

Hands-on GenAI in Action 2025

Session 1: Foundation

Shubin Yu



Dr. Shubin Yu

Email: Yu@hec.fr
Address: W1-407, HEC Paris

Virtual office hours: by appointment only

- ❖ Founder of MimiTalk.app
- ❖ Associate professor of Marketing, HEC Paris, France
- ❖ Associate professor at BI Norwegian Business School, Norway
- ❖ Assistant professor at Peking University, China
- ❖ Ph.D. in Communication Sciences at Ghent University, Belgium
- ❖ Research fields: GenAI for business, Consumer-technology interaction



*Lille My (Mimi)
Girl
3 years old
Coton de tulear
Cheese*

Have you used any GenAI tools?
What are they?
When do you use them?
What do you like to know more about GenAI?



Course description



This intensive two-day course provides business executives with a comprehensive understanding of Generative AI (GenAI) and its transformative potential across various business functions.

You will gain **practical** knowledge of GenAI tools, prompt/context engineering techniques, and real-world applications, enabling you to leverage this cutting-edge technology for strategic decision-making and driving business growth.

Reading list

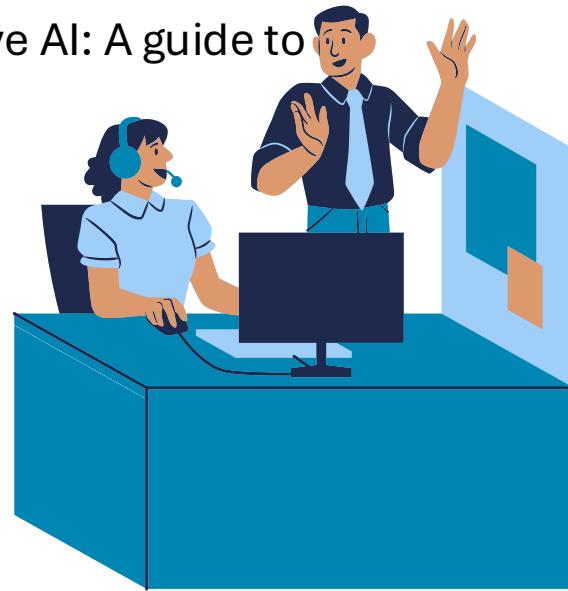
Shubin, Yu (2025), Generative AI for Business, <https://gaiforbusiness.com>

Scott Cook, Andrei Hagiu, and Julian Wright (2024), Turn Generative AI from an Existential Threat into a Competitive Advantage, ***Harvard Business Review***.

H. James Wilson and Paul R. Daugherty (2024), Embracing Gen AI at Work, ***Harvard Business Review***.

Christian Stadler, Martin Reeves (2024), Three Lessons from Chatting about Strategy with ChatGPT, ***Sloan Management Review***.

Andrew McAfee, Daniel Rock, and Erik Brynjolfsson (2024), How to Capitalize on Generative AI: A guide to realizing its benefits while limiting its risks, ***Harvard Business Review***.



Course Objectives:

- Define GenAI and differentiate it from other types of AI.
- Articulate the history, evolution, and future trends of GenAI.
- Identify key GenAI models, tools, and their capabilities.
- Understand the principles and best practices of prompt/context engineering.
- Apply GenAI tools for specific business tasks like information search, business communication, and content generation.
- Analyze the ethical implications and potential risks associated with GenAI.



Practical information

Time

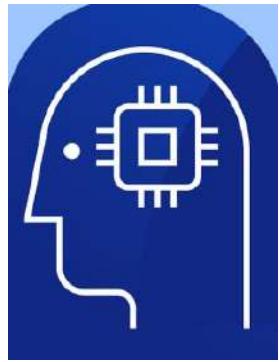
2 days

Course activities

Lecture, case studies, and hands-on exercises, group activities

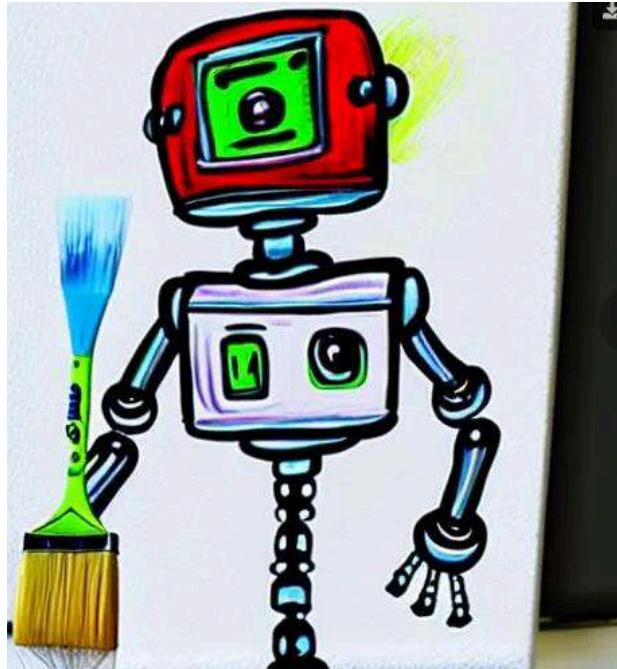
<https://github.com/lanceyuu/GenAI-for-Business>





Introduction to Generative AI

What is Generative AI (GenAI)



A robot is painting (March 2023 vs Jan 2024)



Generative AI refers to artificial intelligence systems that can create new content, such as text, images, audio, or video, based on patterns learned from existing data.

<https://openart.ai/create>

<https://hotpot.ai/art-generator>



Future workplace:

Collaboration with GenAI for various tasks

- Content Creation and Marketing
 - Software Development
 - Legal Work
 - Human Resources
 - Project Management
 - Data Analysis
 - Financial Analysis
 - Customer Service
 - Education
 - Training
 - Consulting
- ...

The impact of GenAI (productivity)

- OpenAI's research estimates that 80% of current work activities can integrate generative AI technologies and capabilities.
- McKinsey's research indicates that generative AI and other technologies have the potential to automate 60% to 70% of the tasks that currently occupy employees' time.
- Accenture reports that GenAI can “transform productivity,” with modeling indicating it can save over 12% of working hours for the average company while improving output quality by 8.5%
- BCG consultants completed tasks 25% more quickly

Evidence from an increasing number of academic studies

Coding

Software engineers code up to twice as fast using Codex
Peng et al. (2023)

Customer service

AI assistance helps customer-support agents work faster, increasing issues resolved per hour by an average of 15%
Brynjolfsson et al., (2025)

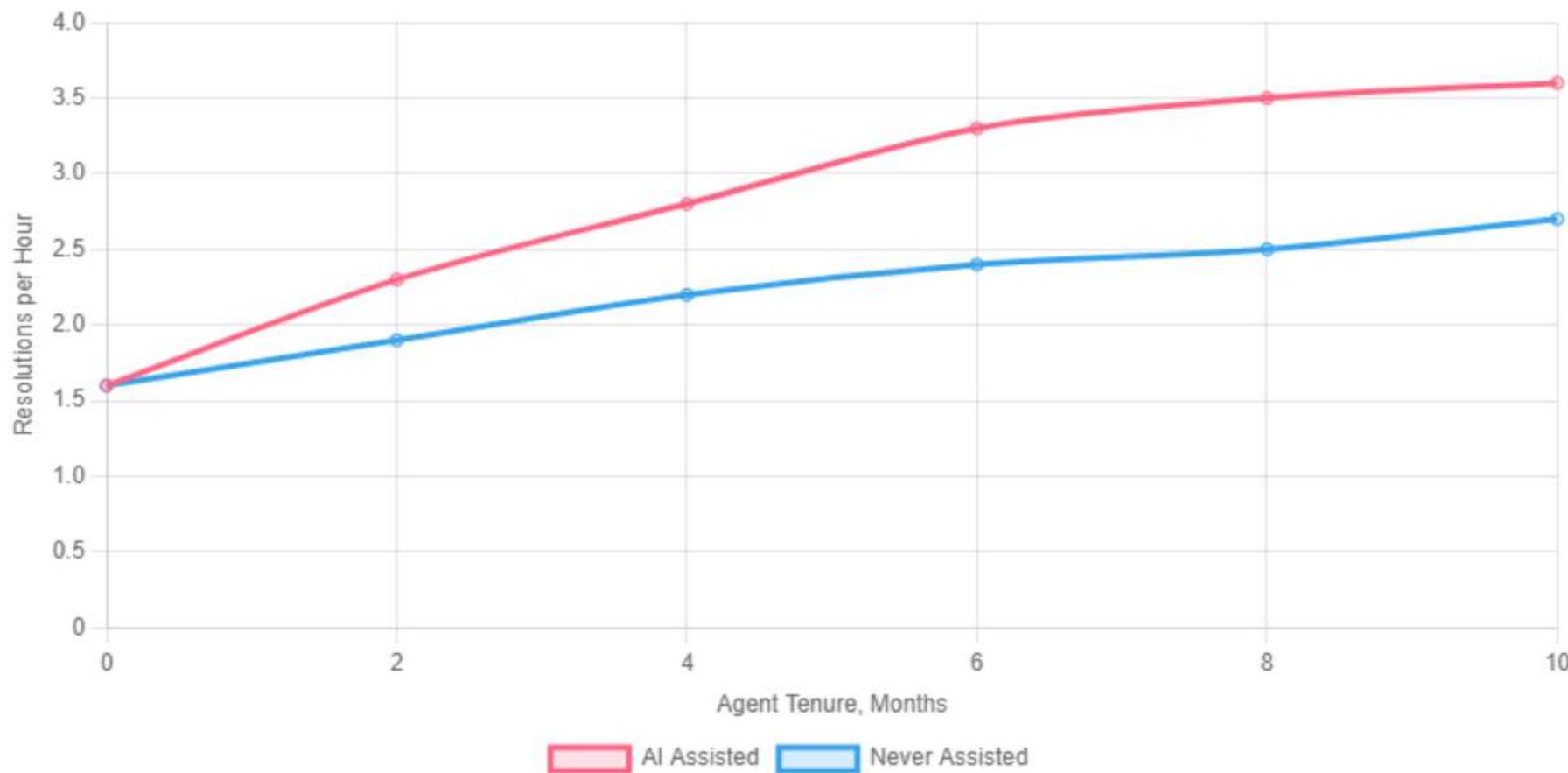
Writing

Writing tasks completed twice as fast
Noy and Zhang (2023)

Diagnosis

Overall reading times shortened when radiologists used AI
Shin et al. (2023)

AI assistance helps newer agents "catch up"



The impact of GenAI (performance)

- **Bain & Company** highlights early successes in sales: AI tools that guide sales reps or generate tailored pitches have delivered **>30% improvements in sales win rates** in pilot programs.
- BCG's experiment found that for a creative ideation task (developing a new product concept and go-to-market plan), **90% of consultants using GPT-4 outperformed those without it**, achieving solutions roughly **40% higher in quality** as rated by experts

Jensen Huang, CEO Nvidia

“ AI IS NOT GOING TO TAKE YOUR JOB. THE PERSON WHO USES AI WILL. USE AI AS FAST AS YOU CAN, SO THAT YOU CAN STAY GAINFULLY EMPLOYED.

”

What is the difference between ChatGPT and GPT?





How is a large language model created?

- 1.Pre-training:
 - Trained on vast amounts of text data from the internet, books, and other sources.
 - Learn patterns, relationships, and structures in language without specific tasks.
- 2.Unsupervised Learning:
 - During pre-training, the model predicts the next word given the previous words.
 - This allows it to learn grammar, facts, reasoning, and even some level of common sense.
- 3.Tokenization:
 - Text is broken down into tokens (words or subwords).
 - The model processes and generates text token by token.
- 4.Contextual Understanding:
 - Uses attention mechanisms to weigh the importance of different parts of the input.
 - Can understand context over long sequences of text.
- 5.Reinforce learning:
 - Learning through trial and error with rewards and penalties.
- 6.Fine-tuning:
 - Can be further trained on specific tasks or domains for better performance.

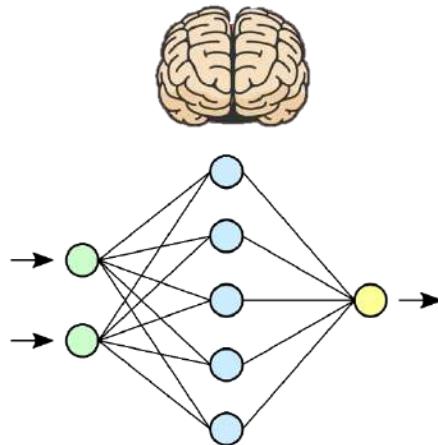
How is a large language model created?

- 1.Pre-training:
 - Trained on vast amounts of text data from the internet, books, and other sources.
 - Learn patterns, relationships, and structures in language without specific tasks.
- 2.Unsupervised Learning:
 - During pre-training, the model predicts the next word given the previous words.
 - This allows it to learn grammar, facts, reasoning, and even some level of common sense.
- 3.Tokenization:
 - Text is broken down into tokens (words or subwords).
 - The model processes and generates text token by token.
- 4.Contextual Understanding:
 - Uses attention mechanisms to weigh the importance of different parts of the input.
 - Can understand context over long sequences of text.
- 5.Reinforce learning:
 - Learning through trial and error with rewards and penalties.
- 6.Fine-tuning:
 - Can be further trained on specific tasks or domains for better performance.

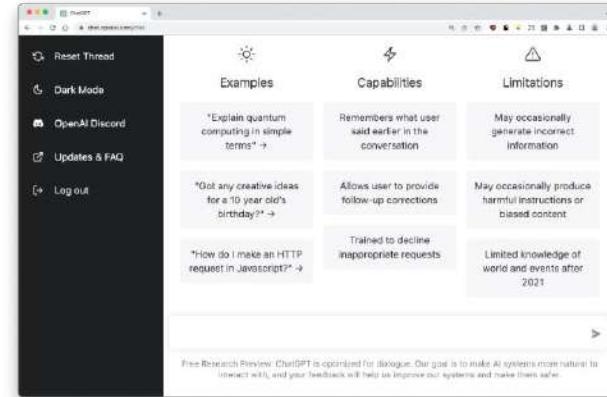


ChatGPT

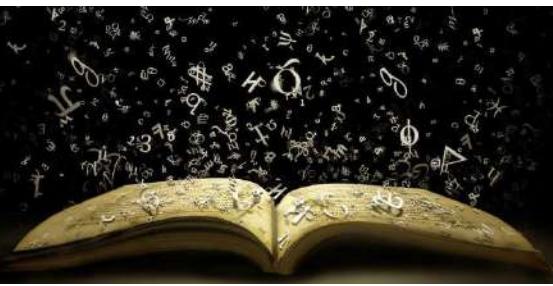
Text from web pages, books, articles, and social gossip
Unsupervised learning



GPT is a series of language models based on neural network architecture.



ChatGPT stands for "Chat Generative Pre-trained Transformer". It is an AI language model developed by OpenAI. The main purpose of ChatGPT is to generate human-like text responses in a conversational manner.



Text from web pages, books,
articles, and social gossip
Unsupervised learning



How is a large language model created?

1. Pre-training:
 - Trained on vast amounts of text data from the internet, books, and other sources.
 - Learn patterns, relationships, and structures in language without specific tasks.
2. Tokenization and embedding:
 - Text is broken down into tokens (words or subwords).
 - The model embeds each word into a vector of numbers that represents its meaning.
3. Unsupervised Learning:
 - During pre-training, the model predicts the next word given the previous words.
 - This allows it to learn grammar, facts, reasoning, and even some level of common sense.
4. Contextual Understanding:
 - Uses attention mechanisms to weigh the importance of different parts of the input.
 - Can understand context over long sequences of text.
5. Reinforce learning:
 - Learning through trial and error with rewards and penalties.
6. Fine-tuning:
 - Can be further trained on specific tasks or domains for better performance.

Input

My dog loves cheese

Tokenization

Break down into small pieces (words or parts of words)

Tokens



Embedding

Turn Token into numerical representation
Capturing their meaning

Embedding

My

0.12,	-0.05,	0.88...
0.33,	-0.06,	0.01...
0.24,	-0.05,	0.88...

Dog

0.12,	-0.05,	0.88...
-34,	0.76,	0.01...
	-0.06	

Loves

0.05	
0.07	
0.26	

Cheese

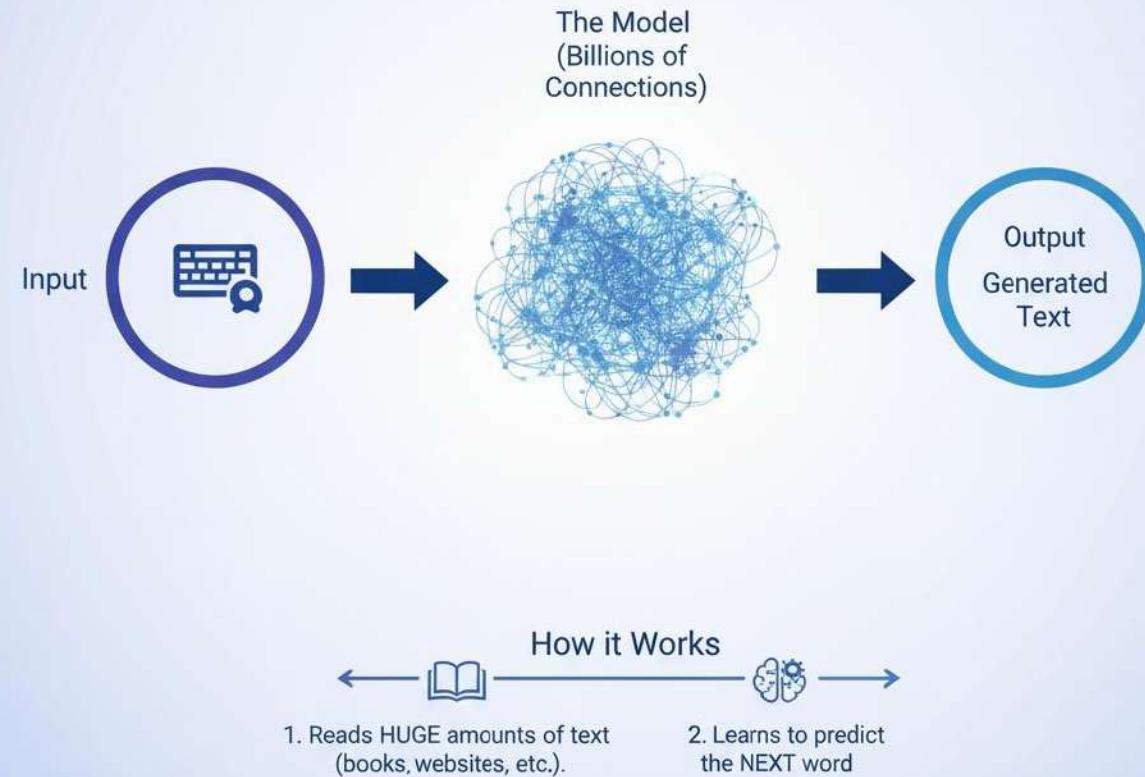
0.98	
-0.45	
1.17	

How is a large language model created?

1. Pre-training:
 - Trained on vast amounts of text data from the internet, books, and other sources.
 - Learn patterns, relationships, and structures in language without specific tasks.
2. Tokenization:
 - Text is broken down into tokens (words or subwords).
 - The model processes and generates text token by token.
3. Unsupervised Learning:
 - During pre-training, the model predicts the next word given the previous words.
 - This allows it to learn grammar, facts, reasoning, and even some level of common sense.
4. Contextual Understanding:
 - Uses attention mechanisms to weigh the importance of different parts of the input.
 - Can understand context over long sequences of text.
5. Reinforce learning:
 - Learning through trial and error with rewards and penalties.
6. Fine-tuning:
 - Can be further trained on specific tasks or domains for better performance.

Neural network

Large Language Model





User Input: "What is the capital of France?"
System processes this and model begins generation...



clojure

Input sequence: "What is the capital of France?"
Model processes → Output layer probabilities:
- "The" (0.05)
- "Paris" (0.78) ← HIGHEST
- "France" (0.02)
- "It" (0.08)
- ...
Selected token: "Paris"

closure



Input sequence: "What is the capital of France? Paris"

Model processes → Output layer probabilities:

- "is" (0.72) ← HIGHEST
- "," (0.15)
- "." (0.08)
- "was" (0.03)
- ...

Selected token: "is"

closure



Input sequence: "What is the capital of France? Paris is"

Model processes → Output layer probabilities:

- "the" (0.85) ← HIGHEST
- "a" (0.08)
- "located" (0.04)
- "indeed" (0.02)
- ...

Selected token: "the"



Input sequence: "What is the capital of France? Paris is the"

Model processes → Output layer probabilities:

- "capital" (0.91) ← HIGHEST
- "largest" (0.04)
- "main" (0.03)
- ...

Selected token: "capital"



Input sequence: "What is the capital of France? Paris is the capital"

Model processes → Output layer probabilities:

- "of" (0.88) ← HIGHEST
- "city" (0.07)
- "." (0.03)
- ...

Selected token: "of"



Input sequence: "What is the capital of France? Paris is the capital of"

Model processes → Output layer probabilities:

- "France" (0.94) ← HIGHEST
- "the" (0.03)
- "this" (0.02)
- ...

Selected token: "France"



Input sequence: "What is the capital of France? Paris is the capital of France"

Model processes → Output layer probabilities:

- "." (0.89) ← HIGHEST
- "," (0.06)
- "and" (0.03)
- ...

Selected token:."

All of these go into the input layer:

1. User Prompt

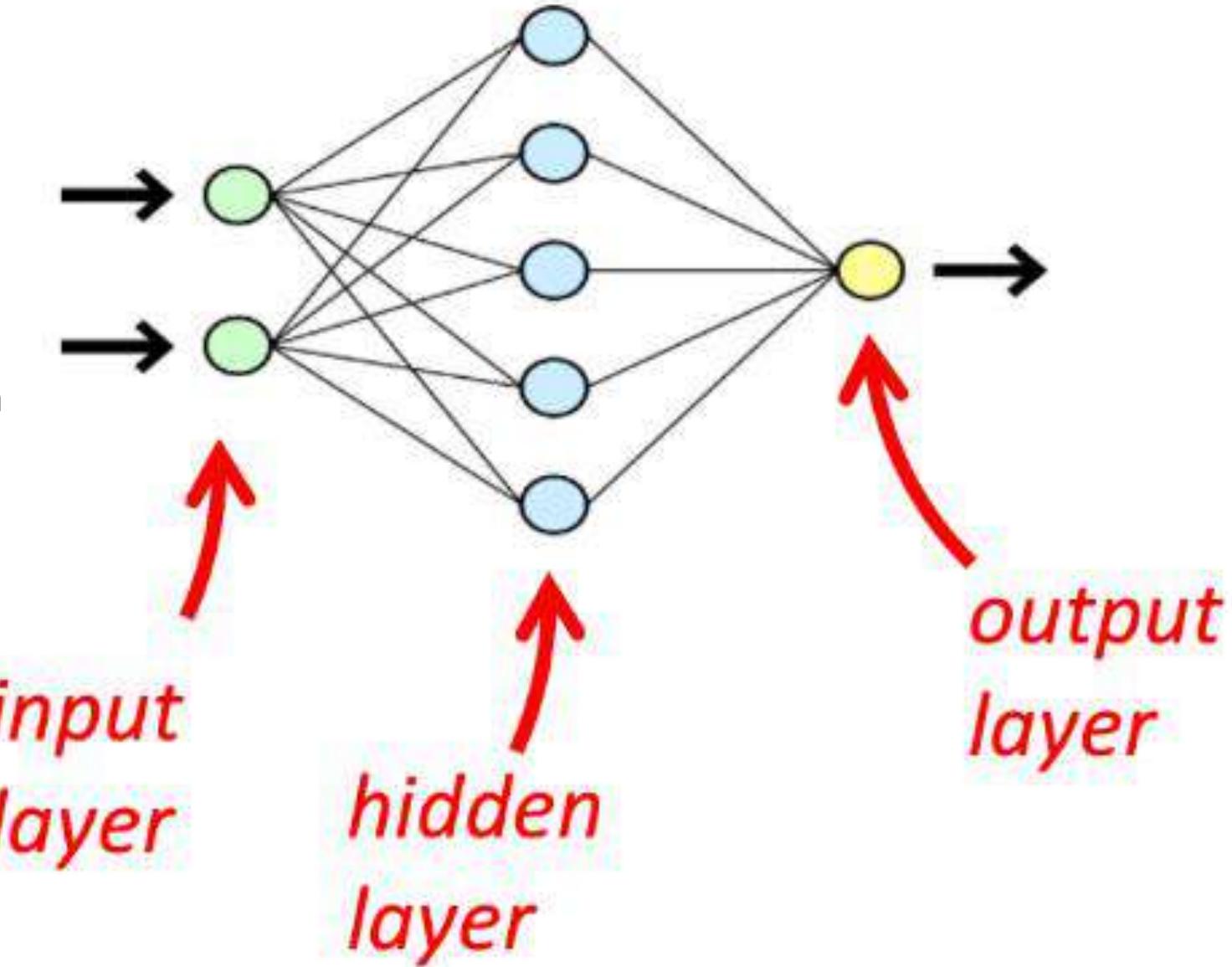
1. "Write a poem about cats"
2. "Explain quantum physics"

2. System Prompt

1. "You are a helpful assistant"
2. "Answer in French"
3. "Be concise and professional"

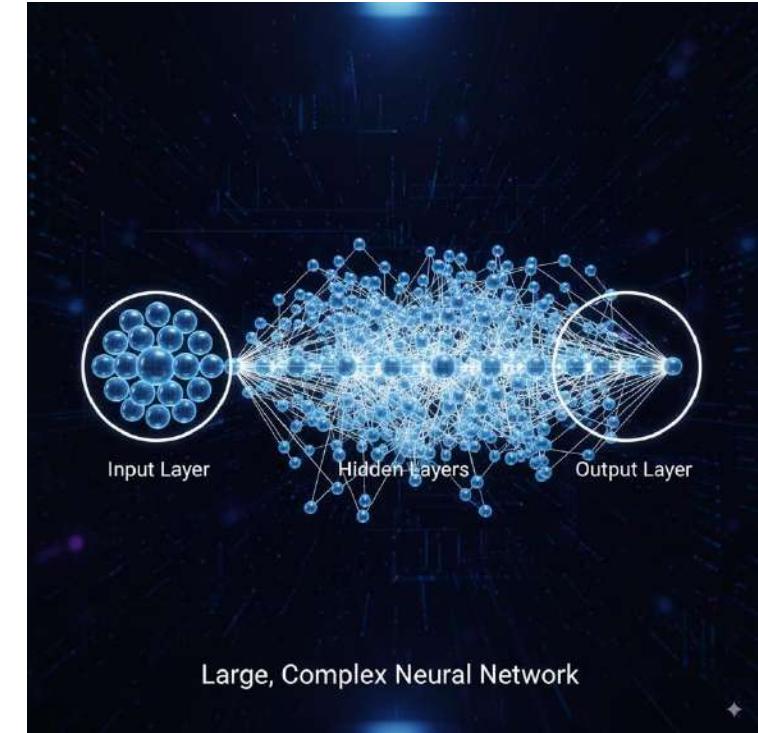
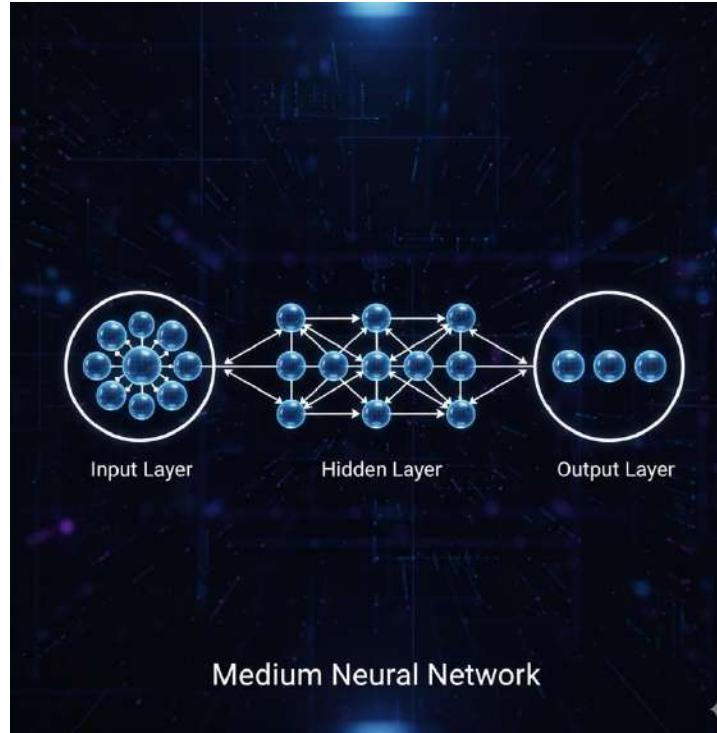
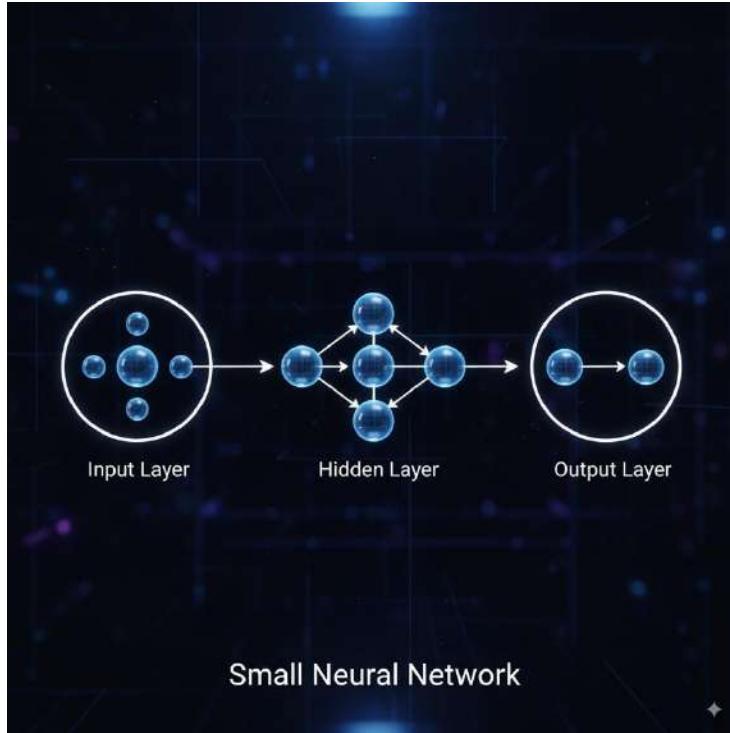
3. Context/Conversation History

1. Previous messages in the conversation
2. Background information
3. Few-shot examples



$$\text{output} = \text{weight}_1 \times \text{input}_1 + \text{weight}_2 \times \text{input}_2 + \dots + \text{bias}$$

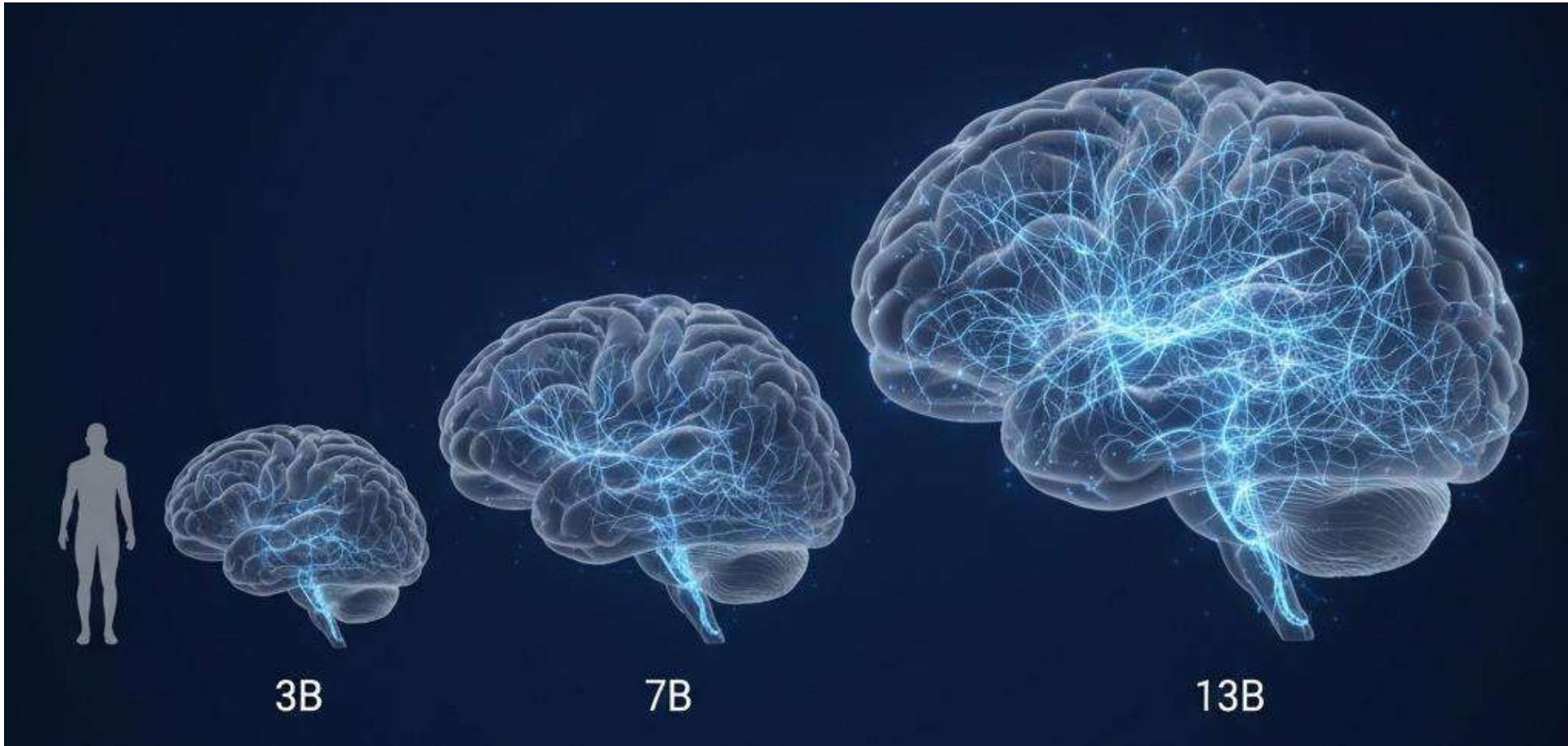
Concept: Model parameters



Model parameters are numerical weights and biases in the neural network that encode:

- Language patterns
- Facts and relationships
- Reasoning abilities
- Response preferences

Concept: Model parameters



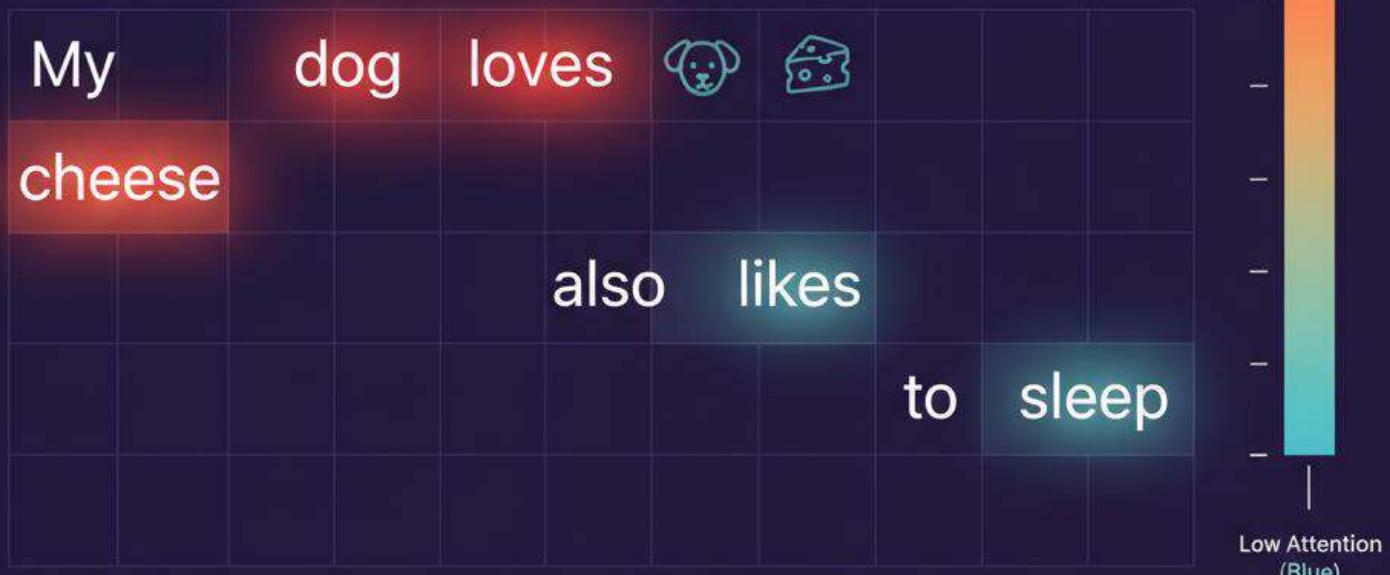
Like synapses in a human brain
Each connection has a "strength" (the parameter value)
7B parameters = 7 billion connection strengths
More parameters \approx more capacity to learn

How is a large language model created?

- 1.Pre-training:
 - Trained on vast amounts of text data from the internet, books, and other sources.
 - Learn patterns, relationships, and structures in language without specific tasks.
- 2.Tokenization:
 - Text is broken down into tokens (words or subwords).
 - The model processes and generates text token by token.
- 3.Unsupervised Learning:
 - During pre-training, the model predicts the next word given the previous words.
 - This allows it to learn grammar, facts, reasoning, and even some level of common sense.
- 4.Contextual Understanding:
 - Uses attention mechanisms to weigh the importance of different parts of the input.
 - Can understand context over long sequences of text.
- 5.Reinforce learning:
 - Learning through trial and error with rewards and penalties.
- 6.Fine-tuning:
 - Can be further trained on specific tasks or domains for better performance.

LLM ATTENTION HEATMAP: FOCUS ON "CHEESE"

My dog loves cheese and also likes to sleep.



How is a large language model created?

- 1.Pre-training:
 - Trained on vast amounts of text data from the internet, books, and other sources.
 - Learn patterns, relationships, and structures in language without specific tasks.
- 2.Tokenization:
 - Text is broken down into tokens (words or subwords).
 - The model processes and generates text token by token.
- 3.Unsupervised Learning:
 - During pre-training, the model predicts the next word given the previous words.
 - This allows it to learn grammar, facts, reasoning, and even some level of common sense.
- 4.Contextual Understanding:
 - Uses attention mechanisms to weigh the importance of different parts of the input.
 - Can understand context over long sequences of text.
- 5.Reinforce learning:
 - Learning through trial and error with rewards and penalties.
- 6.Fine-tuning:
 - Can be further trained on specific tasks or domains for better performance.



Reinforce learning 🐕

Training with rewards & repetition

How is a large language model created?

- 1.Pre-training:
 - Trained on vast amounts of text data from the internet, books, and other sources.
 - Learn patterns, relationships, and structures in language without specific tasks.
- 2.Tokenization:
 - Text is broken down into tokens (words or subwords).
 - The model processes and generates text token by token.
- 3.Unsupervised Learning:
 - During pre-training, the model predicts the next word given the previous words.
 - This allows it to learn grammar, facts, reasoning, and even some level of common sense.
- 4.Contextual Understanding:
 - Uses attention mechanisms to weigh the importance of different parts of the input.
 - Can understand context over long sequences of text.
- 5.Reinforce learning:
 - Learning through trial and error with rewards and penalties.
- 6.Fine-tuning:
 - Can be further trained on specific tasks or domains for better performance..



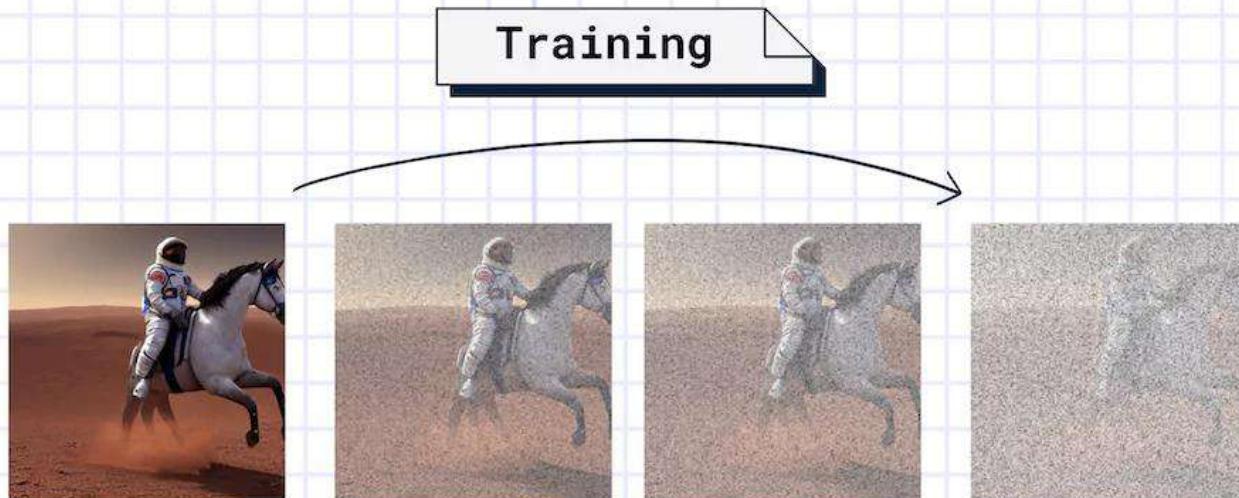
Just as a dog can be further trained for specific tasks, a large language model can be fine-tuned for various domains such as medical diagnosis, email writing, and financial advice.



What about images and sounds? How do multimodal LLMs work?



HOW AI IMAGE GENERATION WORKS



• Hypotenuse AI

Models and techniques

Transformer

Generative Adversarial Networks (GANs)

Diffusion models

Variational Autoencoders (VAEs)

Retrieval-Augmented Generation (RAG)

Recurrent Neural Networks (RNNs)

Autoregressive Models

Convolutional neural networks (CNNs)

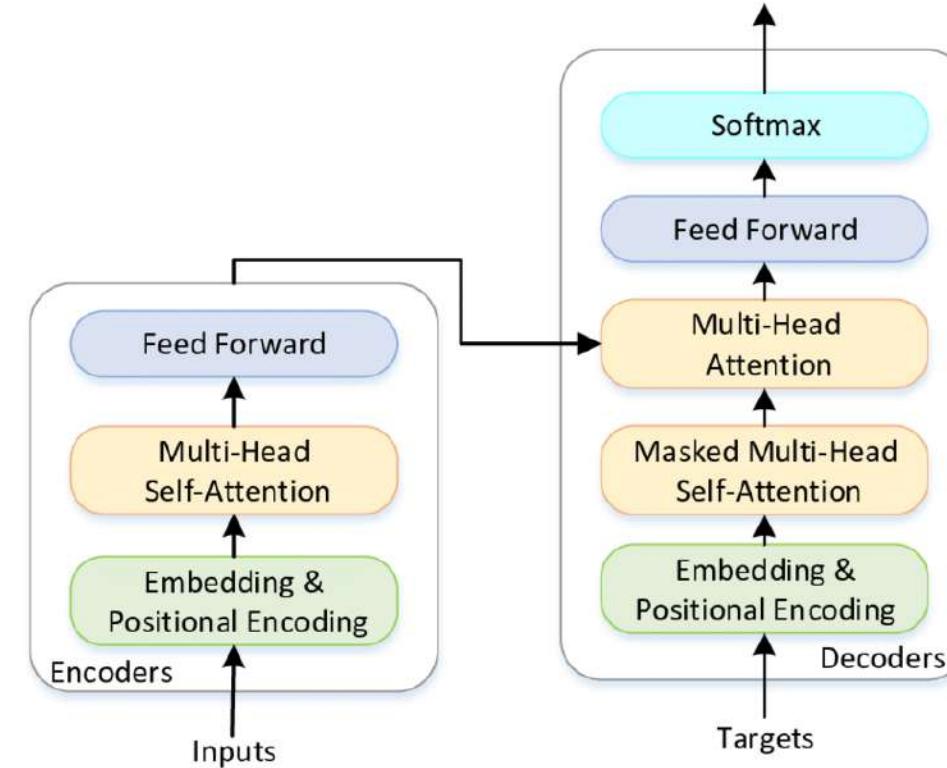
Memory Neural Networks (MNNs)

...

Understanding Transformer Models: The Orchestra Analogy

Transformer Models explained:

- **The Orchestra (Input):** Each word is like a musician
- **The Conductor (Attention Mechanism):**
 - Decides which musicians (words) to focus on
 - Coordinates how they interact with each other
- **The Performance (Processing):**
 - Each musician plays in context of others
 - Creates a harmonious output (understanding)
- **Multiple Conductors (Multi-Head Attention):**
 - Different conductors focus on various aspects
 - Combines multiple perspectives for rich understanding
- **Result:** Ability to understand complex language contexts and generate coherent responses
This allows for powerful language understanding and generation capabilities.



LEARN MORE

<https://blogs.nvidia.com/blog/what-is-a-transformer-model/>

Understanding GANs: The Art Forgery Analogy

Generative Adversarial Networks (GANs)

explained:

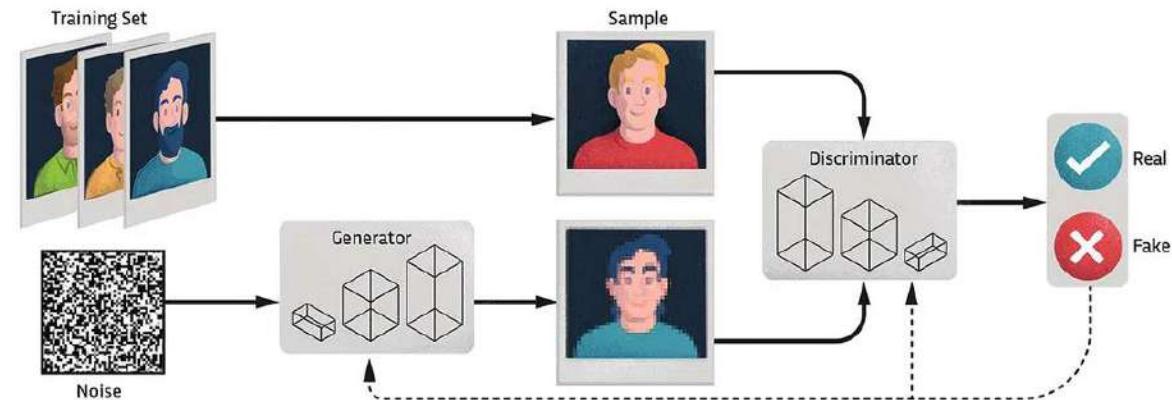
- **The Forger (Generator):** Creates fake artwork
- **The Detective (Discriminator):** Tries to spot fakes

• The Process:

- Forger creates fake art
- Detective examines both real and fake art
- Forger learns from feedback to improve fakes
- Detective gets better at spotting fakes

- **Result:** Forger becomes so good that fakes are indistinguishable from real art

This competitive process leads to the creation of highly realistic artificial data or images.



LEARN MORE

<https://www.oreilly.com/content/generative-adversarial-networks-for-beginners/>

Understanding Diffusion Models: The Dust Cloud Analogy

Diffusion Models explained:

•The Process:

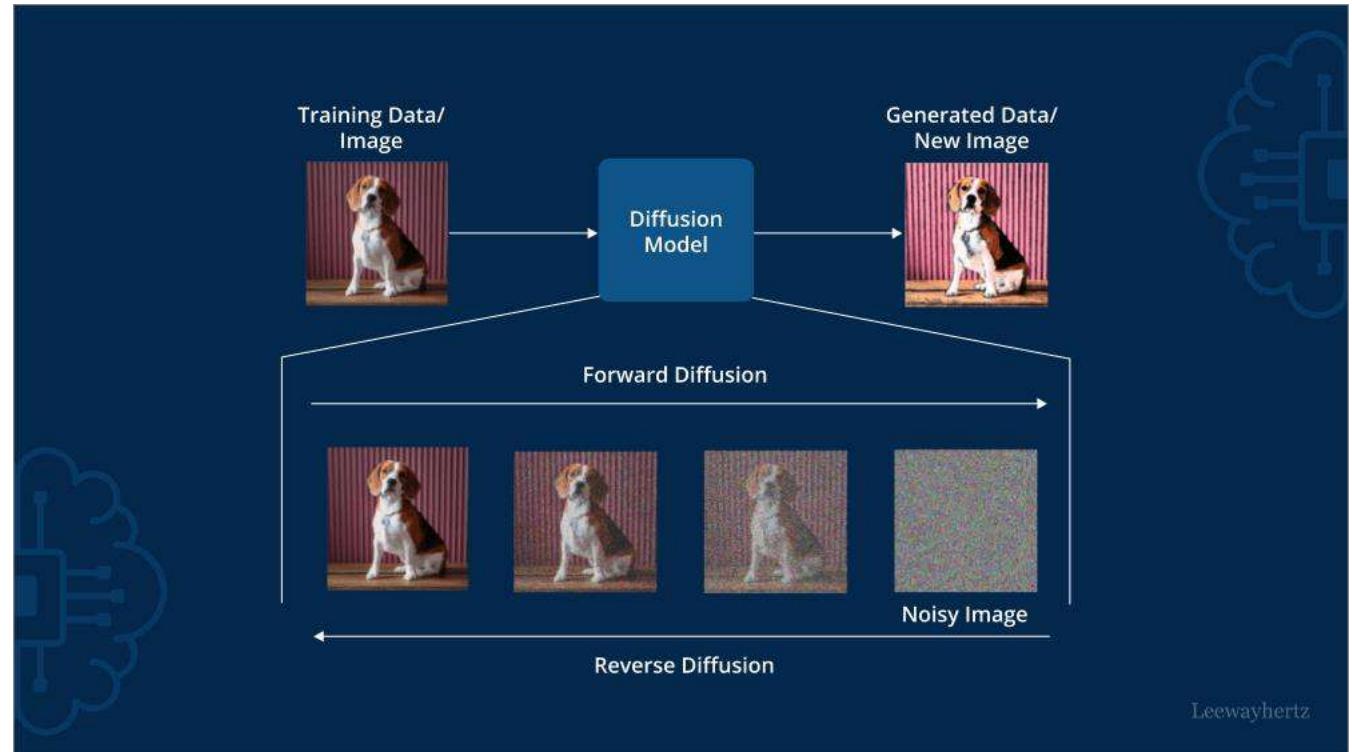
- Start with a clear image
- Gradually add "noise" (like dust)
- Image becomes increasingly blurry
- AI learns to reverse this process

•Generation:

- Begin with pure noise (dust cloud)
- AI gradually removes "dust"
- A clear image emerges step-by-step

•**Result:** AI can create new, realistic images by "de-noising" random noise

This process allows for controlled, high-quality image generation.



LEARN MORE

<https://www.leewayhertz.com/diffusion-models/>

Leewayhertz

Understanding VAEs: The Compression and Decompression Analogy

Variational Autoencoders (VAEs) explained:

•The Encoder (Compressor):

- Takes input data (e.g., images)
- Compresses it into a compact representation

•The Latent Space (Compressed File):

- A compact, probabilistic representation
- Captures essential features of the data

•The Decoder (Decompressor):

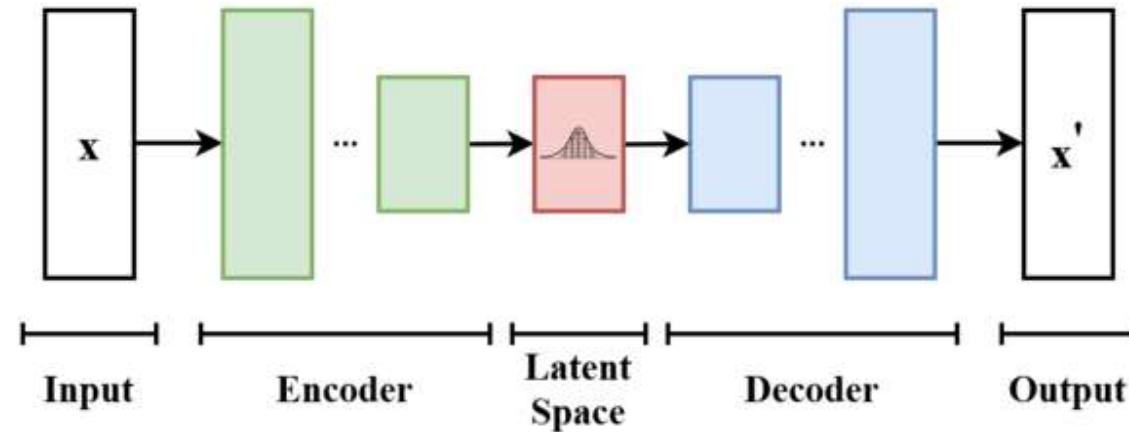
- Takes the compressed representation
- Reconstructs the original-like data

•The Variational Part:

- Adds controlled randomness to compression
- Enables generation of new, similar data

•**Result:** Can reconstruct existing data and generate new, similar data

This allows for both data compression and creative generation capabilities.



Understanding RAG: The Librarian and Author Analogy

Retrieval-Augmented Generation (RAG) explained:

•The Librarian (Retrieval Component):

- Searches a vast library of information
- Finds relevant books or articles for a given topic

•The Author (Generation Component):

- Uses the retrieved information
- Writes original content based on the research

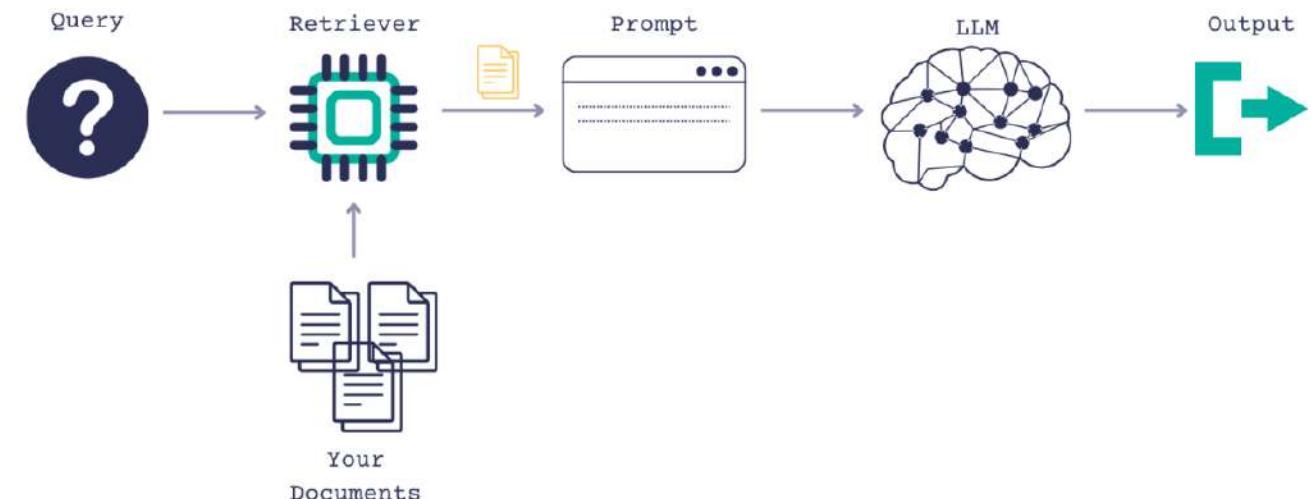
•The Process:

- Receive a question or topic
- Librarian retrieves relevant information
- Author creates response using retrieved info and own knowledge

•The Knowledge Base (Library):

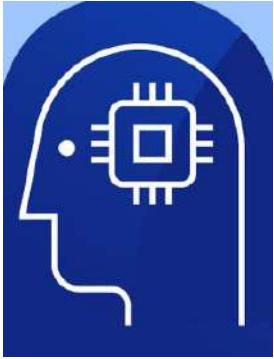
- Can be updated with new information
- Allows for up-to-date and factual responses

•**Result:** Generates informed, accurate responses by combining retrieved information with language generation



LEARN MORE

<https://aws.amazon.com/what-is/retrieval-augmented-generation/>



The history of Generative AI

The history of GenAI

1960s: ELIZA chatbot

1980-1990s: Development of neural networks

2000s: Rise of deep learning

2014: Introduction of Generative Adversarial Networks (GANs)

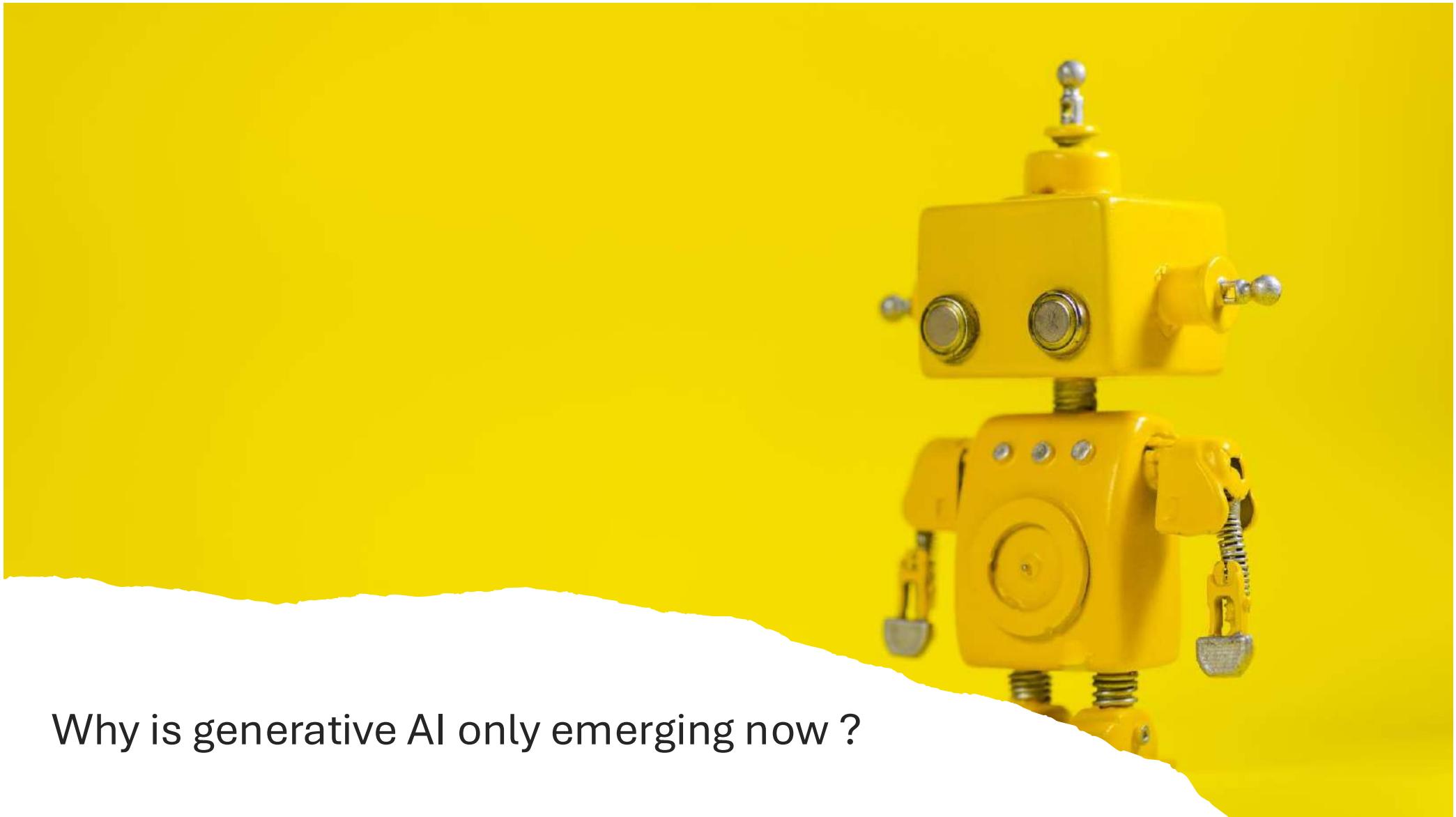
2015: Emergence of Diffusion models

2020: Release of GPT-3

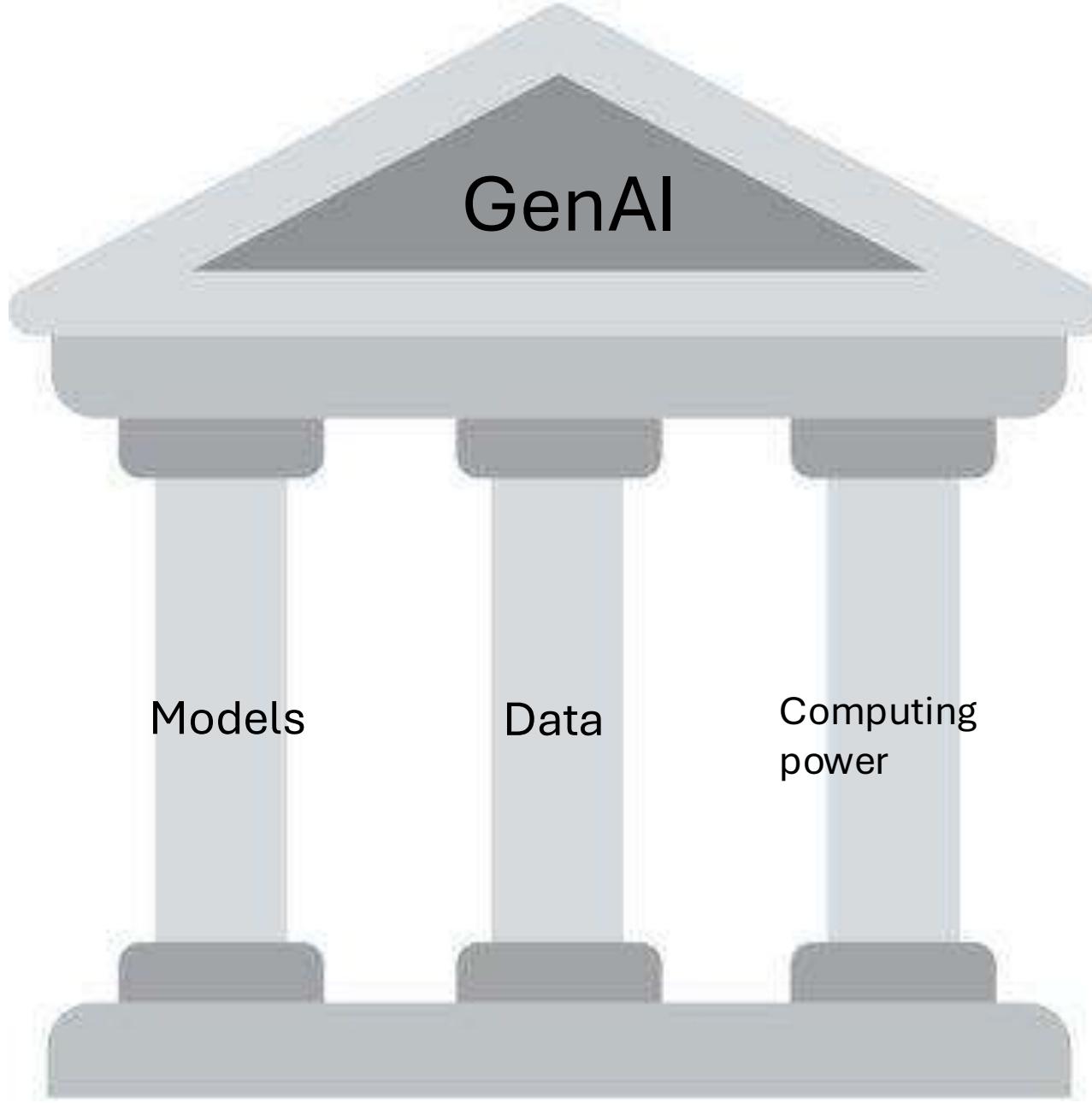
2022: Release of GPT-3.5

2023: Explosion of various GenAI models (Claude, Gemini)

2025: Advanced open-source models (Deepseek r1, Llama, Qwen)



Why is generative AI only emerging now ?



Generative AI: How are the three pillars at present?



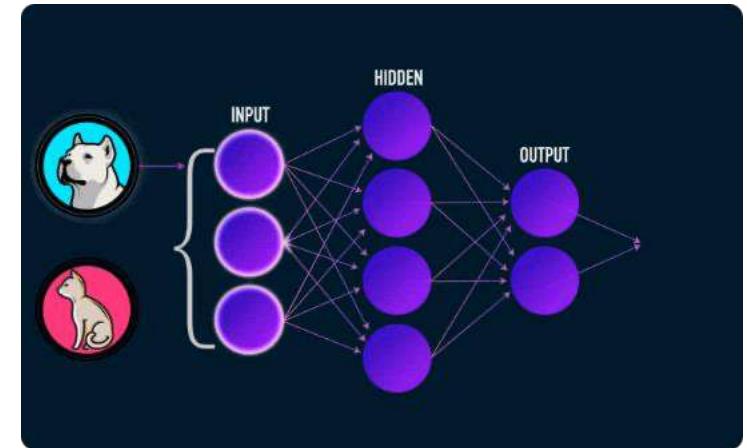
Abundant Data?

AI thrives on data, and the exponential growth of data availability fuels its capabilities.



Strong Computing Power

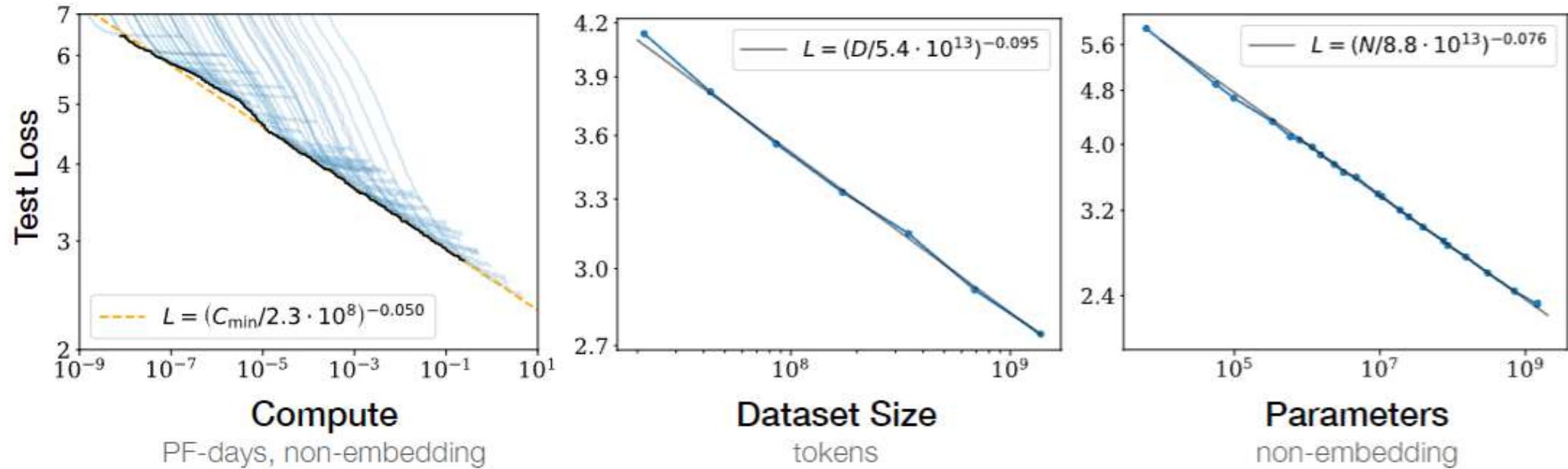
Powerful CPUs and GPUs enable AI models to process massive datasets and generate complex outputs.



Sophisticated AI Models

Generative models like GANs, transformers, and diffusion models are revolutionizing personalized marketing.

- Language modeling performance improves smoothly as we increase the model size, dataset size, and amount of compute used for training. For optimal performance all three factors need to scale up in tandem. Empirical performance has a power-law relationship with each individual factor when not bottlenecked by the other two.
- arXiv: Scaling Laws for Neural Language Models





The current landscape of GenAI models





OpenAI GPT

ChatGPT 4o ▾

What can I help with?

Ask anything

+ Search ...

Create image Get advice Analyze data Make a plan More



GPT 5

Search
Deep Search
Live mode
Canva
Operator
Customized GPT



Gemini

 Gemini ▾
2.0 Flash



Gemini 2.5 Pro

Efficient (low cost
and fast)
Notebook LM
Co-scientist
Learn about

Hello, Michal

Ask Gemini



Deepseek

deepseek

New chat

No chat history

Get App NEW

R1
v3

Open source
Low cost

Hi, I'm DeepSeek.

How can I help you today?

Message DeepSeek

DeepThink (R1) Search

0 ↑

Claude

Good afternoon, Lance

How can Claude help you today?

Claude 3.7 Sonnet ▾ Choose style ▾

Collaborate with Claude using documents, images, and more

Provide stakeholder perspective Extract insights from report Polish your prose

NEW Analysis tool

Upload CSVs for Claude to analyze quantitative data with high accuracy and create interactive data visualizations. [Try it out](#)

Sonnet 4.0
Opus 4.1

Programming
Writing
Low cost
App preview



Grok 4

The screenshot shows the Grok 4 AI interface. At the top left is the Grok logo. On the right are three icons: a gear, a user profile, and a sign-in button. Below the header is a large central text area that says "Welcome to Grok. How can I help you today?". In the middle-left is a search bar with the placeholder "What do you want to know?". Below the search bar are four buttons: a white button with a blue icon, a blue button labeled "DeepSearch", a white button labeled "Think", and a white button with a blue icon. To the right of the search bar is a dropdown menu labeled "Grok 3" with a downward arrow. A small upward arrow icon is positioned next to it. At the bottom of the interface is a dark banner containing the text "Grok 3 Enabled" and "You have the most powerful model available". Below this banner are three buttons: "Research" (with a document icon), "Brainstorm" (with a lightning bolt icon), and "Analyze Data" (with a bar chart icon). At the very bottom, a small note states: "By messaging Grok, you agree to our [Terms](#) and [Privacy Policy](#)".

Text generation performance

Rank (UB)	Model	Score	95% CI (±)	Votes	Organization	License
1	gemini-2.5-pro	1456	±5	35,405	Google	Proprietary
1	gpt-5-high	1447	±7	11,405	OpenAI	Proprietary
1	claude-opus-4-1-20250805-thinking-16k	1447	±7	8,615	Anthropic	Proprietary
2	o3-2025-04-16	1444	±4	40,935	OpenAI	Proprietary
2	chatgpt-4o-latest-20250326	1443	±4	36,773	OpenAI	Proprietary
2	gpt-4.5-preview-2025-02-27	1439	±6	15,271	OpenAI	Proprietary
2	claude-opus-4-1-20250805	1436	±6	11,548	Anthropic	Proprietary
7	gpt-5-chat	1426	±7	8,585	OpenAI	Proprietary
8	grok-4-0709	1422	±6	18,239	xAI	Proprietary
8	kimi-k2-0711-preview	1421	±5	18,588	Moonshot	Modified MIT
8	deepseek-v3.1	1419	±9	4,844	DeepSeek	MIT
8	claude-opus-4-20250514-thinking-16k	1419	±5	23,771	Anthropic	Proprietary
8	qwen3-235b-a22b-instruct-2507	1419	±6	12,971	Alibaba	Apache 2.0
8	deepseek-r1-0528	1417	±6	21,287	DeepSeek	MIT
8	deepseek-v3.1-thinking	1415	±9	4,285	DeepSeek	MIT
8	mistral-medium-2508	1412	±7	7,388	Mistral	Proprietary
9	glm-4.5	1410	±6	11,120	Z.ai	MIT

Generative AI Applications

Text Generation:

- Utilizes large language models to generate contextually relevant text
- Can be used for tasks such as dialogue, explanation, summarization, etc.

Image Generation:

- Uses techniques like GANs and VAEs to generate high-quality, realistic images
- Applied in fields such as art, design, entertainment, etc.

Audio Generation:

- Creates music, text-to-speech, synthesized voices
- Applied in media, entertainment, education, and other fields

Generative AI Applications

Video Generation:

- Creates dynamic videos based on text descriptions or images
- Applied in fields such as art, entertainment, education, healthcare, etc.

Code Generation:

- Generates code snippets, functions, or complete programs
- Assists in software development, debugging, and testing

Data Generation and Augmentation:

- Generates synthetic data, enhances existing datasets
- Applied in healthcare, gaming, education, autonomous driving, and other fields

Virtual World Creation:

- Creates realistic virtual environments and virtual characters
- Applied in gaming, entertainment, education, metaverse platforms, etc.

Tool: LLM arena

LMArena ▾

New Chat

Leaderboard

Take your chats anywhere
Create an account to save your chat history across your devices.

Login

Send Feedback

Report Bugs

Terms of Use Privacy Policy Cookies

Overview Text WebDev Vision Text-to-Image Image Edit Search Text-to-Video Image-to-Video Copilot Start Voting

Text-to-Image Arena

Compare LLMs based on their ability to generate images that match text descriptions.

Last Updated Aug 25, 2025 Total Votes 1,752,043 Total Models 22

Overall

Rank (UB) ↑ Model ↑ Score ↑ 95% CI (±) ↑ Votes ↑ Organization ↑ License ↑

1	gemini-2.5-flash-image-preview (nano-banana)	1147	±2	220,674	Google	Proprietary
2	imagen-4.0-ultra-generate-preview-06-06	1135	±2	193,895	Google	Proprietary
3	gpt-image-1	1129	±3	128,710	OpenAI	Proprietary
4	imagen-4.0-generate-preview-06-06	1119	±2	196,696	Google	Proprietary
5	qwen-image-prompt-extend	1082	±3	123,596	Alibaba	Apache 2.0
5	seedream-3	1077	±3	159,028	Bytedance	Proprietary
6	flux-1-kontext-max	1075	±3	78,017	Black Fores...	Proprietary
8	imagen-3.0-generate-002	1062	±3	256,225	Google	Proprietary
9	flux-1-kontext-pro	1056	±2	165,377	Black Fores...	Proprietary

Hands-on GenAI in Action 2025

Session 2: Communicating with GenAI

Shubin Yu

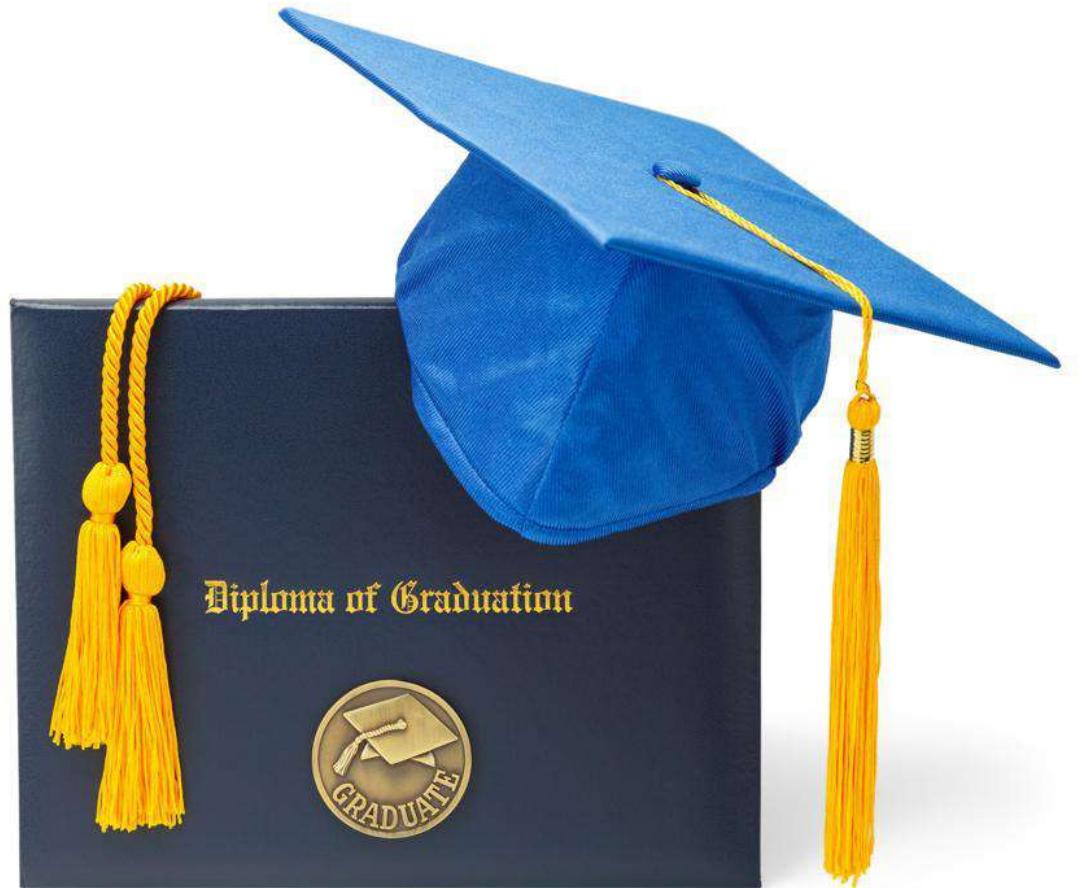
Human-GenAI communications



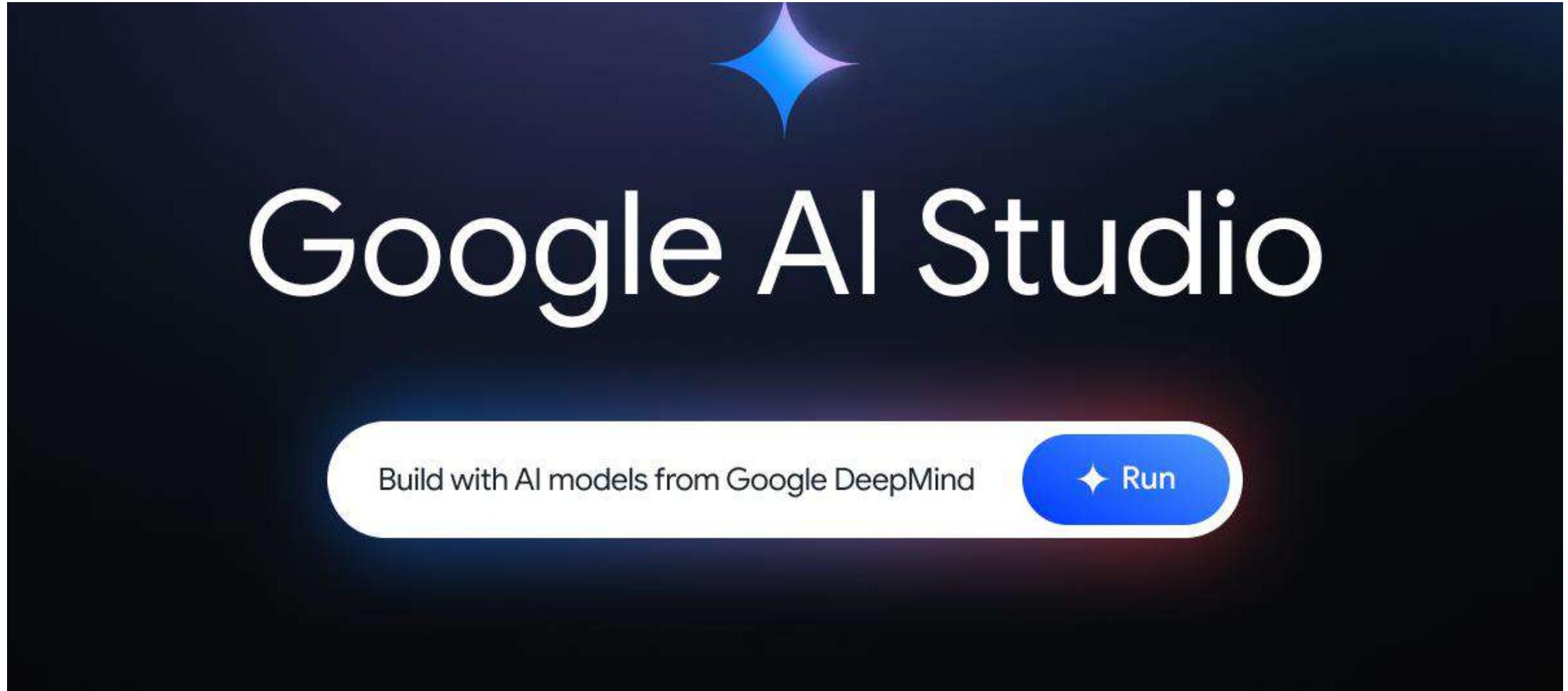
Communication is the key in the future

LLMs aren't mind readers.
They only know what you tell them.

→ Quality input = quality output.

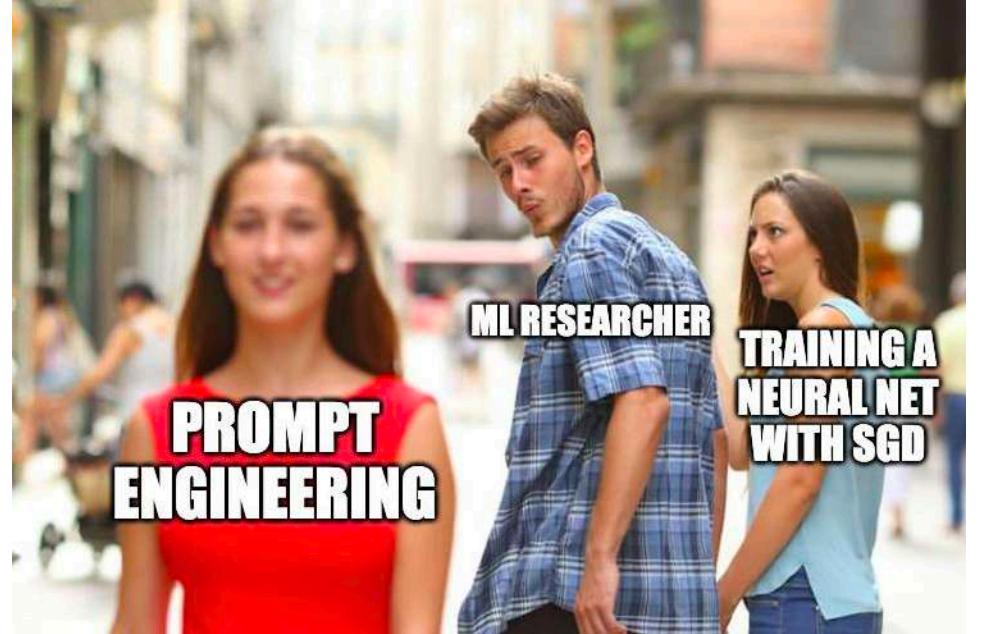


Tool: AI Studio



Concept: Prompt engineering

- Prompt engineering is the practice of designing inputs for generative AI tools that will produce optimal outputs (McKinsey, 2024)
- Say something that AI can easily understand and follow!



[persona]

You are a senior product marketing manager at Apple and you have just unveiled the latest Apple product in collaboration with Tesla, the Apple Car, and received 12,000 pre-orders, which is 200% higher than target

[context]

Write an email to your boss, Tim Cook, sharing this positive news

[task]

[format]

The email should include a tl; dr (too long, didn't read) section, project background (why this product came into existence), business results section (quantifiable business metrics), and end with a section thanking the product and engineering teams.

[exemplar]

[tone]

Use clear and concise language and write in a confident yet friendly tone

Technique	Description
Zero-shot Prompting	Directly asking questions without providing examples, relying on the AI's pre-trained knowledge.
Few-shot Prompting	Providing a few examples before the main question to guide the AI's response format and style.
Chain of Thought (CoT)	Encouraging AI to show step-by-step reasoning for problem-solving, improving transparency and accuracy.
Step-by-step Prompting	Breaking complex tasks into a series of simple steps, guiding AI through each stage.
Tree of Thought (ToT)	Exploring multiple reasoning paths, similar to a decision tree, for more comprehensive solutions.
Self-consistency	Generating multiple independent solutions and choosing the most common or reasonable answer.
Structural prompts	Using standardized prompt structures with placeholders for specific content, ensuring consistency.
Reverse Prompting	Asking AI to generate prompts that would lead to a specific output, exploring AI's associative logic.
AI Interview Technique	Simulating an interview process where AI asks a series of questions to gather detailed information.
Thought Provocation	Using open-ended questions or hypothetical scenarios to stimulate creative thinking.
Meta-prompting	Using prompts to generate or improve other prompts, optimizing AI interaction strategies.

Structuring prompts (JSON prompting)

- A well-structured prompt makes it easier for the AI to follow along, especially for complex requests.
- **Use labeled sections or prefixes:** For simple prompts, you can use labels followed by a colon. For example:
 - **Background:** *[Describe context]*
 - **Objective:** *[State the goal]*
 - **Instructions:** *[Tell the AI what to do]*
 - **Data:** *[Provide any input data]*

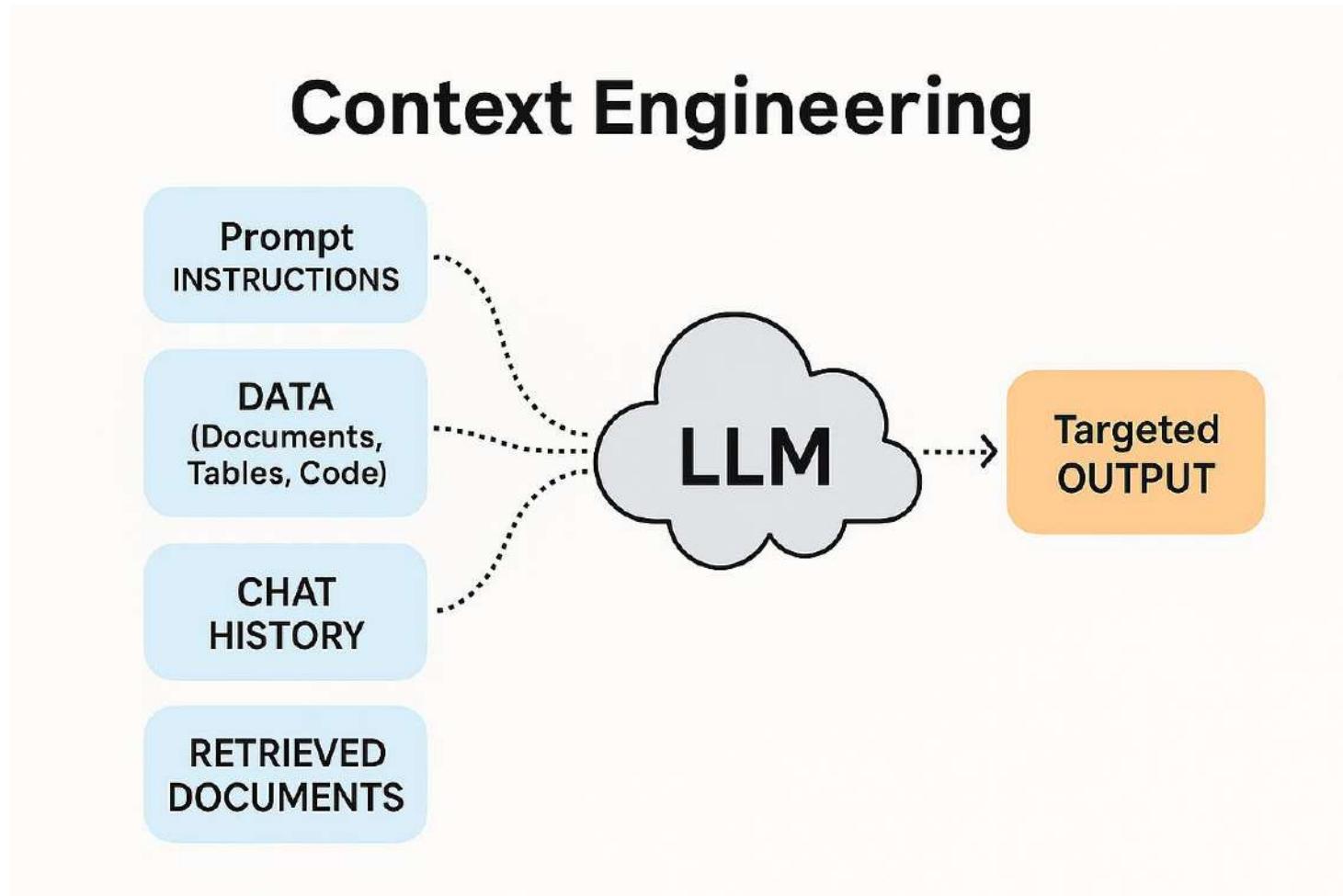
Structuring prompts

- **Use delimiters or XML-like tags for complex prompts:**

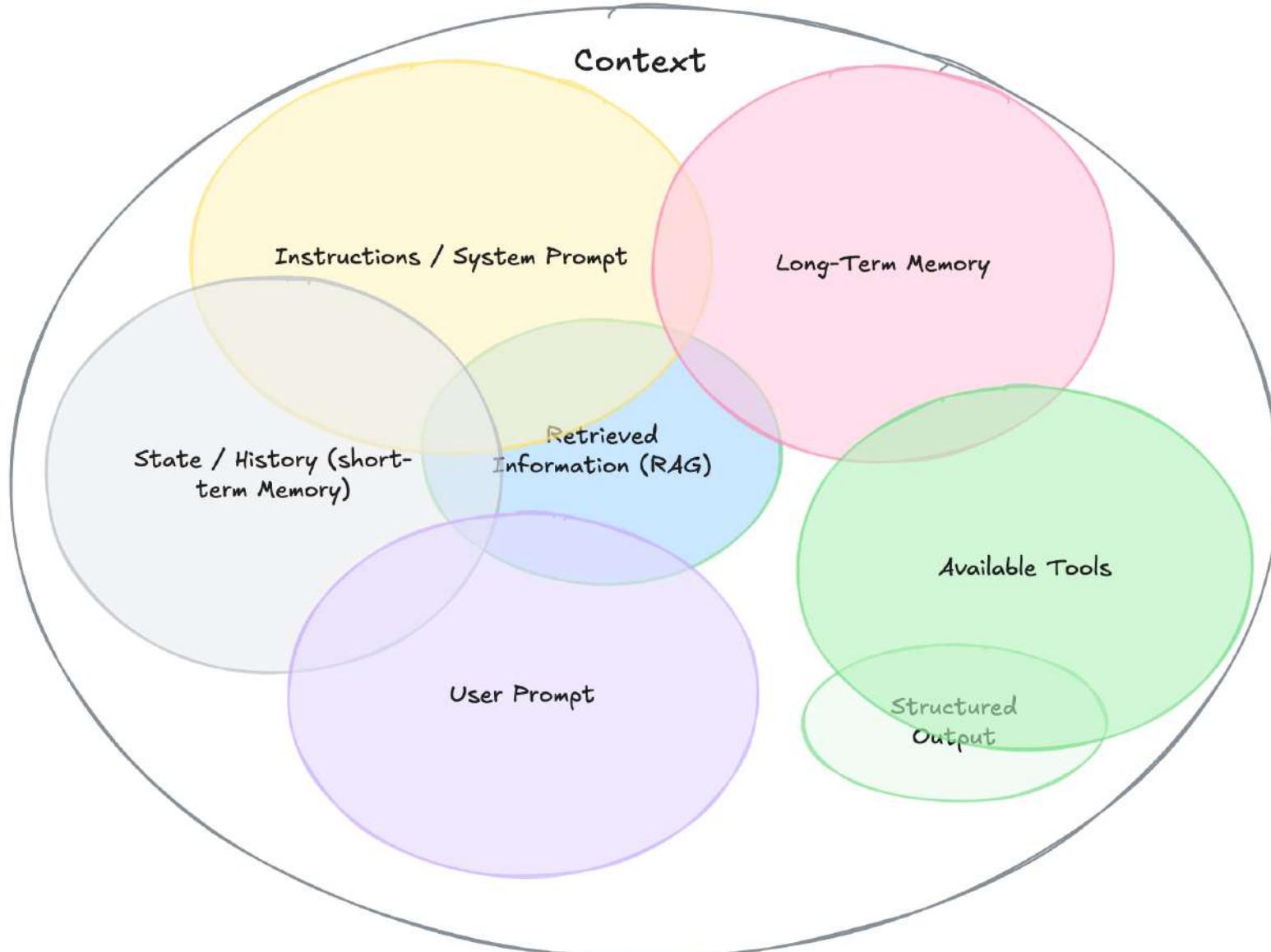
- For example:

```
<instructions> ... </instructions>
<background> ... </background>
<data> ... </data>
<output_format> ... </output_format>
```

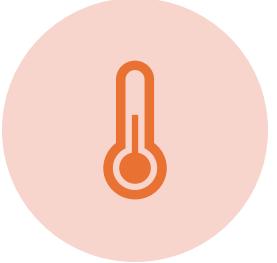
Concept: Context engineering



Context Engineering



Concept: Model settings



TEMPERATURE



TOP-K



TOP-P



TOKENS

Model setting

Temperature

- Temperature controls randomness in predictions. Lower values (e.g., 0.2) make outputs more deterministic, while higher values (e.g., 1.0) make outputs more creative and varied.

Top-k Sampling

- Top-k limits the model to selecting the next token from the top k most probable options. For example, $k=50$ means only the 50 most likely tokens are considered, reducing randomness while keeping diversity.

Top-p Sampling

- Top-p controls how adventurous the word choice is. Imagine all possible next words sorted from most likely to least. Starting at the top, you keep adding words to the “allowed” list until their combined likelihood reaches p . With $p = 0.9$, the model only picks from the smallest set of likely words whose probabilities add up to 90%, ignoring the long tail of unlikely options. Lower p sticks to very safe, predictable words; higher p lets in more variety while still favoring likely choices.

Tokens

- Tokens are the basic units of text that the model processes. They can be as short as a single character or as long as a word (e.g., "apple" is one token, while "applesauce" might be split into multiple tokens).

Tool: Poe



- Most advanced models
- Personalized apps
- Cost-efficiency

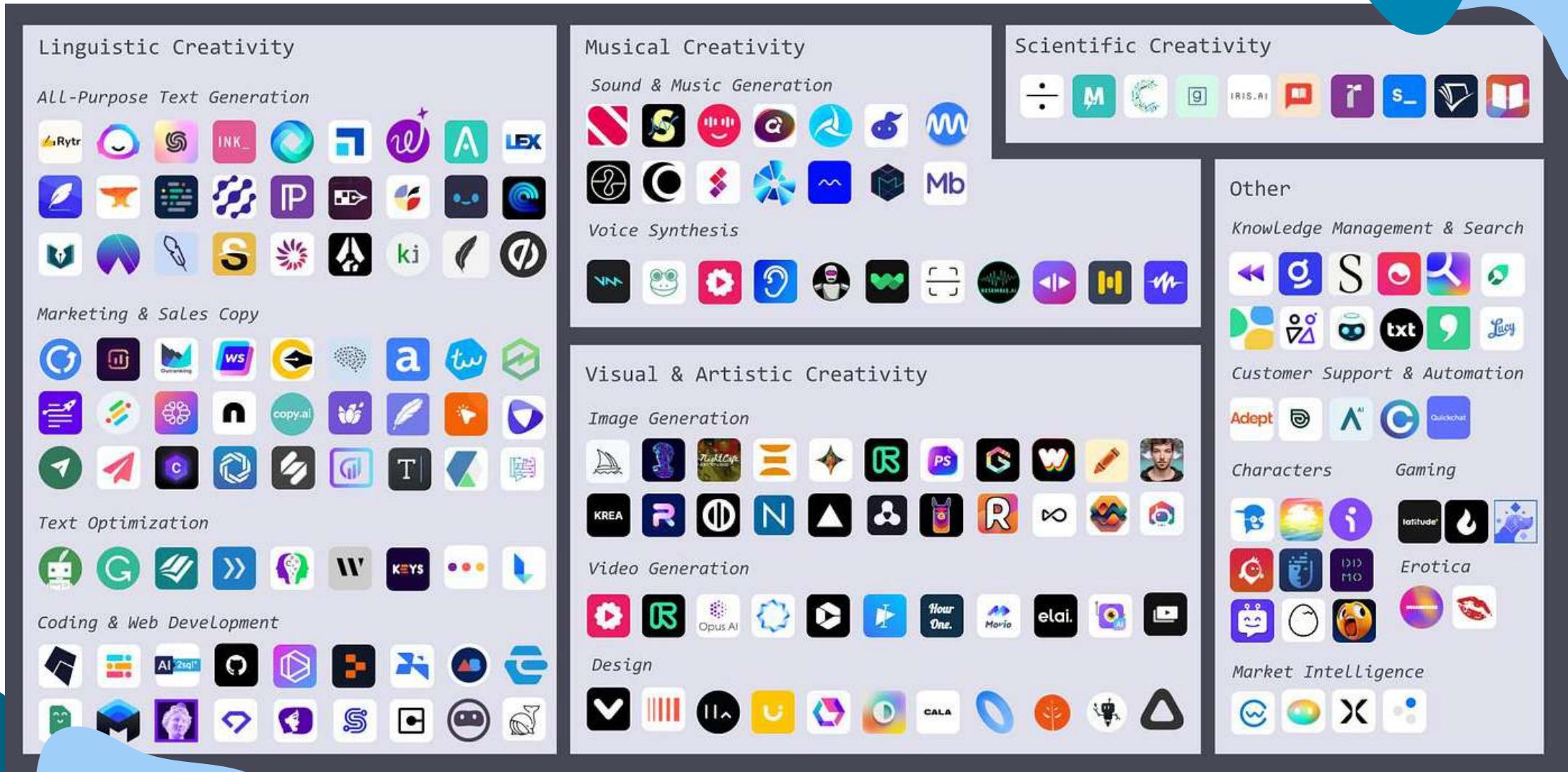
Hands-on GenAI in Action 2025

Session 3: Mimi's Toolkit

Shubin Yu



Toolbox



GenAI applications for content generation

Topic	Examples
Text Generation	Claude, gemini, gpt
Image Generation	Nano Banana, Imagen, midjourney
Image editing	Nano Banana, flux, Seedream
Video Generation	Veo, Runway, Dreamina, Kling, Hailuo
GenAI-Driven Avatar Creation	HeyGen, Synthesia
Generative Audio and Music Production	SUNO, AIVA, Udio
Game Generation	Rosebud AI, layer.ai
Code Generation	Cursor, Windsurf, CodeX, Replit
Virtual World Generation	Convai, Skybox

Tool: image generation



Core prompt structure for image generation

- **Subject:** Specific main focus (age, appearance, clothing, expression)
- **Environment:** Setting, location, time of day, atmosphere
- **Composition:** Perspective, framing, camera angle, depth of field
- **Style & Aesthetic:** Artistic direction, mood, realism level
- **Technical Details:** Camera specs, lighting, resolution

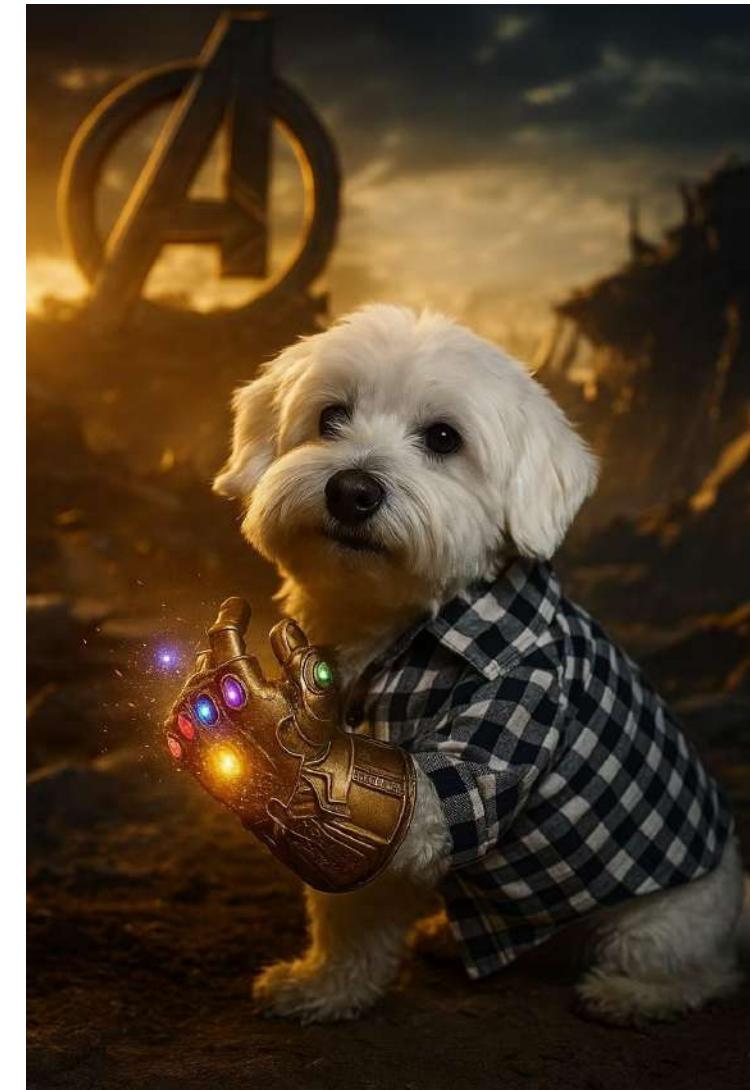
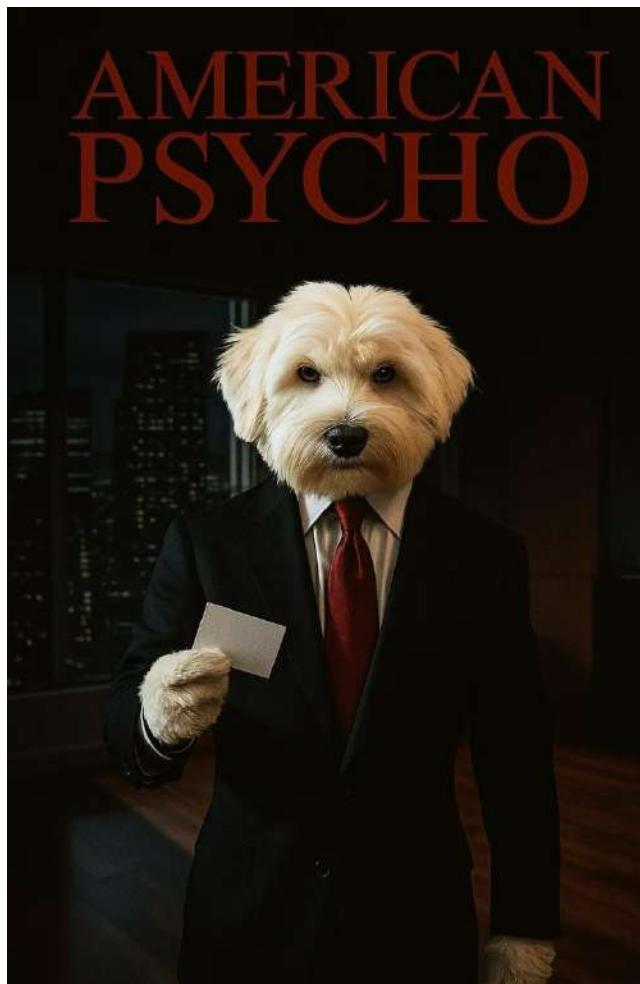
- **Random words:** add some random adjectives and adverbs to make it more creative and realistic.





Key techniques for image generation

- **Be specific but not overwhelming** - Focus on meaningful details
- **Use realism keywords:** "photorealistic, 8K UHD, cinematic lighting, HDR"
- **Reference real photographers/styles:** "Annie Leibovitz style", "film noir aesthetic"
- **Control composition:** "rule of thirds", "shallow depth of field", "low-angle shot"
- **Iterate systematically** - Change one element at a time
- **Make it not perfect –** Imperfection makes perfect



Exercise

- Generate and edit an image of Mimi in a movie scene
- Mimi's picture can be downloaded on our GitHub page
- Make use of the prompt technique
- Need to be highly realistic
- Upload it to Padlet: <https://padlet.com/binbs/mimi>

Hands-on GenAI in Action 2025

Session 4: Assistant, automation, and agent

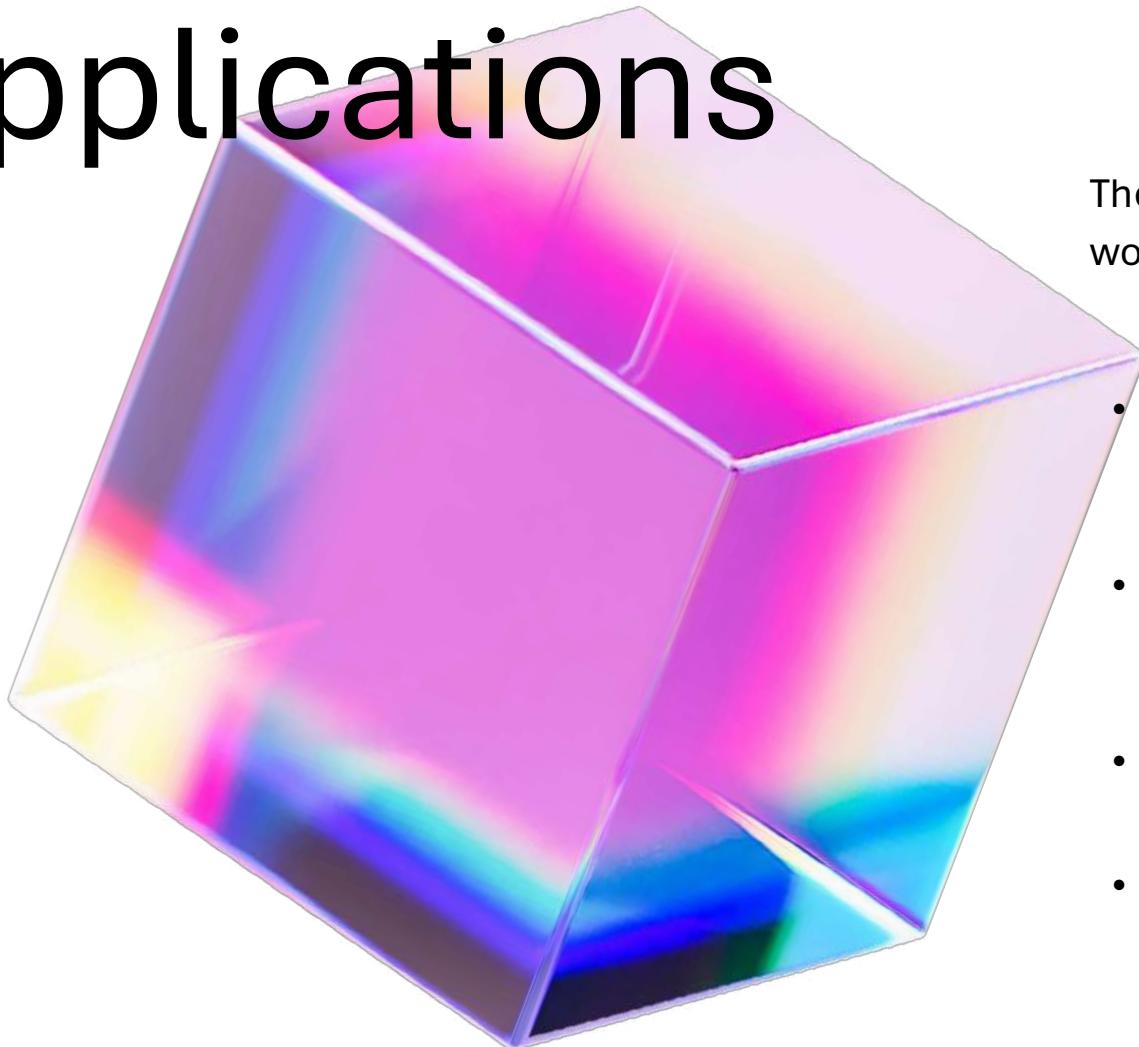
Shubin Yu

Tool: Perplexity



- Search
- Deep search
- Lab

Applications



There are four different types of GenAI applications for individuals at works.

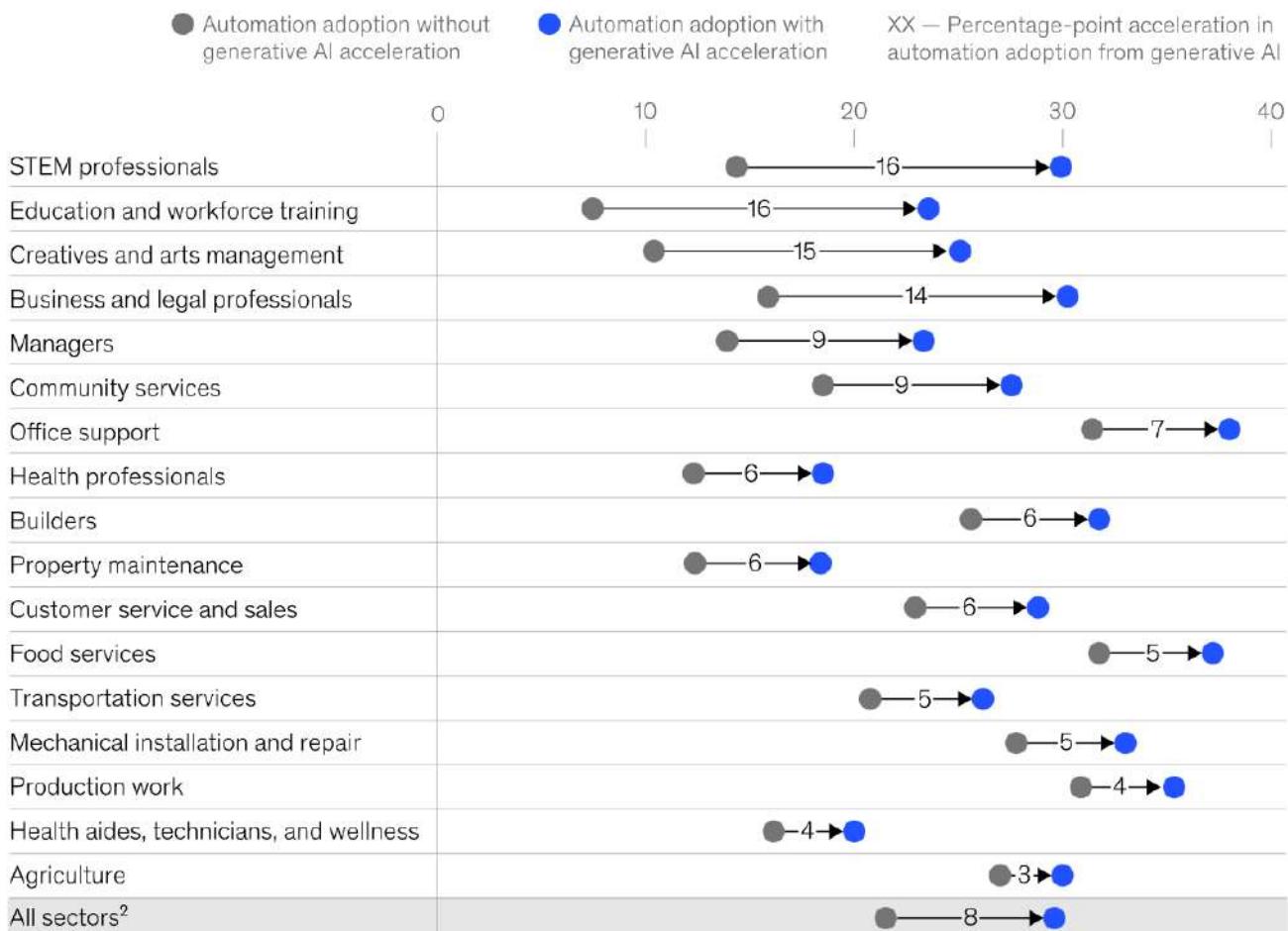
- Foundation models/ Non-Agentic chatbots (e.g., ChatGPT, Gemini, Claude, Deepseek, Grok, Poe)
- Customized assistants (e.g., customized GPT, Coze, Gem, Poe)
- Workflow automation (e.g., n8n, Zapier, Make.com, Dify)
- Agents (e.g., Manus, Replit, Deep Research*, godmode.space)

Tool: Coze



With generative AI added to the picture, 30 percent of hours worked today could be automated by 2030.

Midpoint automation adoption¹ by 2030 as a share of time spent on work activities, US, %



¹Midpoint automation adoption is the average of early and late automation adoption scenarios as referenced in *The economic potential of generative AI: The next productivity frontier*, McKinsey & Company, June 2023.

²Totals are weighted by 2022 employment in each occupation.

Source: O*NET; US Bureau of Labor Statistics; McKinsey Global Institute analysis

Content Creation Workflow (Create an article for my website)

1. Idea Generation:

GenAI Role: Suggests topics based on trending subjects, audience interests, and SEO keywords.

Example: Using tools like Poe to brainstorm blog post ideas using 30 LLMs

2. Draft Writing:

GenAI Role: Creates initial drafts of articles, reports, or social media posts.

Example: Webscrapping and automatically collecting 100 relevant articles online, generating a 1,500-word summary article on “How to use LLMs responsibly” (always cite the original articles)

3. Editing and Proofreading:

GenAI Role: Reviews content for grammar, style, and coherence, providing suggestions for improvement.

Example: Utilizing DeepL to refine the draft for readability and error-free content, translating into 100 languages

4. Generating images and podcast:

GenAI Role: Generating images and podcast based on the text and also the prompt

Example: Using Luma to create images and LM notebook to create the podcast with two virtual hosts

5. SEO Optimization:

GenAI Role: Integrates relevant keywords and optimizes meta descriptions to improve search engine ranking.

Example: Automatically embedding target keywords and crafting SEO-friendly titles and headers.

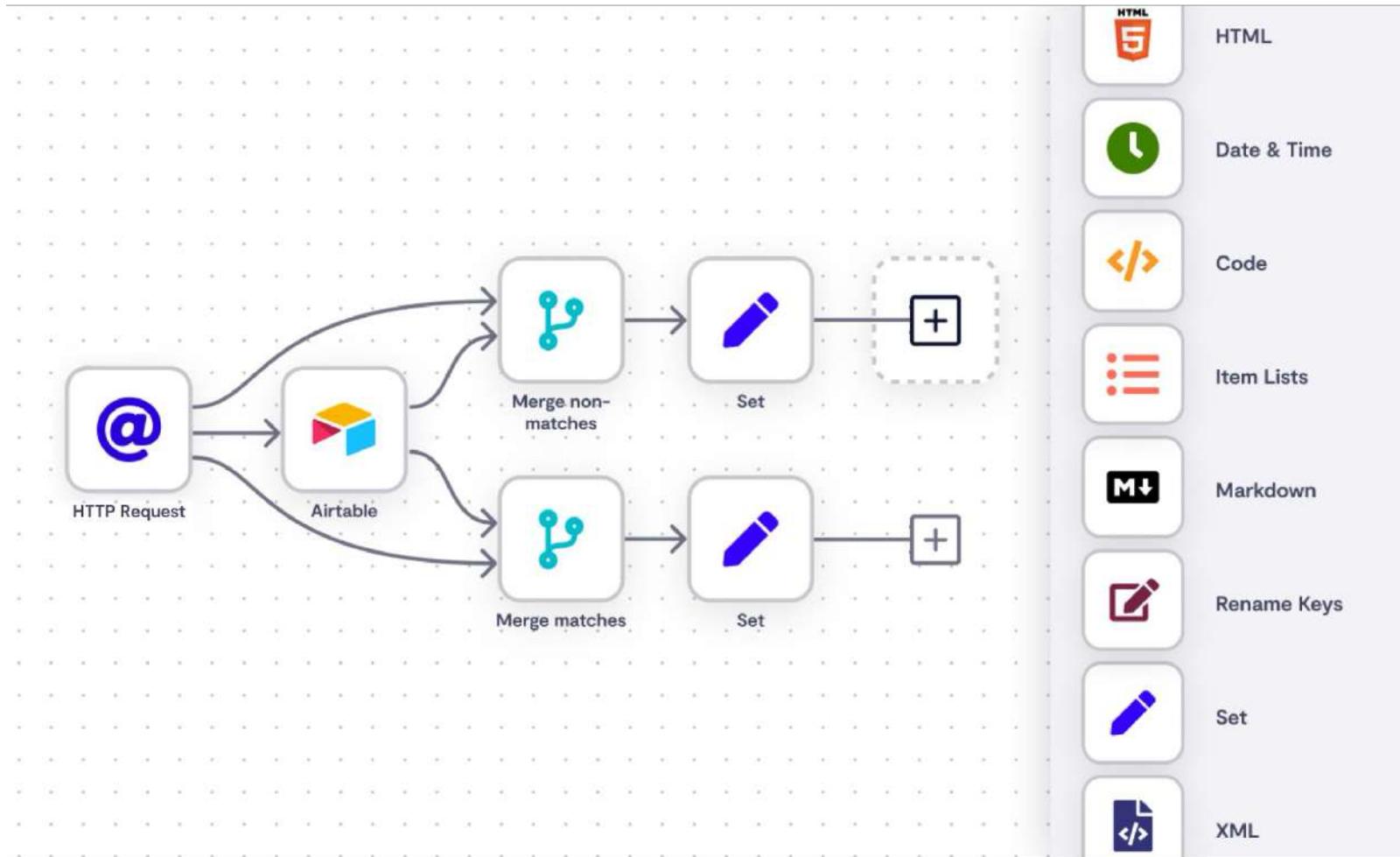
Workflow automation using GenAI

- Content Creation and Marketing
- Customer Service
- Software Development
- Human Resources
- Legal and Compliance
- Financial Services
- Healthcare
- Project Management
- Product Design
- Data Analysis



How can I automate
my current tasks at
work?

Tool: n8n



Agents



MCP (Anthropic)

Model context protocol

An open standard that enables developers to build secure, two-way connections between their data sources and AI-powered tools

A2A (Google)

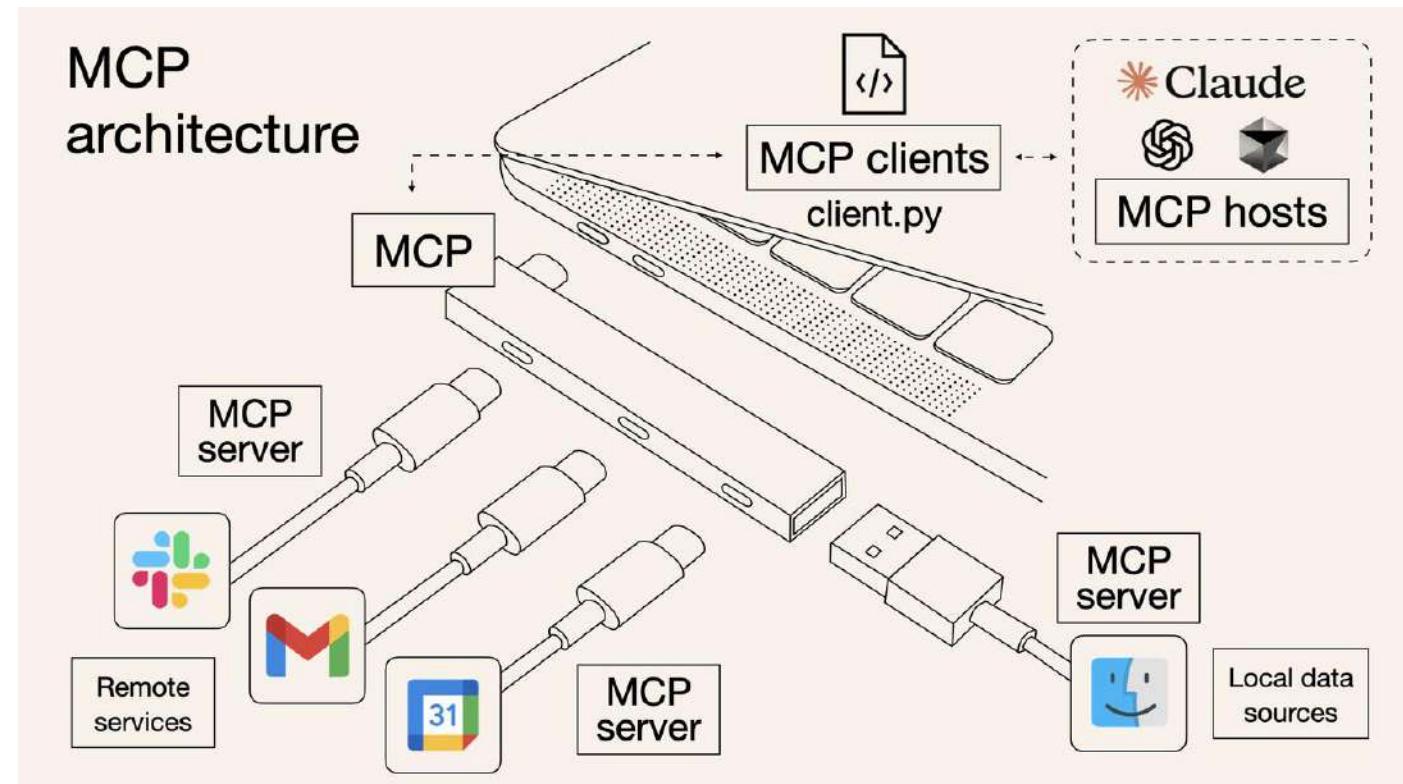
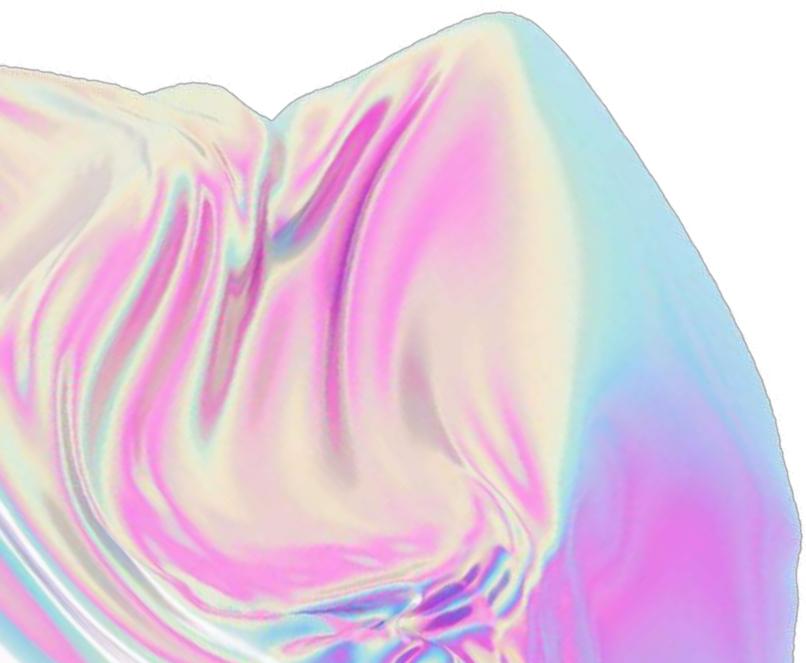
Agent2Agent Protocol

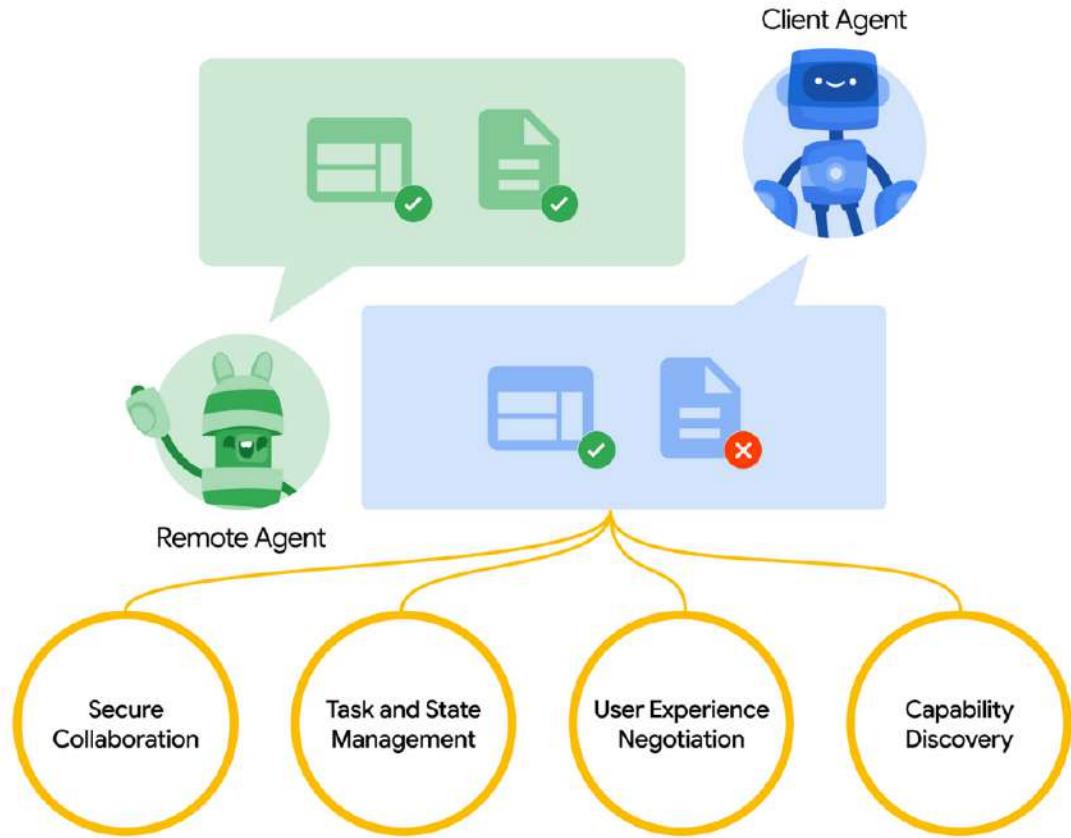
An open protocol that provides a standard way for agents to collaborate with each other, regardless of the underlying framework or vendor.

MCP (Anthropic)

Model context protocol

An open standard that enables developers to build secure, two-way connections between their data sources and AI-powered tools

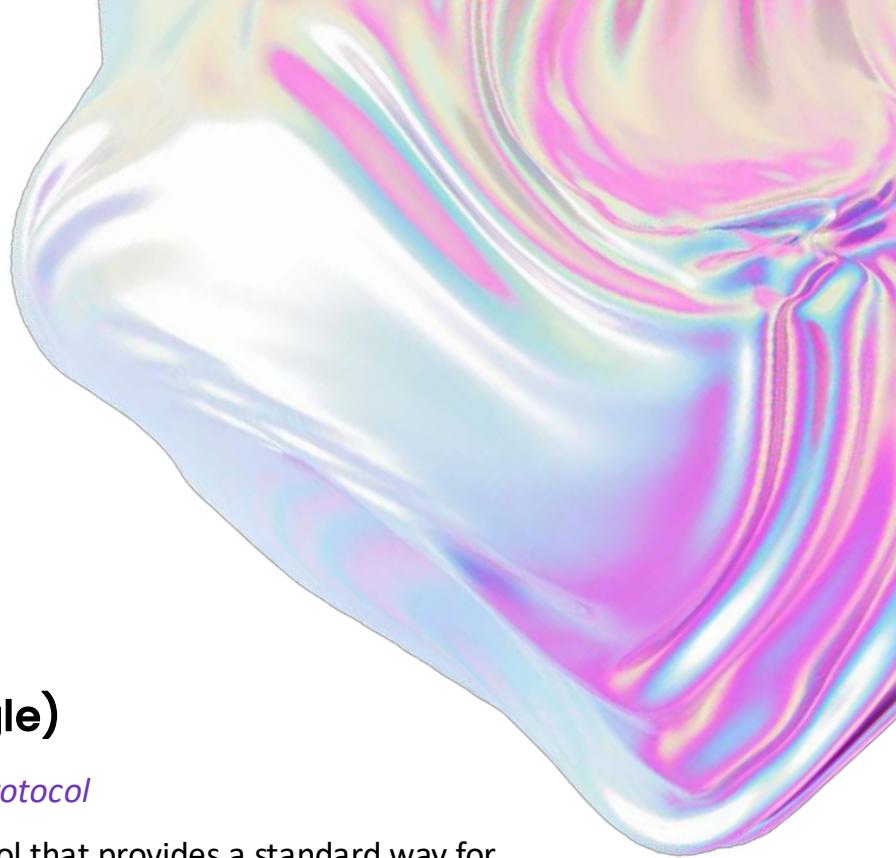


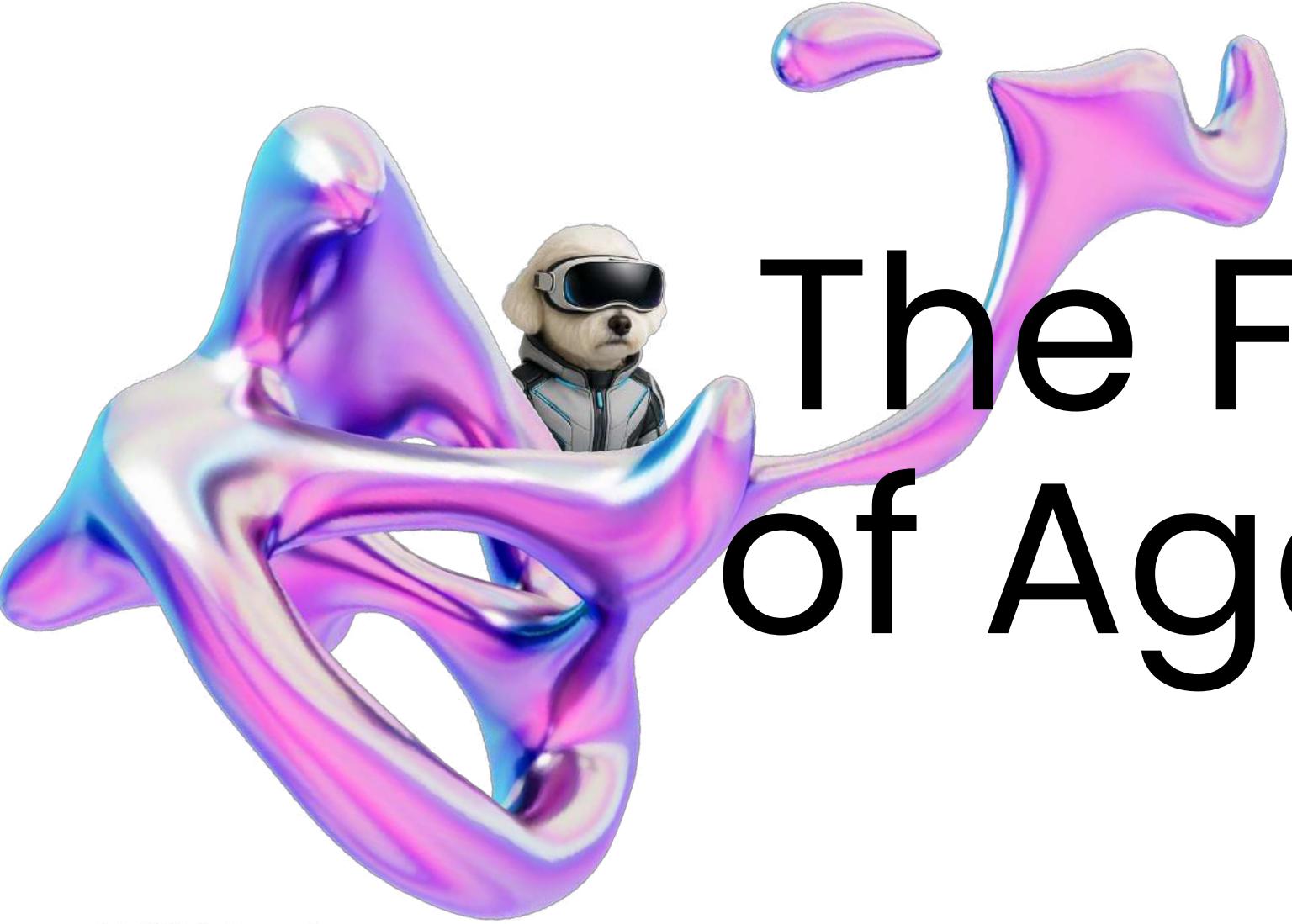


A2A (Google)

Agent2Agent Protocol

An open protocol that provides a standard way for agents to collaborate with each other, regardless of the underlying framework or vendor.





The Future of Agents

Tool: Manus



Hands-on GenAI in Action 2025

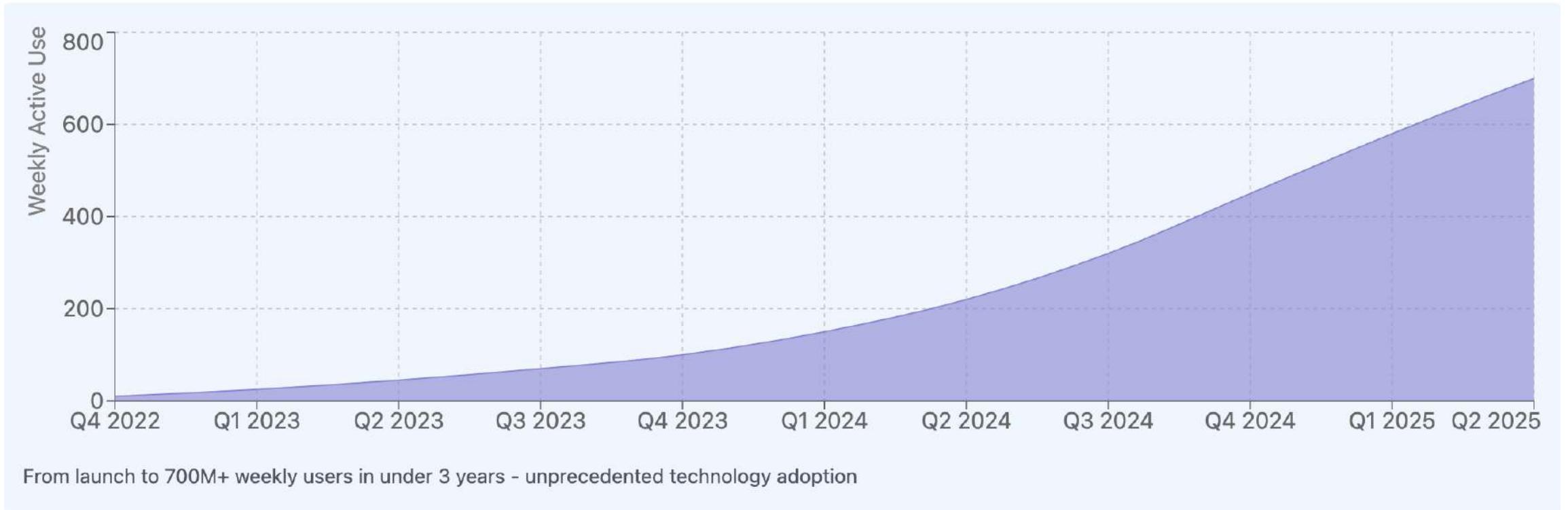
Session 5: Strategic integrations

Shubin Yu

How are people using GenAI?

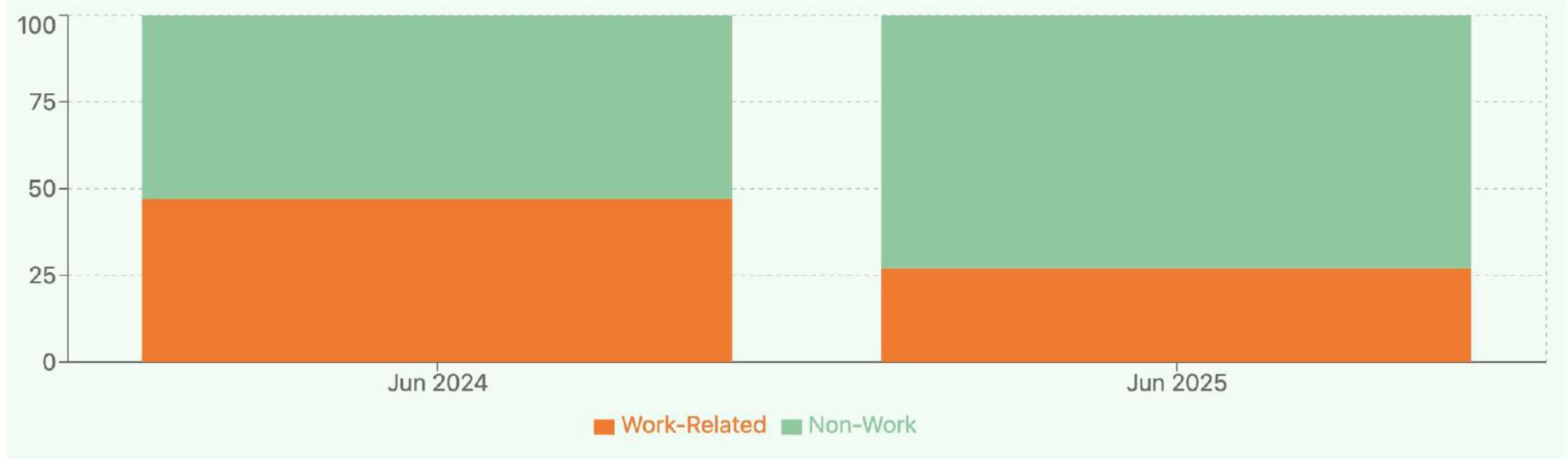
- Reports from OpenAI
- Reports from Anthropic
- Reports from the consulting company

OpenAI



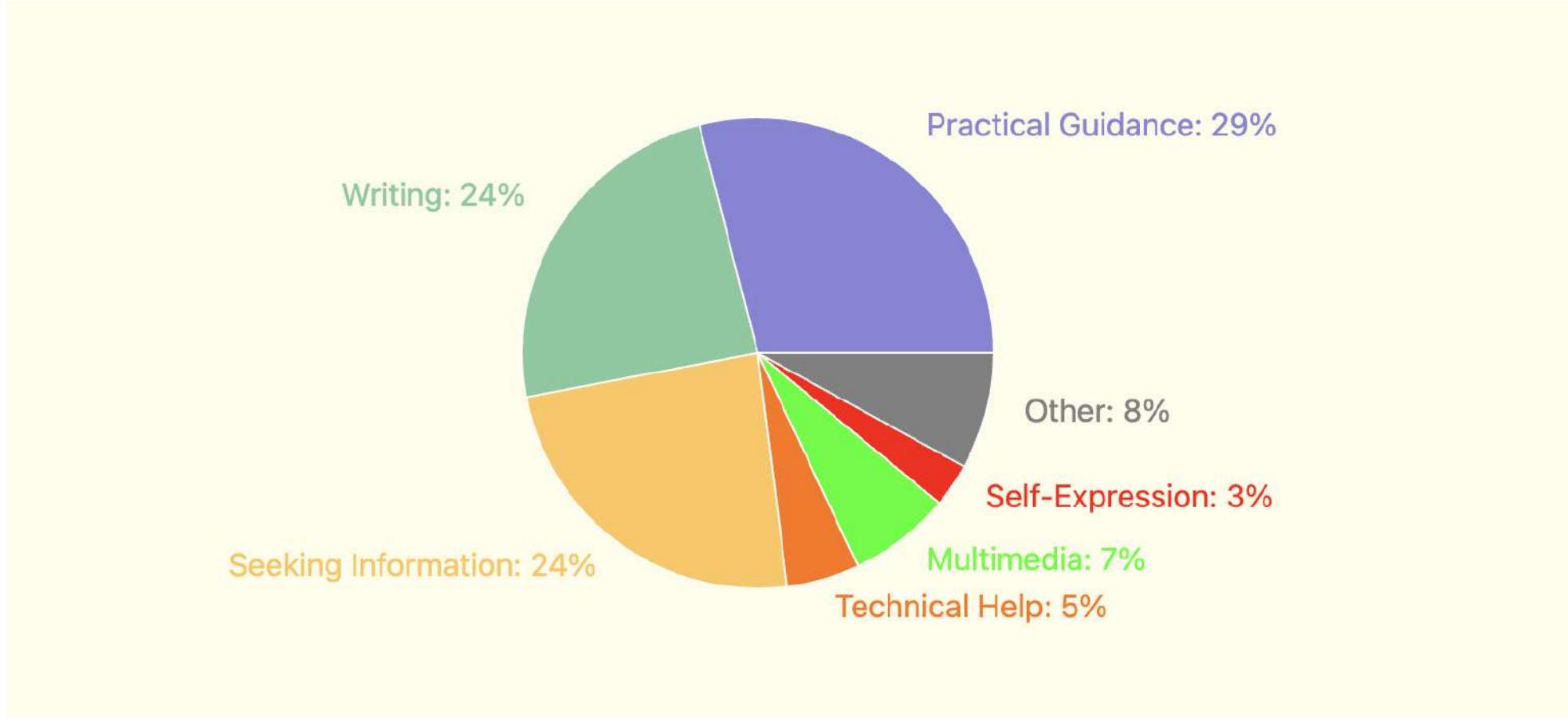
Weekly active ChatGPT users on consumer plans (Free, Plus, Pro)

OpenAI



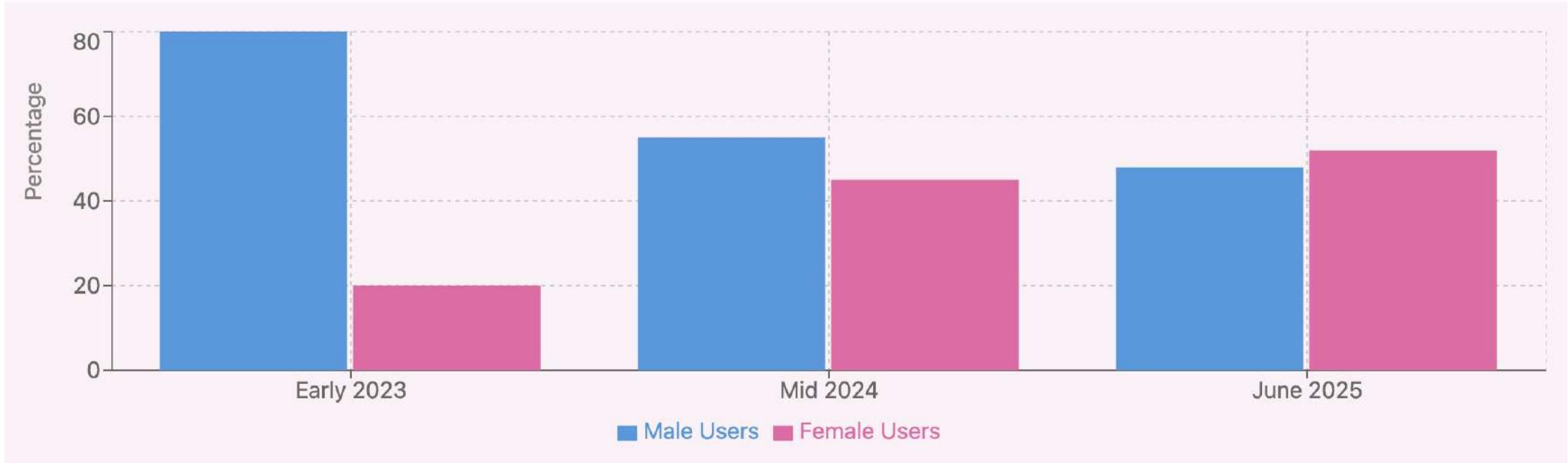
Non-work usage grew from 53% to 73% - ChatGPT becoming more of a personal assistant

OpenAI



Three categories dominate: Practical Guidance, Writing, and Information Seeking (~78% total)

OpenAI

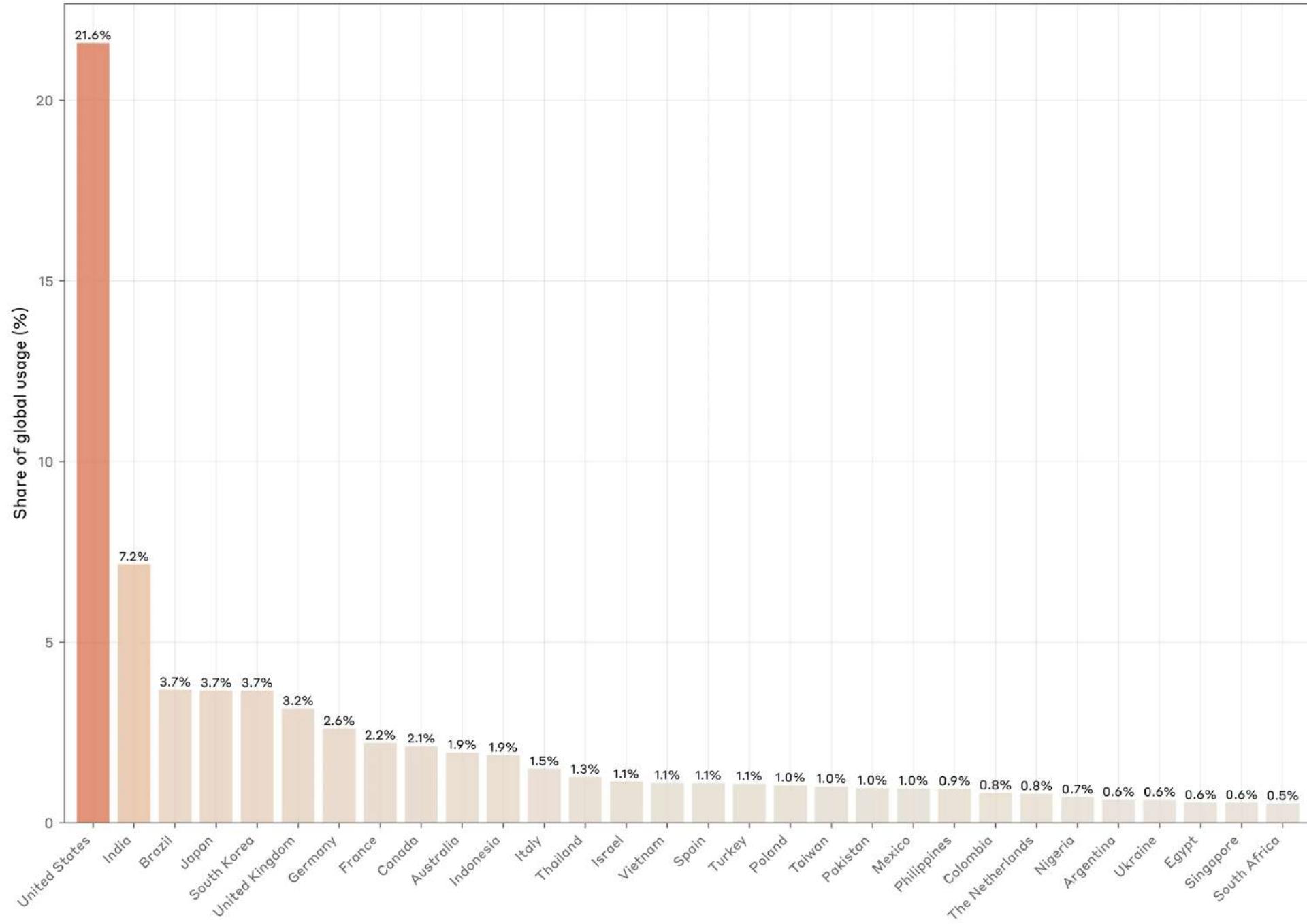


From 80% male early adopters to slight female majority by 2025

OpenAI

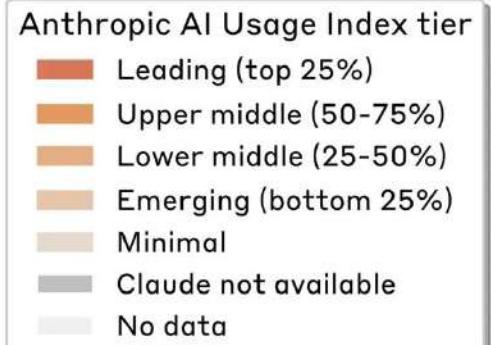
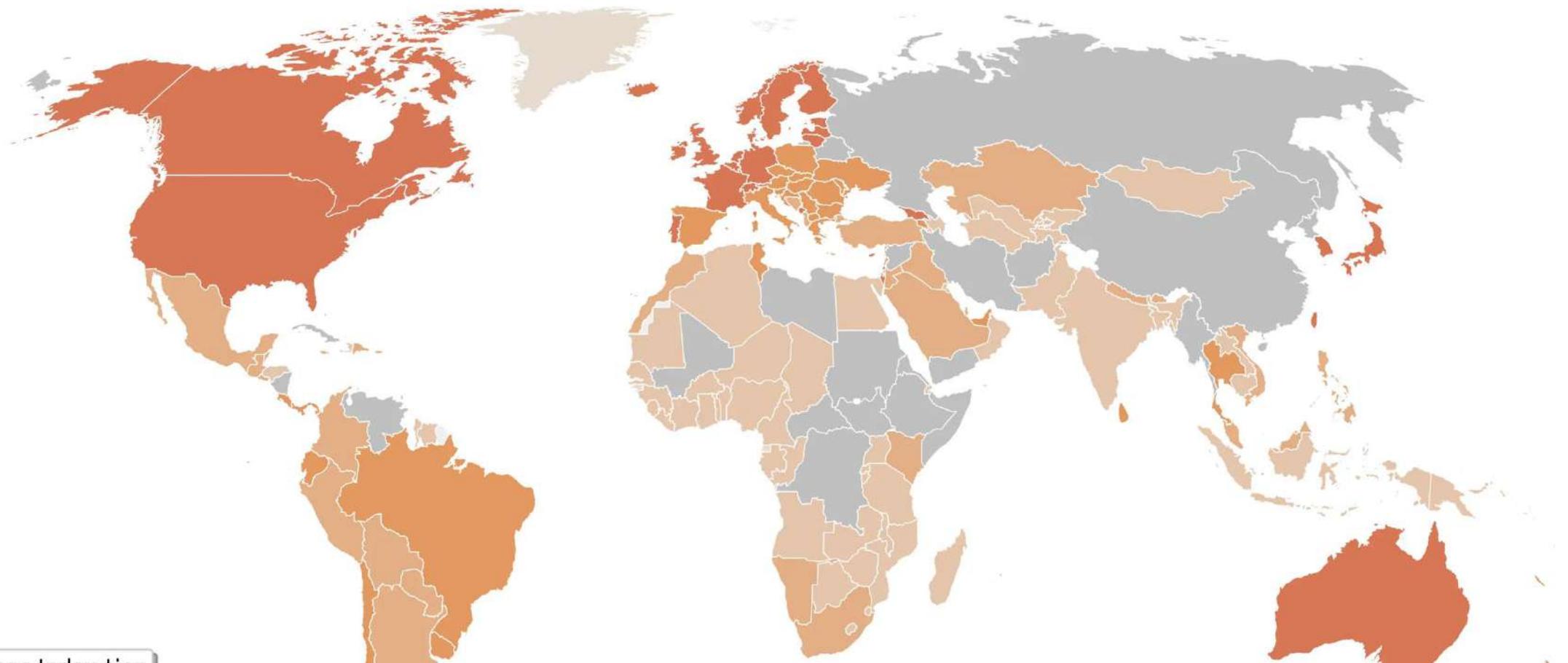
- 700 million users
- Approximately 40% of employees in the U.S. have used ChatGPT.
- 18- to 25-year-olds make up 46% of the users.
- Only 27% is used for work, of which 40% is for writing.
- 52% of users are female.

Top 30 countries by share of global Claude usage

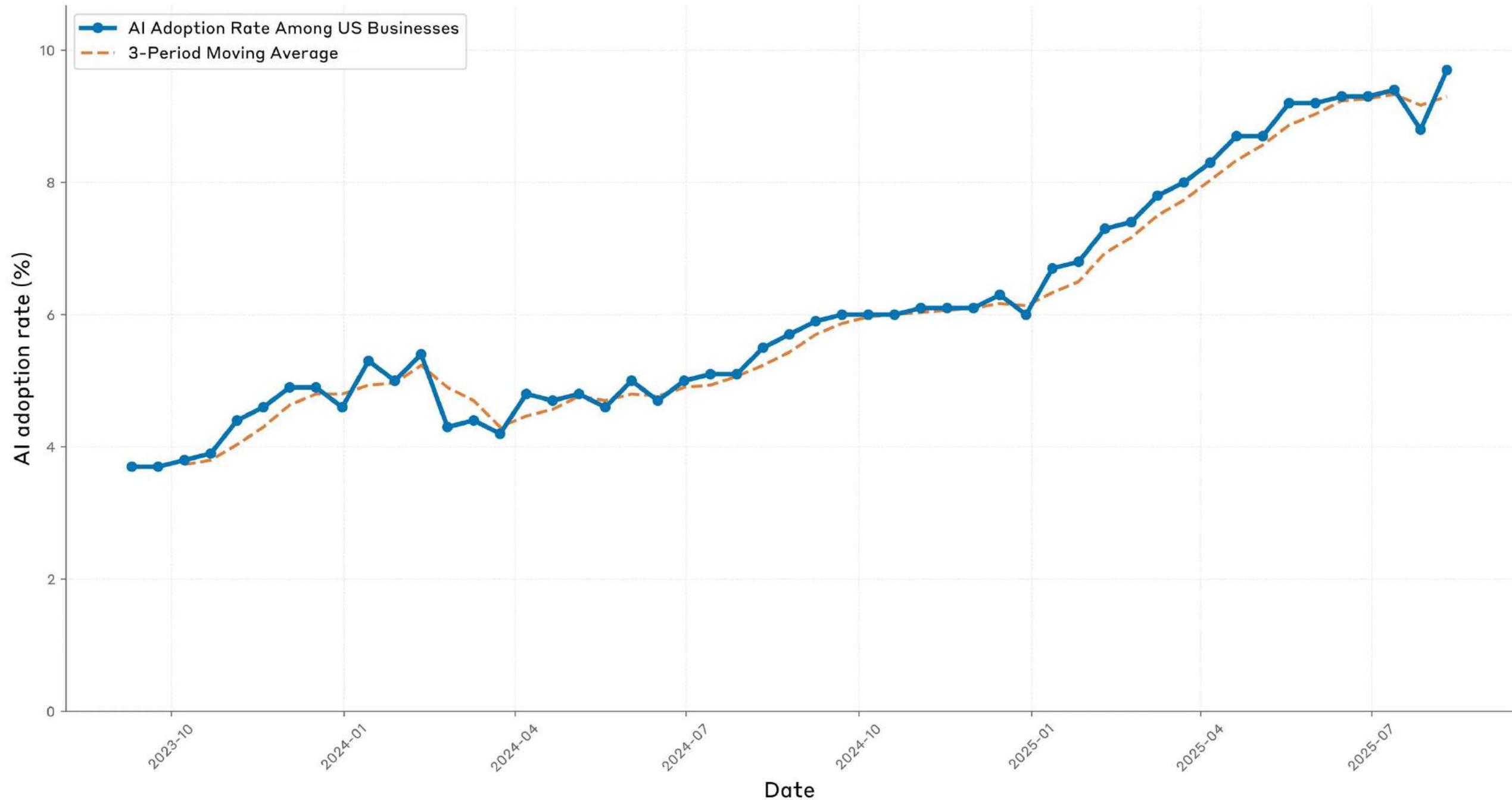


Claude adoption overall is highly geographically concentrated.

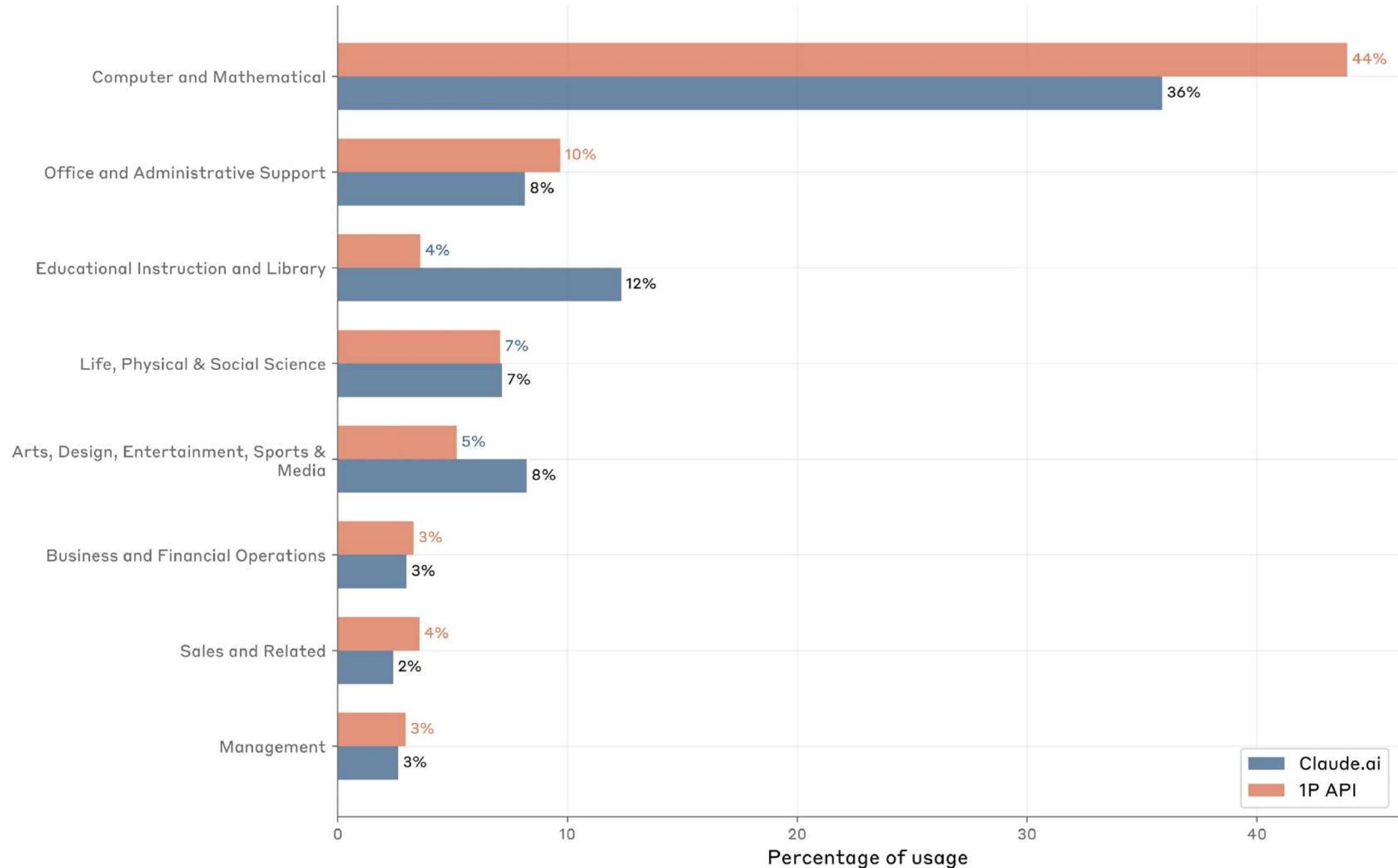
Anthropic AI Usage Index tiers by country



Census reported AI adoption rates among US businesses from the Business Trends and Outlook Survey



Usage shares across top occupational categories: Claude.ai vs 1P API



Claude

- Developed countries used more. This may result in inequality across various regions.
- 77% of enterprise usage is dedicated to automation.
- Enterprise usage prefers more advanced models.

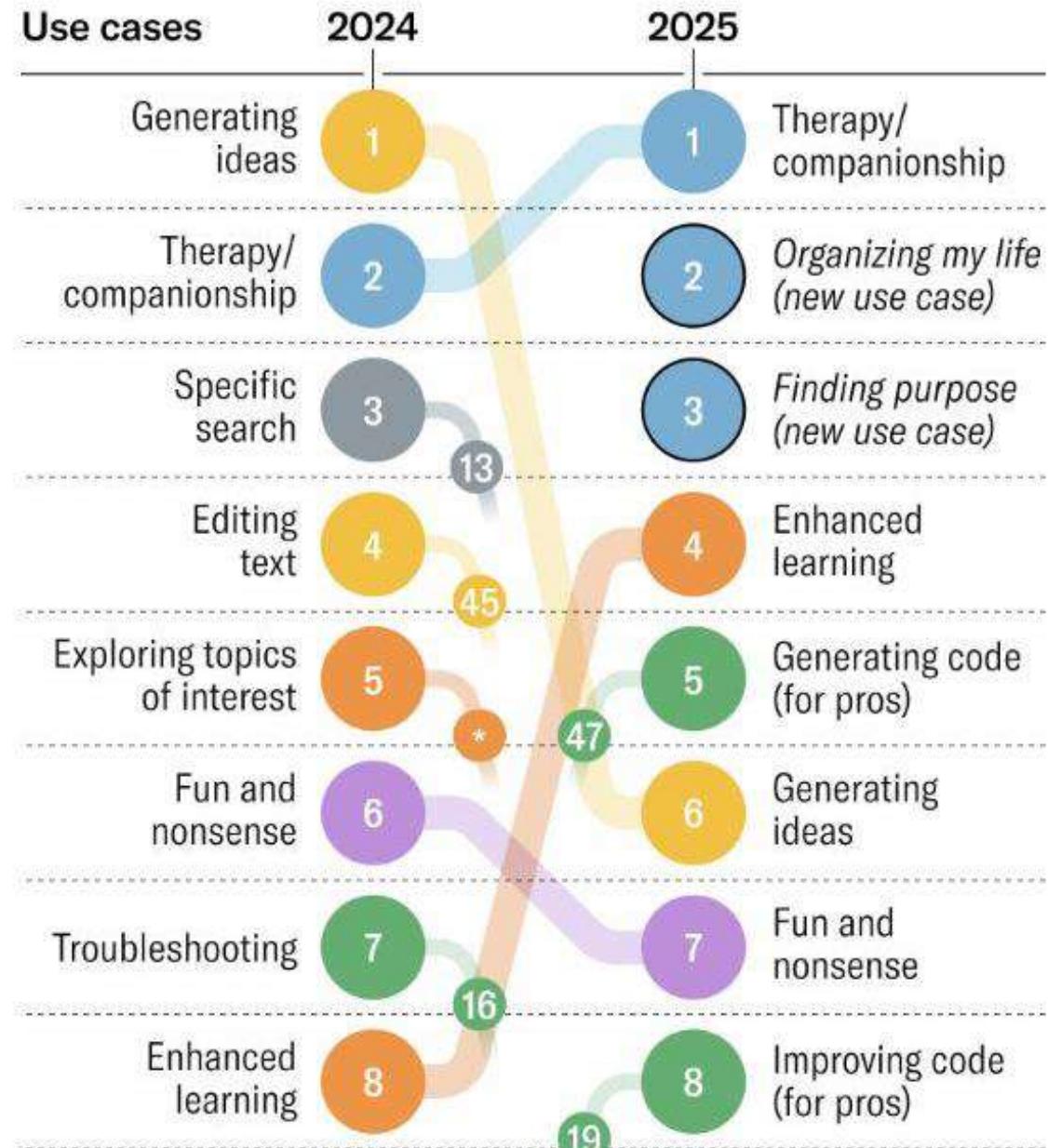
Top 10 Gen AI Use Cases

The top 10 gen AI use cases in 2025 indicate a shift from technical to emotional applications, and in particular, growth in areas such as therapy, personal productivity, and personal development.

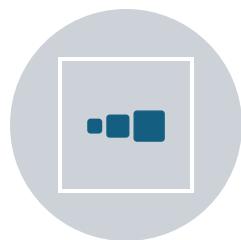
Themes

- PERSONAL AND PROFESSIONAL SUPPORT
- CONTENT CREATION AND EDITING
- LEARNING AND EDUCATION

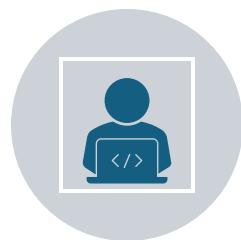
- TECHNICAL ASSISTANCE AND TROUBLESHOOTING
- CREATIVITY AND RECREATION
- RESEARCH, ANALYSIS, AND DECISION-MAKING



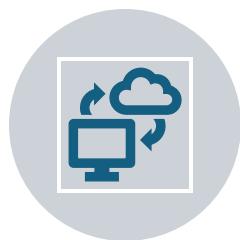
Four levels of integration



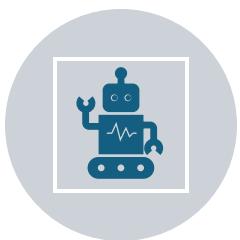
ADOPT PUBLICLY
AVAILABLE TOOLS



CUSTOMIZE THE
TOOL, API



CREATE AUTOMATIC
AND CONTINUOUS
DATA FEEDBACK
LOOPS



DEVELOP OWN
MODELS

Source: Scott Cook, Andrei Hagiu, and Julian Wright (2024), Turn Generative AI from an Existential Threat into a Competitive Advantage, *Harvard Business Review*.
IBM online course: <https://www.coursera.org/learn/generative-ai-for-executives-business-leaders>

Level 1: Adopt publicly available tools

- Use off-the-shelf generative AI tools
- Improve internal processes and efficiency
- No customization of AI models
- Temporary advantage, soon becomes table stakes
- Privacy issues

Source: Scott Cook, Andrei Hagiu, and Julian Wright (2024), Turn Generative AI from an Existential Threat into a Competitive Advantage, *Harvard Business Review*.
IBM online course: <https://www.coursera.org/learn/generative-ai-for-executives-business-leaders>

Level 2: Customize the tools

- Create customized AI tools using company data and know-how
- Enhance customer experience and add new capabilities
- Some customization through fine-tuning with company data
- Potential for personalization and improved user interfaces

Source: Scott Cook, Andrei Hagiu, and Julian Wright (2024), Turn Generative AI from an Existential Threat into a Competitive Advantage, *Harvard Business Review*.
IBM online course: <https://www.coursera.org/learn/generative-ai-for-executives-business-leaders>

Level 3: Create automatic and continuous data feedback loops

- AI tools produce reliable signals from customer usage
- Feedback automatically improves the model with minimal human intervention
- Creates a compounding competitive advantage
- Requires redesigning products/services to integrate AI throughout

Source: Scott Cook, Andrei Hagiu, and Julian Wright (2024), Turn Generative AI from an Existential Threat into a Competitive Advantage, *Harvard Business Review*.
IBM online course: <https://www.coursera.org/learn/generative-ai-for-executives-business-leaders>

Level 4: Develop own models

- Highly costly
- Flexible, tailored to specific problems
- Difficult to be copied, competitive advantages
- Requires strong tech resources

Source: Scott Cook, Andrei Hagiu, and Julian Wright (2024), Turn Generative AI from an Existential Threat into a Competitive Advantage, *Harvard Business Review*.
IBM online course: <https://www.coursera.org/learn/generative-ai-for-executives-business-leaders>



GenAI in Marketing

Generative AI: Powering Data-Driven Marketing

Marketing Communications

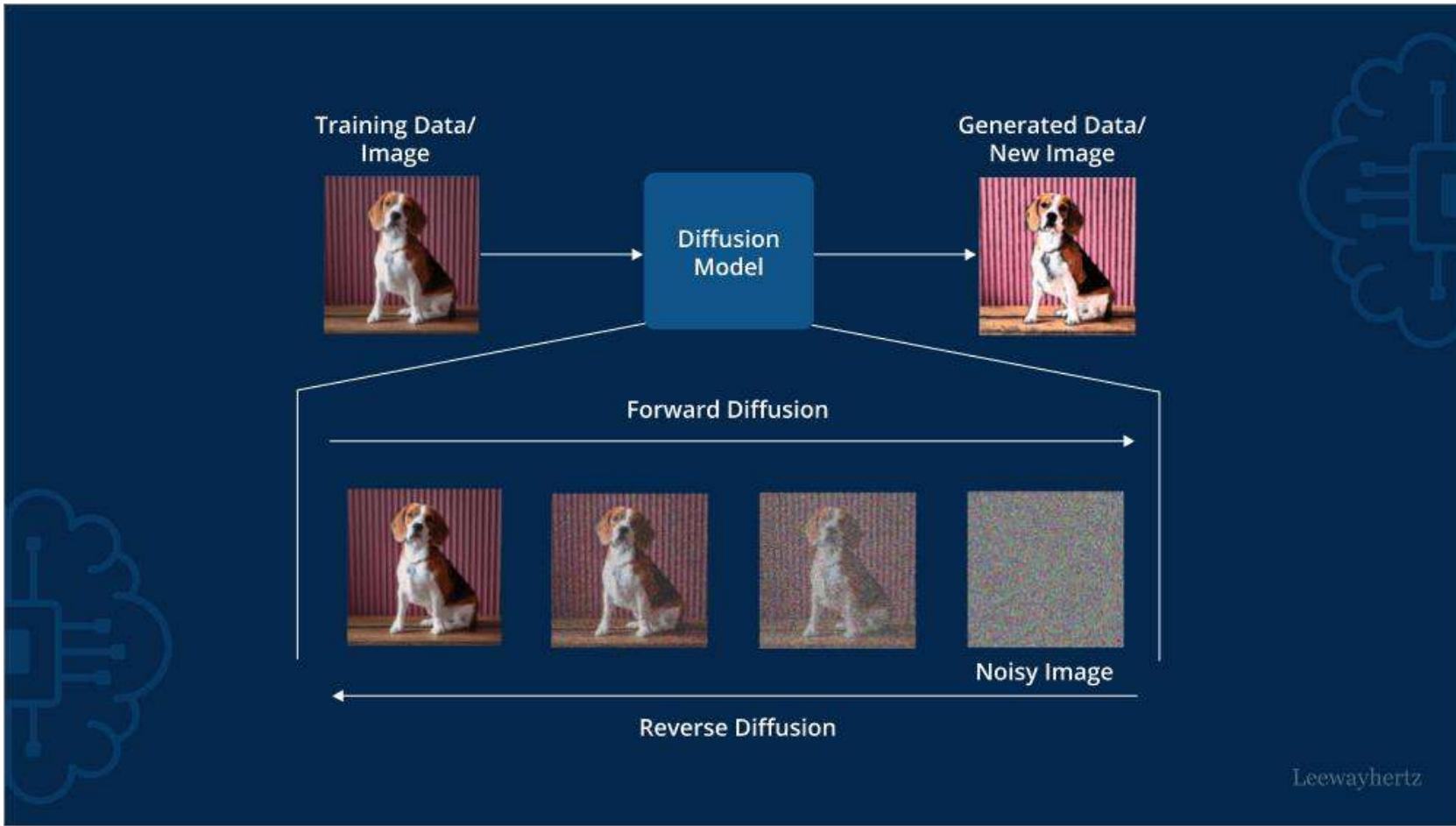
GenAI generates engaging content, personalized campaigns, and segmented audiences.

Customer Service

GenAI chatbots provide instant responses, personalize interactions, and improve customer satisfaction.

Product Development and research

GenAI optimizes designs, prototypes features, and creates innovative products based on customer insights.



Generative AI: Powering Creative Advertising

Efficient Content Creation

AI can automate many aspects of content creation, saving time and resources for marketers. It can generate high-quality ideas, ad copy, design visuals, and even produce audios and videos.

Fine-tune models

Advertisers can fine-tune open-source GenAI models (e.g., stable diffusion) using their own training materials.

Performance

Preliminary research findings show that AI-generated ads have much better performance than freelancers generated ads*.





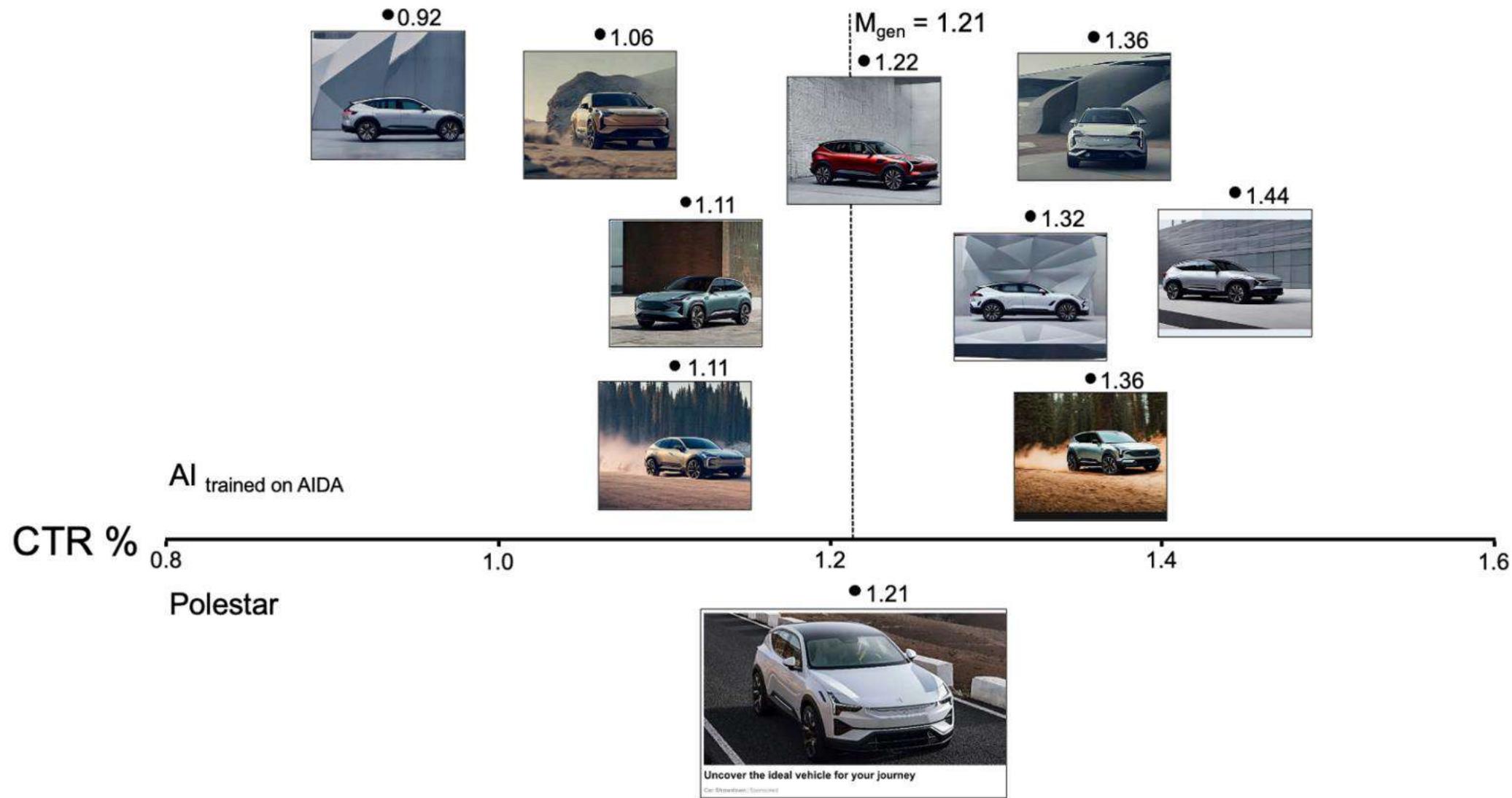
Actual Polestar



AI trained on AIDA



Figure 8: Click-through rates of actual and AI-generated ads

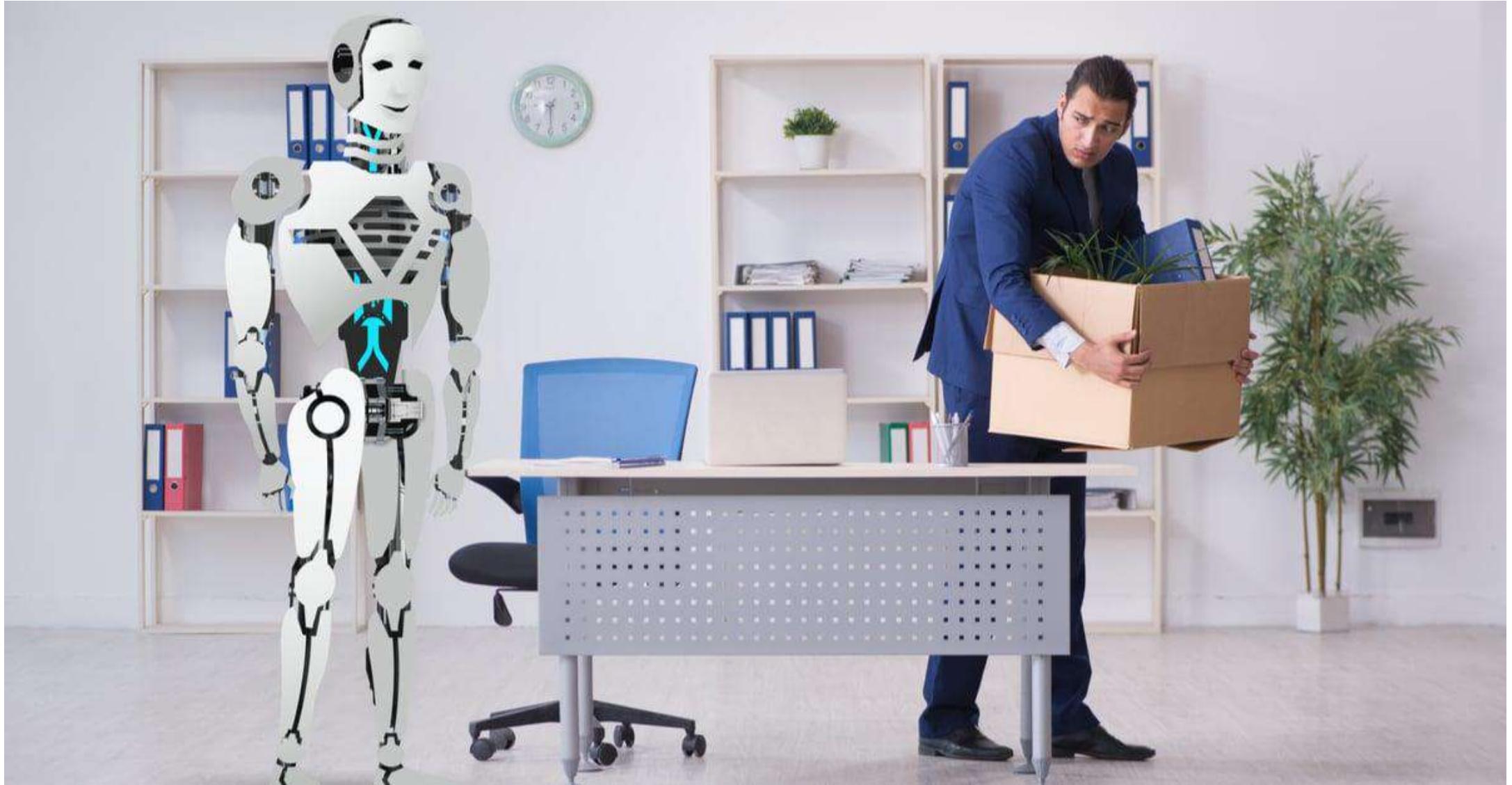


Which one was
made by AI?



It's only \$150 to
create this video
ad.





Tool: video generation



Core prompt structure for video generation

- **Subject:** Main character/object with detailed appearance and emotions
- **Action:** Specific movements and behaviors (vivid verbs)
- **Scene/Environment:** Physical setting, atmosphere, lighting conditions
- **Camera Motion:** Movement type (pan, dolly, tracking, static)
- **Composition:** Shot type (close-up, wide, overhead, medium)
- **Style/Aesthetic:** Cinematic references, color grading, mood
- **Temporal Elements:** Timing, transitions, sequence flow

Key techniques

- **Cinematic Realism Focus** - Include camera specs (ARRI ALEXA, RED EPIC), lens choices (24mm wide, 85mm portrait)
- **Dynamic Camera Movement** - Specify tracking shots, slow pans, dolly movements
- **Atmospheric Details** - Add environmental effects (fog, rain, sunbeams, steam)
- **Professional References** - Citation cinematographers/directors for style guidance
- **Temporal Specifications** - Include slow-motion, transitions, sequence timing
- **Color Grading** - Specify mood-appropriate palettes (warm/vibrant vs cool/muted)

Formula template

- **Subject + Action + Scene + Camera Movement + Lighting + Style**
- **Example:** "Professional chef with white hat, chopping vegetables with focused concentration, in a sunlit stainless-steel kitchen, slow pan left to right, warm cinematic lighting, shallow depth of field"

Cinematic Movie Studio Style

Hollywood-level production values with professional camera work

"Epic wide shot of ancient temple, ARRI ALEXA camera, dramatic golden hour lighting, sweeping crane movement, Hans Zimmer-style atmosphere"

POV (Point-of-View) Style

Character perspective showing what they see and experience

"POV shot walking through bustling market, handheld camera movement, vendors calling out, colorful spices and fabrics, immersive audio"

FPV (First-Person View) Style

Drone-like perspective with fluid, dynamic movement

"FPV drone racing through forest canopy, weaving between trees, sunlight filtering through leaves, smooth gimbal movement, adrenaline rush"

Director Style References

Christopher Nolan

Complex narratives, practical effects, IMAX quality

Wes Anderson

Symmetrical composition, pastel colors, whimsical

James Cameron

Epic scale, cutting-edge technology, immersive

Greta Gerwig

Natural lighting, intimate character focus

Cinematographic Styles

Film Noir: "high contrast shadows, dramatic lighting, urban setting"

Golden Hour: "warm sunlight, soft shadows, romantic mood"

Cyberpunk: "neon lights, rain-slicked streets, futuristic cityscape"

Documentary: "handheld camera, natural lighting, authentic moments"