB enclosures and heat shrink tubing.



<https://contexteng.com.au...>

Context Engineering - Engineering Services for Project ...

Context Engineering is a small civil engineering consultancy operating out of Brisbane, QLD. We are primarily focused on inner city development.

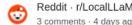
[Who We Are](#) [Contact](#) [What We Do](#)



<https://blog.langchain.com/context-engineering-for-a...>

Context Engineering

2 days ago — Context engineering is the art and science of filling the context window with just the right information at each step of an agent's trajectory.



Reddit - r/LocalLLaMA
3 comments · 4 days ago

What Is Context Engineering? My Thoughts.. : r/LocalLLaMA

Context engineering is about providing the LLM with all the info it needs to answer the query accurately. The difference with plain prompt ...

What's Context Engineering and How Does it Apply Here? 7 posts Jun 30, 2025

What Is Context Engineering? My Views... : r/OpenAI ... 11 posts Jun 29, 2025

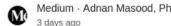
More results from [www.reddit.com](#)



<https://www.llmindex.ai/blog/context-engineering...>

Context Engineering - What it is, and techniques to consider

14 hours ago — What is Context Engineering. AI agents require the relevant context for a task, to perform that task in a reasonable way. We've known this for a ...



Medium - Adnan Masood, PhD.
3 days ago

Context Engineering: Elevating AI Strategy from Prompt ...

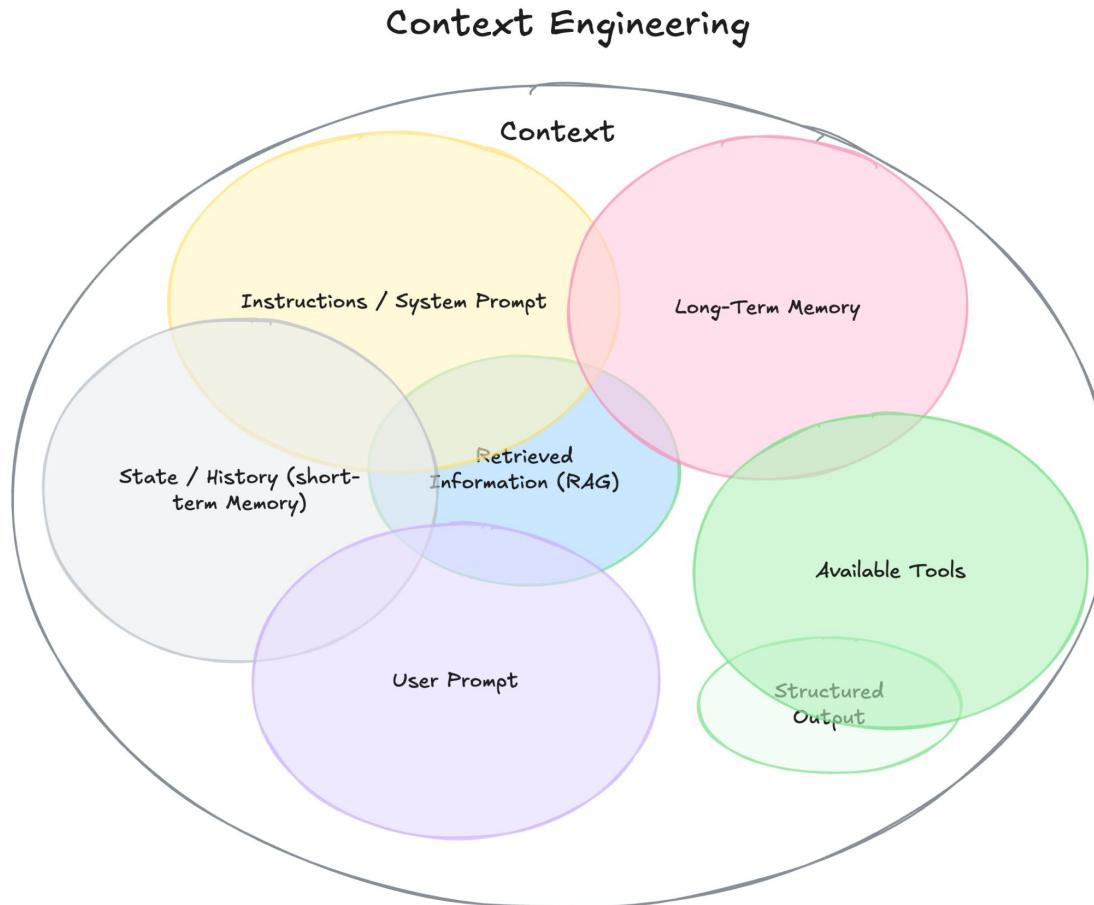
In contrast, context engineering is about providing all the necessary background and information dynamically so the AI can respond effectively [...

People also search for

Context engineering AI

Context engineering LLM

Context Engineering



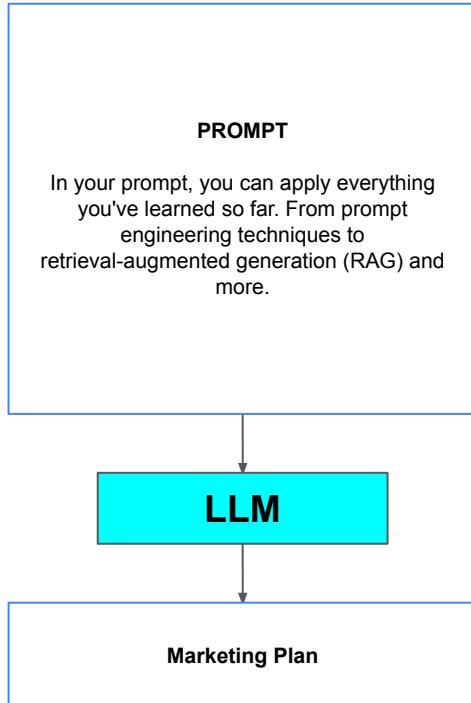
Nikko Carlo Yabut

Compounding LLMs

Intro to Agentic AI

Single LLM

Task: Create a Marketing Plan

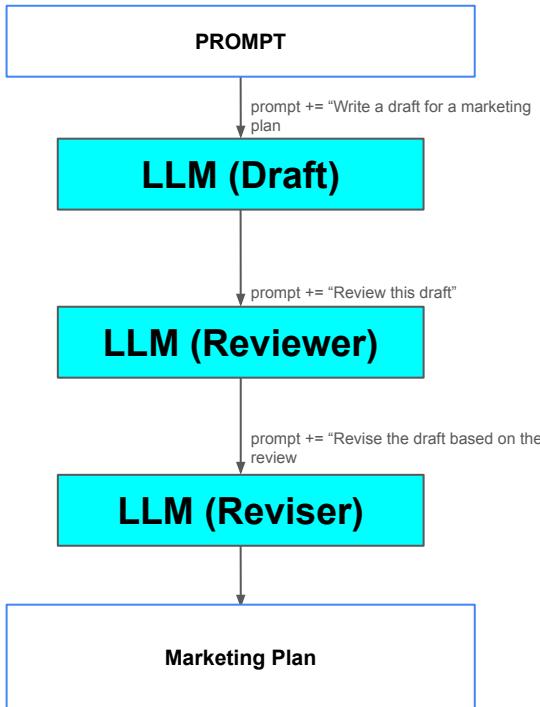


IRL

In real life, we don't just write something once and call it done. There's usually a **reviewer**, followed by a **revision**, until we arrive at a final version.

Compounding LLM Approach

Task: Create a Marketing Plan



IRL

In real life, we don't just write something once and call it done. There's usually a **reviewer**, followed by a **revision**, until we arrive at a final version.

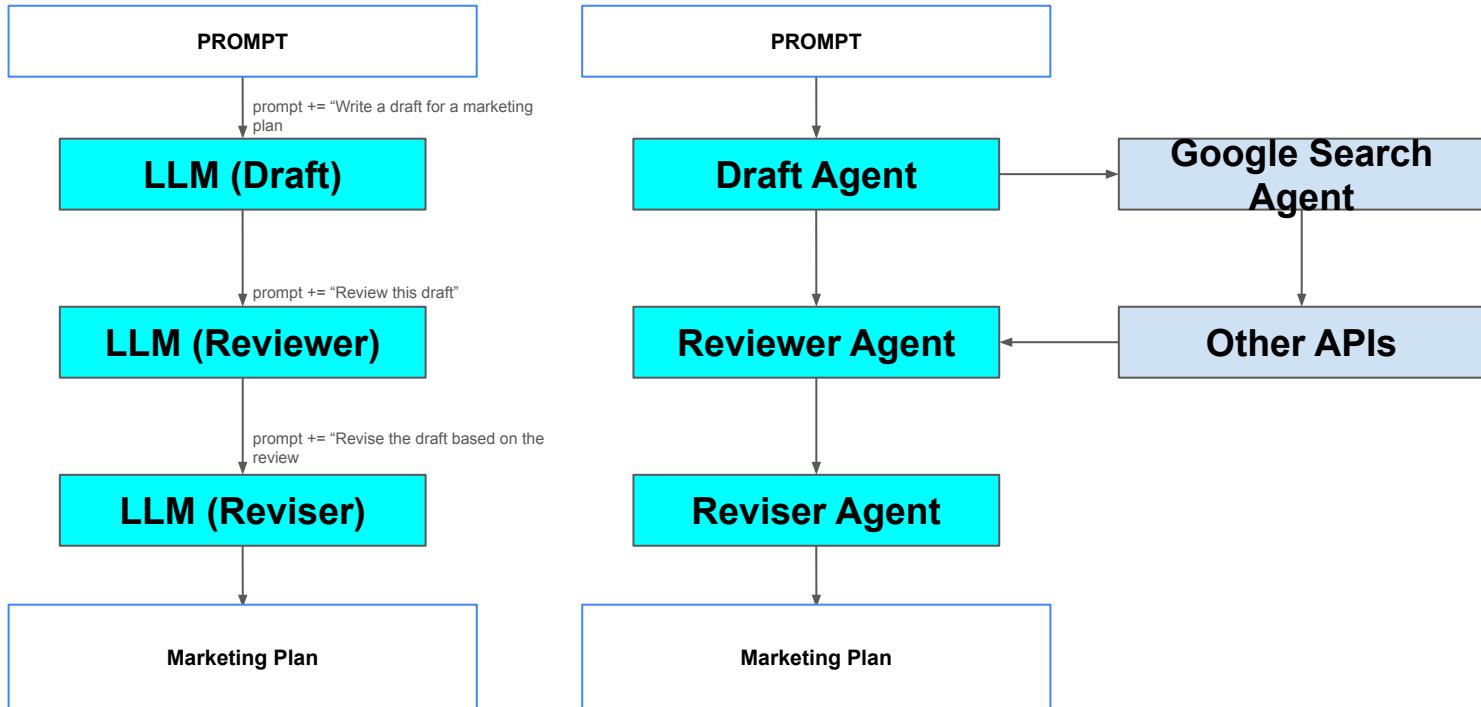


As you can see, this **compounding LLM approach** often yields better results. Your initial draft is **reviewed**, **critiqued**, and then **refined** — just like in real collaborative work.

Compounding LLM Approach as an **Agentic AI**

1 Agents

2 Agent can be a tool

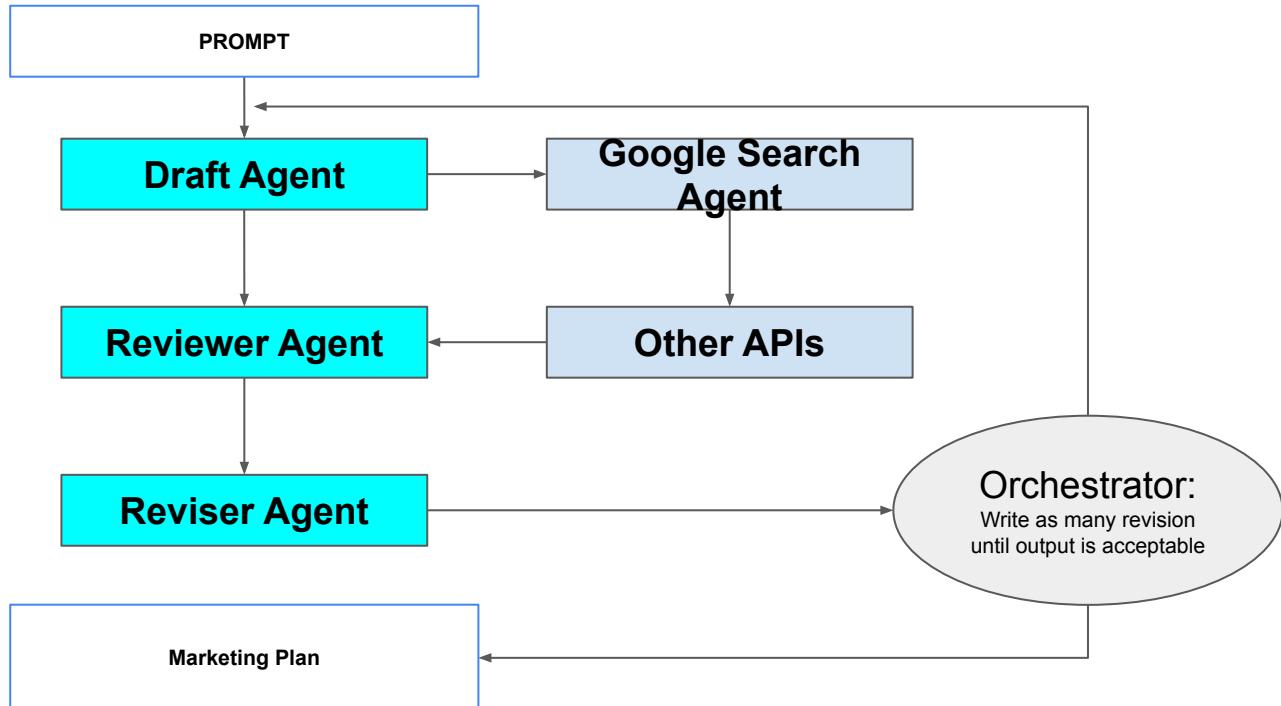
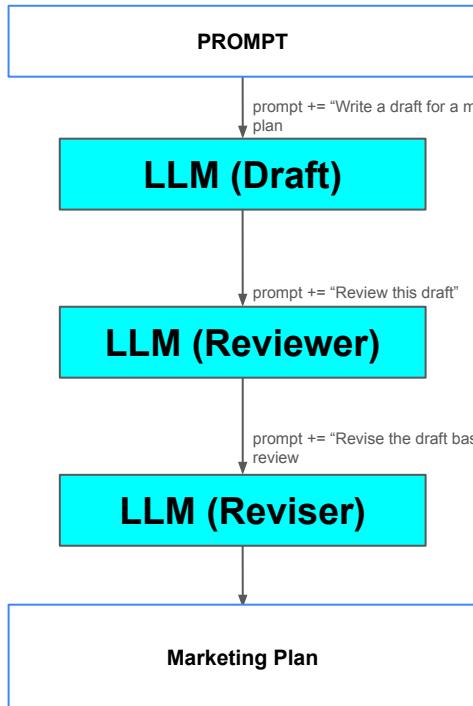


Compounding LLM Approach as an **Agentic AI**

1 Agents

2 Agent can be a tool

3 Path is not predetermined



[Why go through all this?]

Task Adaptability

To deliver better results.

By giving each LLM a specific role, we improve the quality of its context — and in turn, the quality of its output.

Other examples

1. Legal Document Review Agent System

Task: Draft and finalize a legal contract.

Agentic Workflow:

-  **Drafting Agent:** Writes the initial version of a contract.
-  **Compliance Agent:** Reviews for legal inconsistencies, missing clauses, or regulatory gaps.
-  **Risk Assessment Agent:** Flags liability risks or vague terms.
-  **Finalization Agent:** Decides if the contract is ready for use or needs revision.

Why it's agentic:

Each agent has a **defined responsibility**, and the output **loops** until the contract passes the legal bar set by the final decision agent.

2. Research Paper Generator

Task: Generate a literature review or technical research summary.

Agentic Workflow:

-  **Retriever Agent:** Uses RAG to fetch relevant academic papers.
-  **Summarizer Agent:** Extracts key findings and contributions.
-  **Critique Agent:** Reviews the summary for gaps, outdated studies, or bias.
-  **Refinement Agent:** Improves coherence, flow, and adds missing references.
-  **Final Check Agent:** Assesses whether the paper is complete and publish-ready.

Why it's agentic:

This chain transforms **raw data** into a refined academic artifact through **role-specific reasoning**, with **multi-turn refinement**.

Other examples

3. 🛍️ AI Personal Shopper Agent System

Task: Help a user find the perfect outfit for an event based on preferences, budget, and availability.

Agentic Workflow:

-  **User Profiling Agent:** Gathers user style, size, color preferences, and budget.
-  **Trend Scanner Agent:** Analyzes current fashion trends and seasonal styles.
-  **Inventory Scout Agent:** Searches online stores for matching items in stock.
-  **Negotiator Agent:** Applies discount codes or selects optimal bundles.
-  **Decision Agent:** Presents 3 curated outfit options, justifying the selection.

Why it's agentic:

Each agent **acts autonomously** on a **different reasoning layer** (style, availability, price), then composes a practical, real-world action outcome — a shopping decision.

4. 🌟 Disaster Response Planning Agent Team

Task: Generate a real-time action plan in response to a natural disaster alert.

Agentic Workflow:

-  **Geolocation Agent:** Identifies affected regions and population density using real-time maps.
-  **Resource Allocation Agent:** Matches available emergency supplies and personnel to locations in need.
-  **Communication Agent:** Drafts and sends out alerts to local officials and the public.
-  **Time Optimization Agent:** Prioritizes response routes and supply chains based on urgency and accessibility.
-  **Supervisor Agent:** Validates the plan and triggers autonomous deployment or asks for human approval.

Why it's agentic:

This system **acts in the physical world**, not just with words — agents gather, reason, plan, and act, combining LLMs with geospatial tools and real-world data streams.

Other examples

5. 📽️ Video Editor Agent Team (for YouTubers or Content Creators)

Task: Turn a raw video recording into an engaging, publish-ready YouTube video.

Agentic Workflow:

- **🧠 Script Analyzer Agent:** Parses spoken audio or transcript to understand the story structure and key highlights.
- **✂️ Editing Agent:** Cuts silences, filler words, or irrelevant segments; applies pacing adjustments.
- **🎨 Visual Enhancer Agent:** Adds dynamic text overlays, transitions, zoom effects, or B-roll suggestions.
- **🎵 Audio Cleaner Agent:** Reduces background noise, balances volume, and inserts background music based on tone.
- **📋 Summary Agent:** Generates video title, description, tags, and thumbnail caption based on final content.

Why it's agentic:

The agents **collaborate asynchronously** to refine different aspects of the video, and their decisions are **context-aware** — e.g., matching visuals to speech content.

6. 🎙️ Real-Time Podcast Producer Agent

Task: Assist live podcast recording with enhancements and automation.

Agentic Workflow:

- **🗣️ Speech-to-Text Agent:** Transcribes live conversation into text in real-time.
- **🧠 Insight Agent:** Highlights memorable quotes, topic shifts, and emotional tone changes.
- **🔊 Audio Engineering Agent:** Automatically adjusts levels, removes mic pops, and applies compression as you speak.
- **🎙️ Fact-Checking Agent:** Flags inaccurate claims or prompts the host with missing info during breaks.
- **📝 Post-Production Agent:** Generates show notes, timestamps, social media snippets, and a title.

Why it's agentic:

This system doesn't just process — it **listens, thinks, and prepares output** across modalities (speech, text, metadata), enhancing both live quality and downstream publishing.

Demo and Activity

https://colab.research.google.com/drive/1pmY5U2TnEegQ_m_sWbmXALmhu3tazSc6#scrollTo=NieN8Fe7Oi0Y

Business Case Ideation

Build an Agentic Business Case

1. Define the problem

What's the core issue?

Who experiences it?

4. Data Flow - Why is it Agentic?

Map-out the whole process

2. Who is your user?

Who feels this problem?

Ex: *Support reps, Students, HR staff, etc.*

3. Design the Agent Team

Create 3–5 agents with distinct roles.

For each agent, describe:

- Name and function
 - What modality it works with (text, audio, video, etc.)
- What value it adds to the overall workflow

Lecture Audio/Video

--> *Speech-to-Text Agent*

--> *Transcript*

--> *Topic Segmentation Agent*

--> *Segmented Transcript (by topic)*

--> *Insight Extraction Agent*

--> *Key Points, Definitions, Takeaways*

--> *Visual Slide Generator Agent*

--> *Auto-Generated Slides*

--> *Summary & Quiz Agent*

--> *1-Page Summary + Quiz*