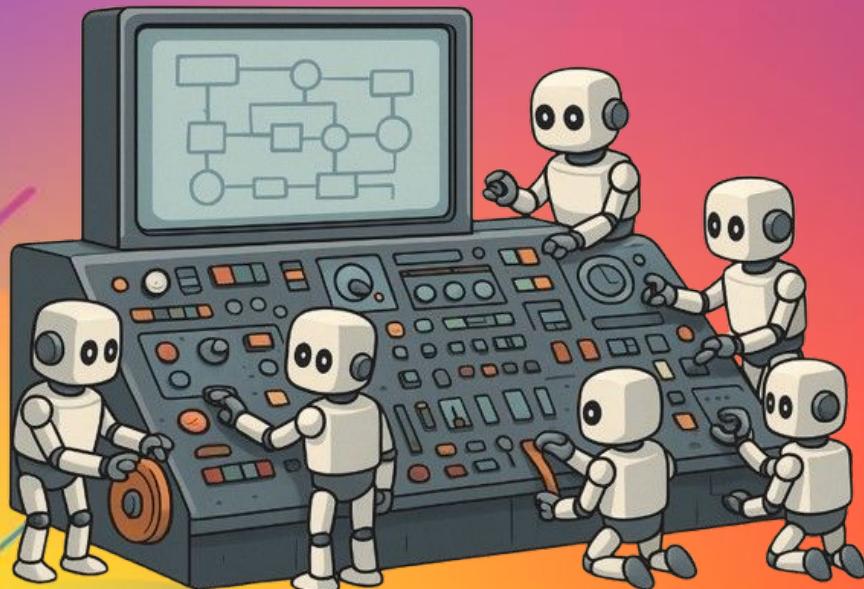


# AI Agents using Small Language Models



Hi, welcome to the tutorial.  
We're happy to have you here!



**Jeanfed Ramírez-Lima**

INAOE, Ford, jeanfed.ramirez@gmail.com

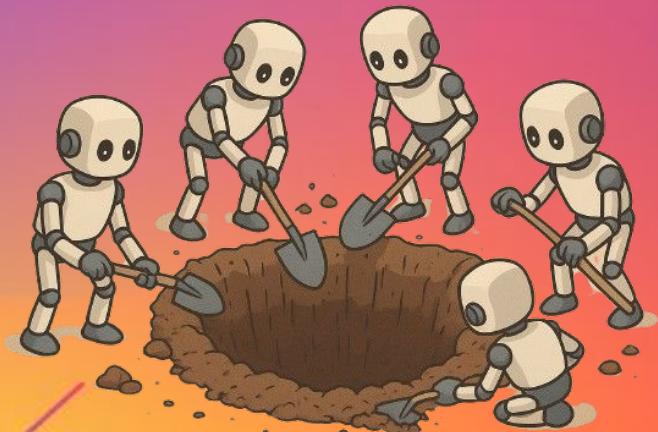
**Luis Joaquín Arellano**

INAOE, arellano.luis.3c@gmail.com

**Hugo Jair Escalante**

INAOE, hugo.jair@gmail.com





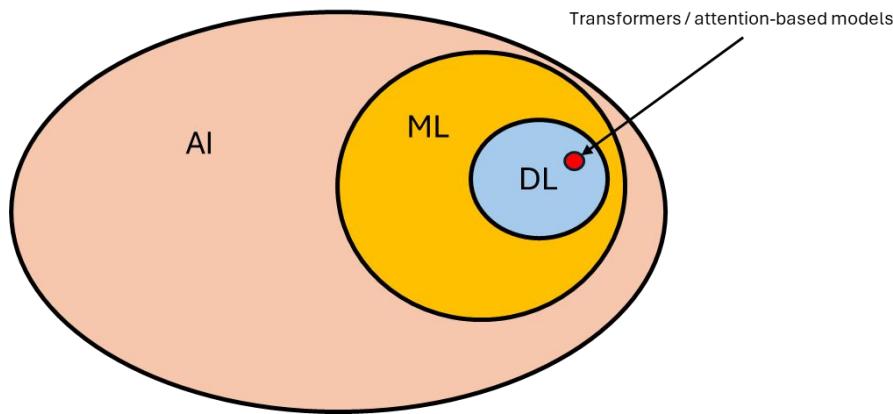
## Content

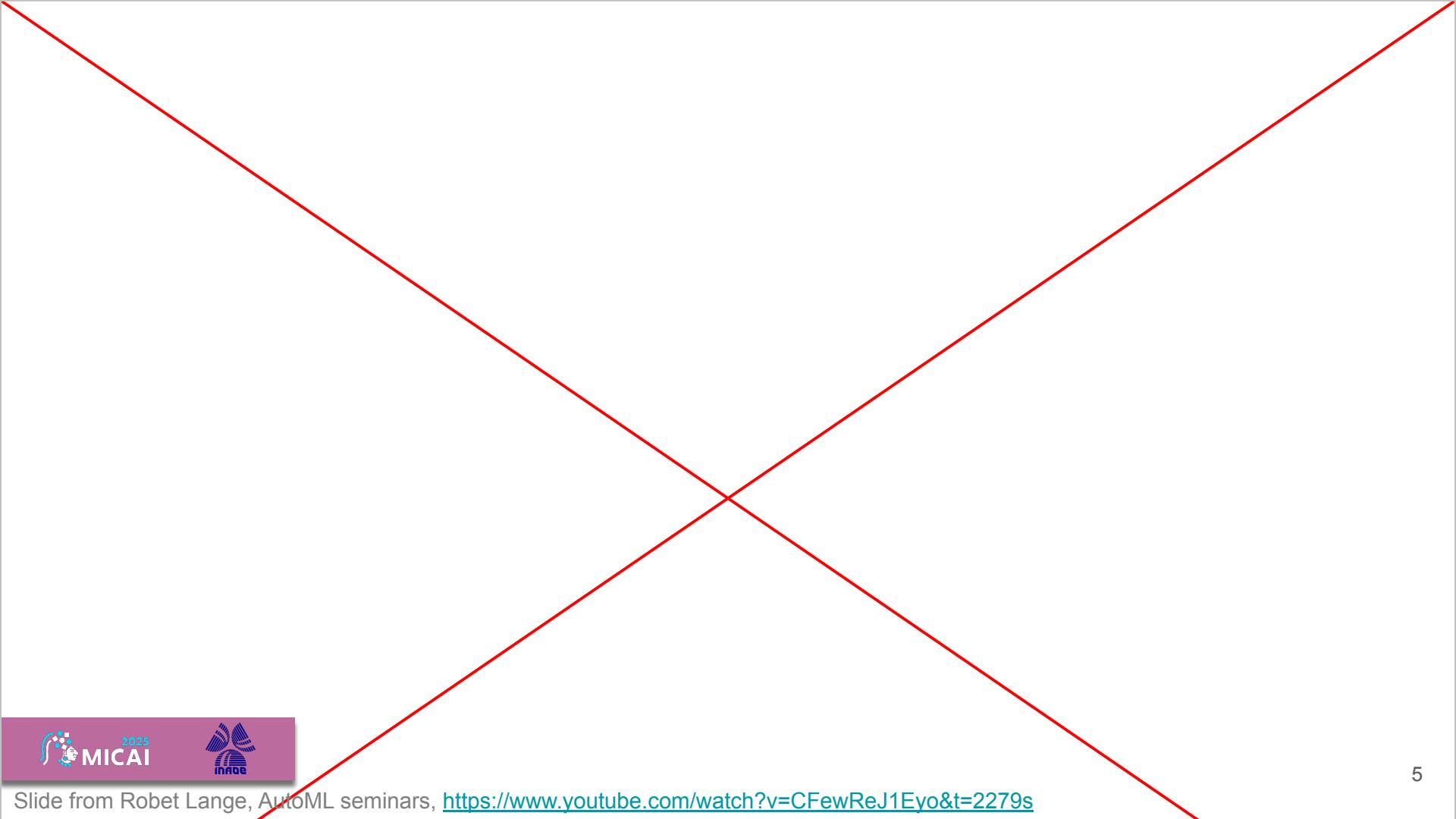
- Machine learning definition
- Deep learning history
- Neural architectures

# Deep Learning Intro

# Machine Learning and Artificial Intelligence

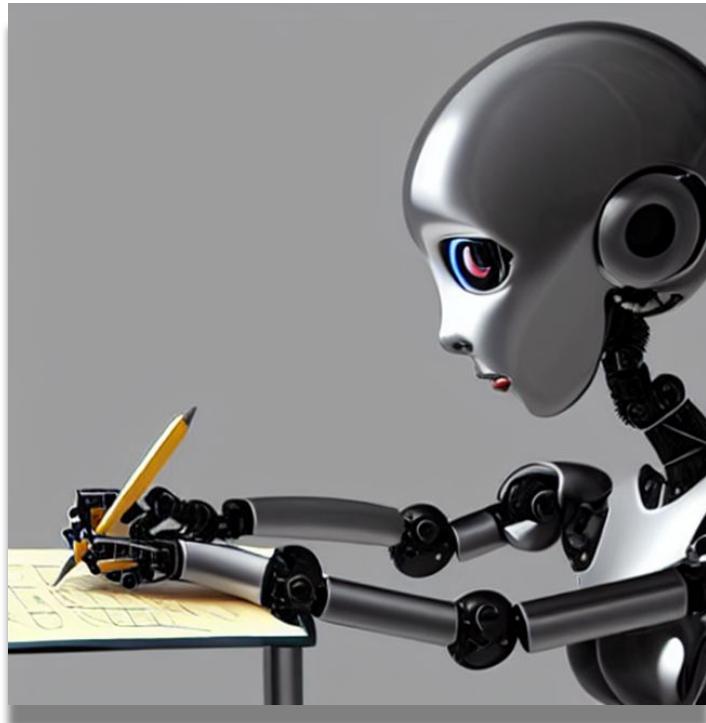
ML is a subfield of AI focusing on the development of computer programs that automatically adapt their behavior from their interaction with data. They can “*learn*” without being programmed explicitly



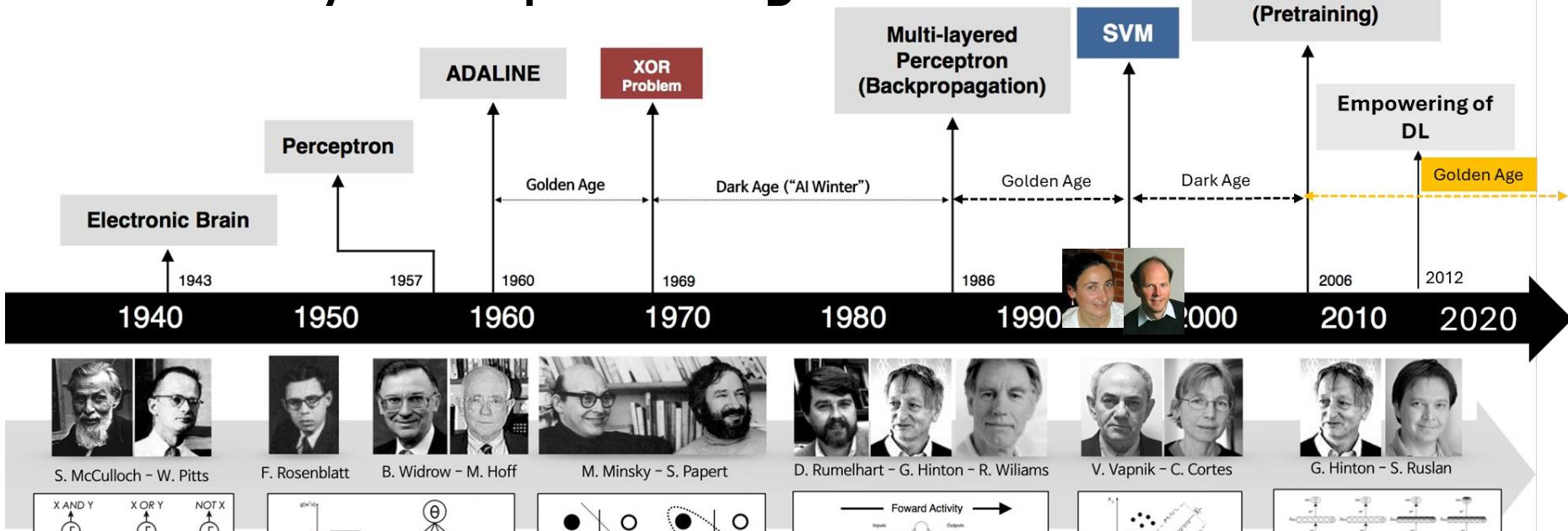


# Deep Learning

- Methodology of ML approaching problems through models organized in layers of parameters
- Layers capture descriptive & discriminative information from raw data directly
- Can be used for supervised / unsupervised learning, reinforcement learning, feature extraction, probabilistic modeling, etc.
- Some DL models are: *MLPs, DNNs, CNNs, Transformers, GANs, DBNs, AEs*, etc.



# Brief History of Deep Learning



- Adjustable Weights
- Weights are not Learned

- Learnable Weights and Threshold

- XOR Problem

- Solution to nonlinearly separable problems
- Big computation, local optima and overfitting

- Limitations of learning prior knowledge
- Kernel function: Human Intervention

- Hierarchical feature Learning

# THE NOBEL PRIZE IN PHYSICS 2024

Illustrations: Niklas Elmehed



John J. Hopfield

Geoffrey E. Hinton

Recent history  
(2024)



# THE NOBEL PRIZE IN CHEMISTRY 2024



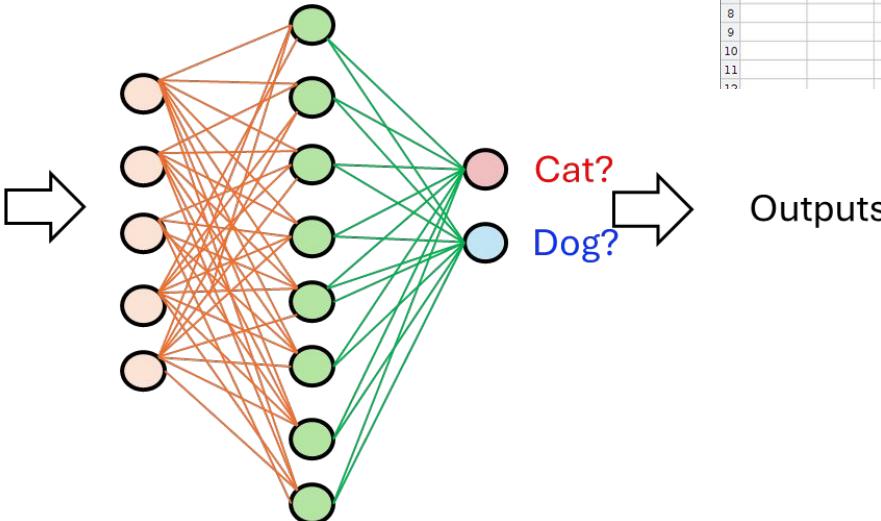
David  
Baker

Demis  
Hassabis

John M.  
Jumper

"for protein structure prediction"

# Multi-layer Perceptron

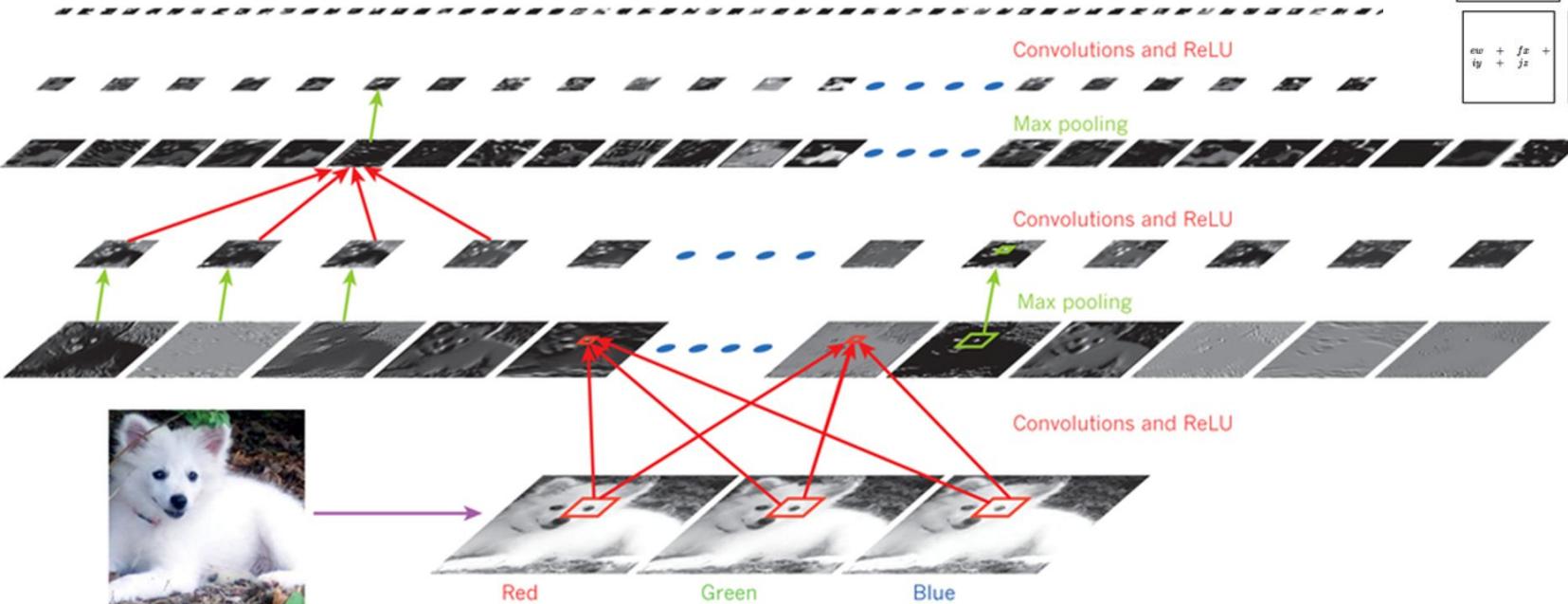
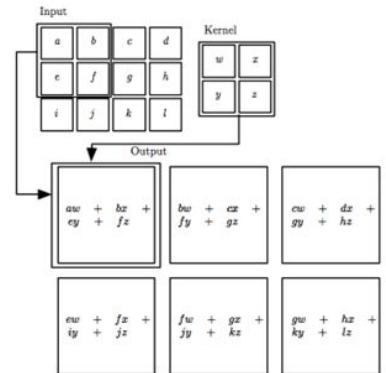


Variables - U									
1	2	3	4	5	6	7	8	9	
1	-0.7203	1.1714	-1.5029	3.3133	0.0338	-0.3660	-0.2886	-1.0836	
2	-0.5102	0.2660	0.3880	0.8725	-0.5068	-1.6706	1.5959	1.0228	
3	-1.6597	-0.5933	-0.4178	1.8309	1.0331	0.1132	-0.3820	1.0903	
4	0.3005	0.3909	-0.2304	-2.3196	3.2656	-0.0441	0.4254	0.0461	
5	0.2737	-0.6014	0.0314	0.5855	0.3692	2.1029	1.5813	-0.9019	
6									
7									
8									
9									
10									
11									
12									
13									

Variables - W		
1	2	3
1	0.5884	0.1427
2	0.9775	0.6085
3	0.2137	0.6286
4	0.6291	0.0238
5	0.7810	0.7787
6	0.3571	0.4912
7	0.9407	0.4116
8	0.9616	0.7775
9		
10		
11		
12		
13		

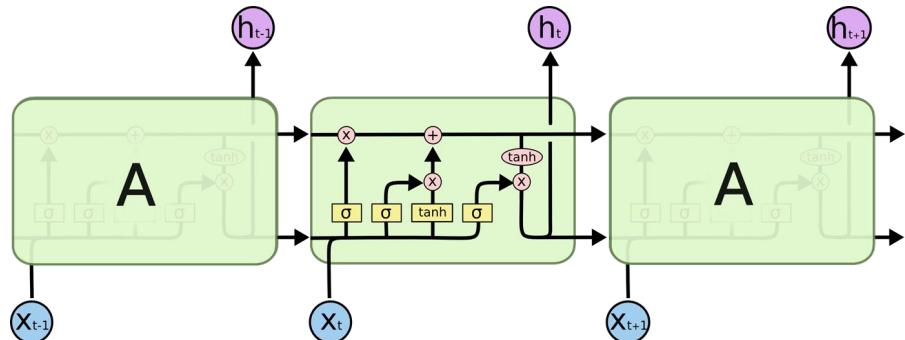
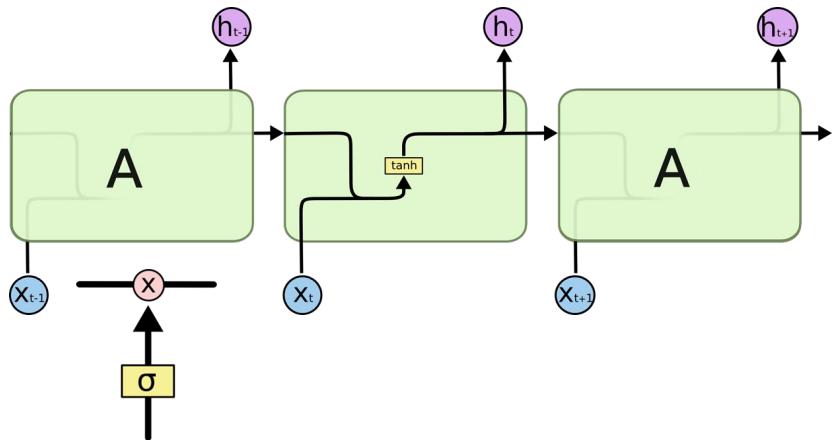
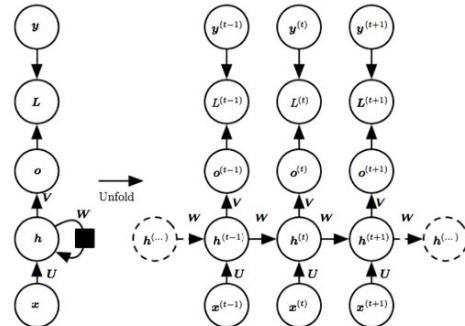
# Convolutional Neural Networks

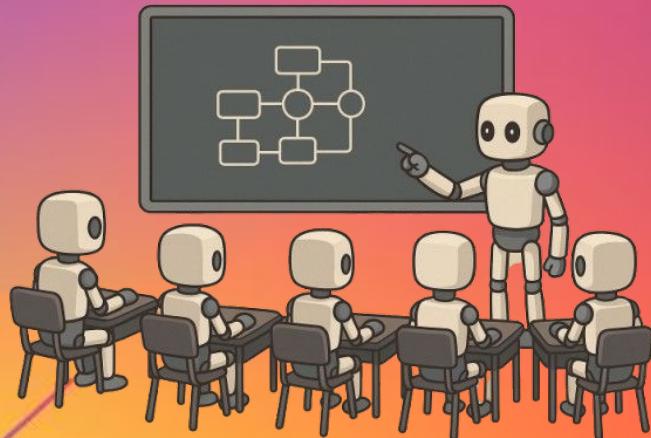
Samoyed (16); Papillon (5.7); Pomeranian (2.7); Arctic fox (1.0); Eskimo dog (0.6); white wolf (0.4); Siberian husky (0.4)



# Sequential Models

## RNNs, LSTMs, and GRUs





## Content

- Attention mechanism
- Transformer architecture
- Encoders and decoders

# Attention and Transformers

# Machine Learning and Artificial Intelligence

*.... If modern artificial intelligence has a founding document, a sacred text, it is Google's 2017 research paper "Attention Is All You Need."...*

Rob Toews, Forbes, 2023

## Attention Is All You Need

Ashish Vaswani\*  
Google Brain  
avaswani@google.com

Noam Shazeer\*  
Google Brain  
noam@google.com

Niki Parmar\*  
Google Research  
nikip@google.com

Jakob Uszkoreit\*  
Google Research  
usz@google.com

Llion Jones\*  
Google Research  
llion@google.com

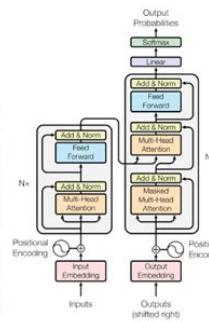
Aidan N. Gomez\* †  
University of Toronto  
aidan@cs.toronto.edu

Lukasz Kaiser\*  
Google Brain  
lukaszkaiser@google.com

Illia Polosukhin\* ‡  
illia.polosukhin@gmail.com

### Abstract

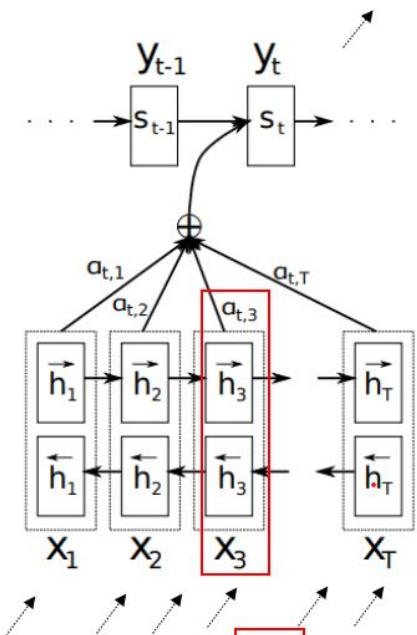
The dominant sequence transduction models are based on complex convolutional neural networks that include an encoder and a decoder performing models also connect the encoder and decoder through a mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolution. Experiments on two machine translation tasks show that the Transformer outperforms existing methods by a large margin while being more parallelizable and requiring less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results by over 2 BLEU. On the WMT 2014 English-to-French task our model establishes a new single-model state-of-the-art BLEU score training for 3.5 days on eight GPUs, a small fraction of the training time required by the best models from the literature.



# Attention Mechanisms

- Neural machine translation:
  - A weight for the decoder considering, inputs, hidden states and a (learnable) **context vector** derived from hidden states
- The context aligns input and output sequences, indicating what input words are more likely to be related to each output word

A dog is standing on a hardwood floor



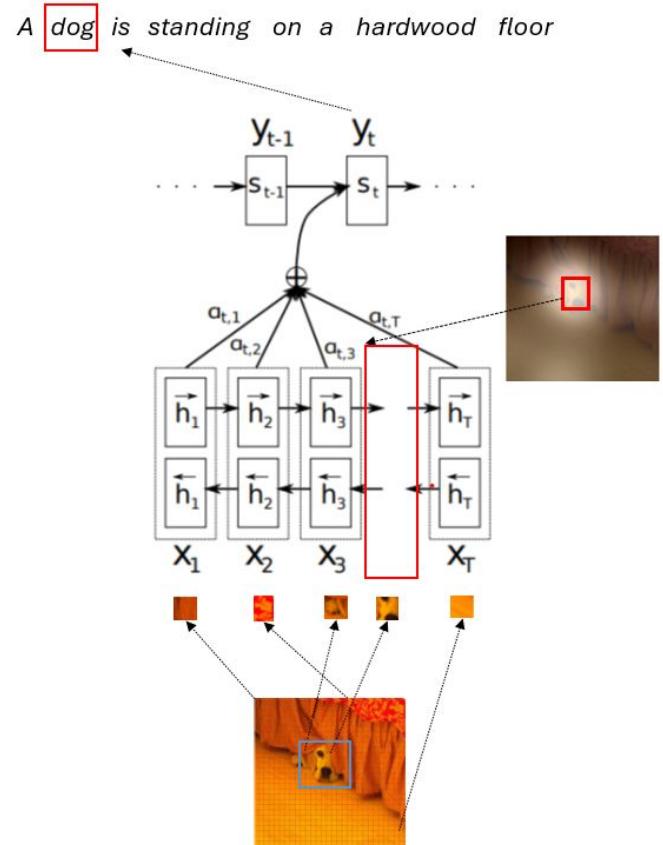
Un perro está parado sobre un piso de madera

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j.$$

$$p(y_i | y_1, \dots, y_{i-1}, \mathbf{x}) = g(y_{i-1}, s_i, c_i),$$

# Attention Mechanisms

- Neural machine translation:
  - A weight for the decoder considering, inputs, hidden states and a (learnable) **context vector** derived from hidden states
- The context aligns input and output sequences, indicating what input words are more likely to be related to each output word

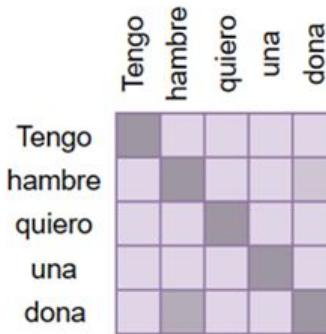


$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j. \quad p(y_i | y_1, \dots, y_{i-1}, \mathbf{x}) = g(y_{i-1}, s_i, c_i),$$

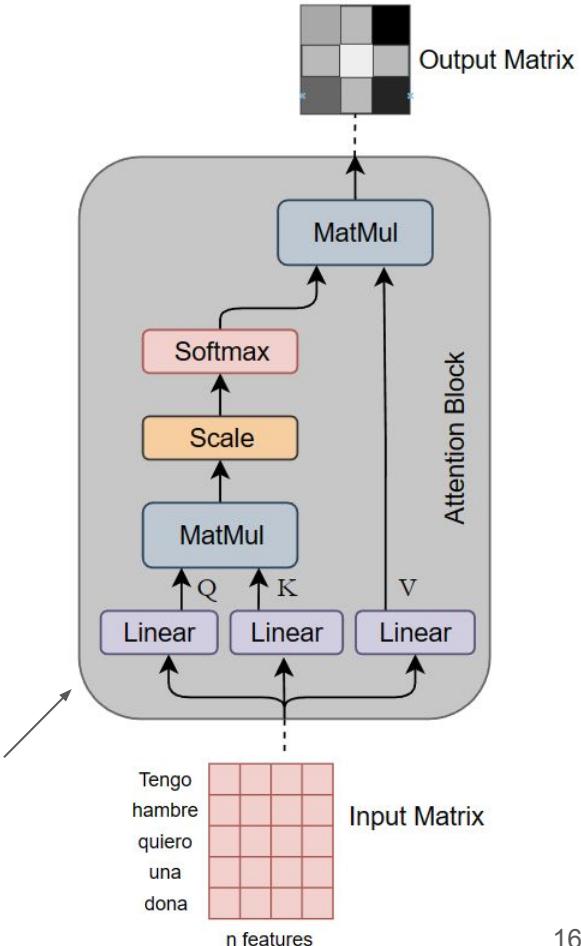
# Transformers

Self attention: how relevant are the words in a sentence to each other?

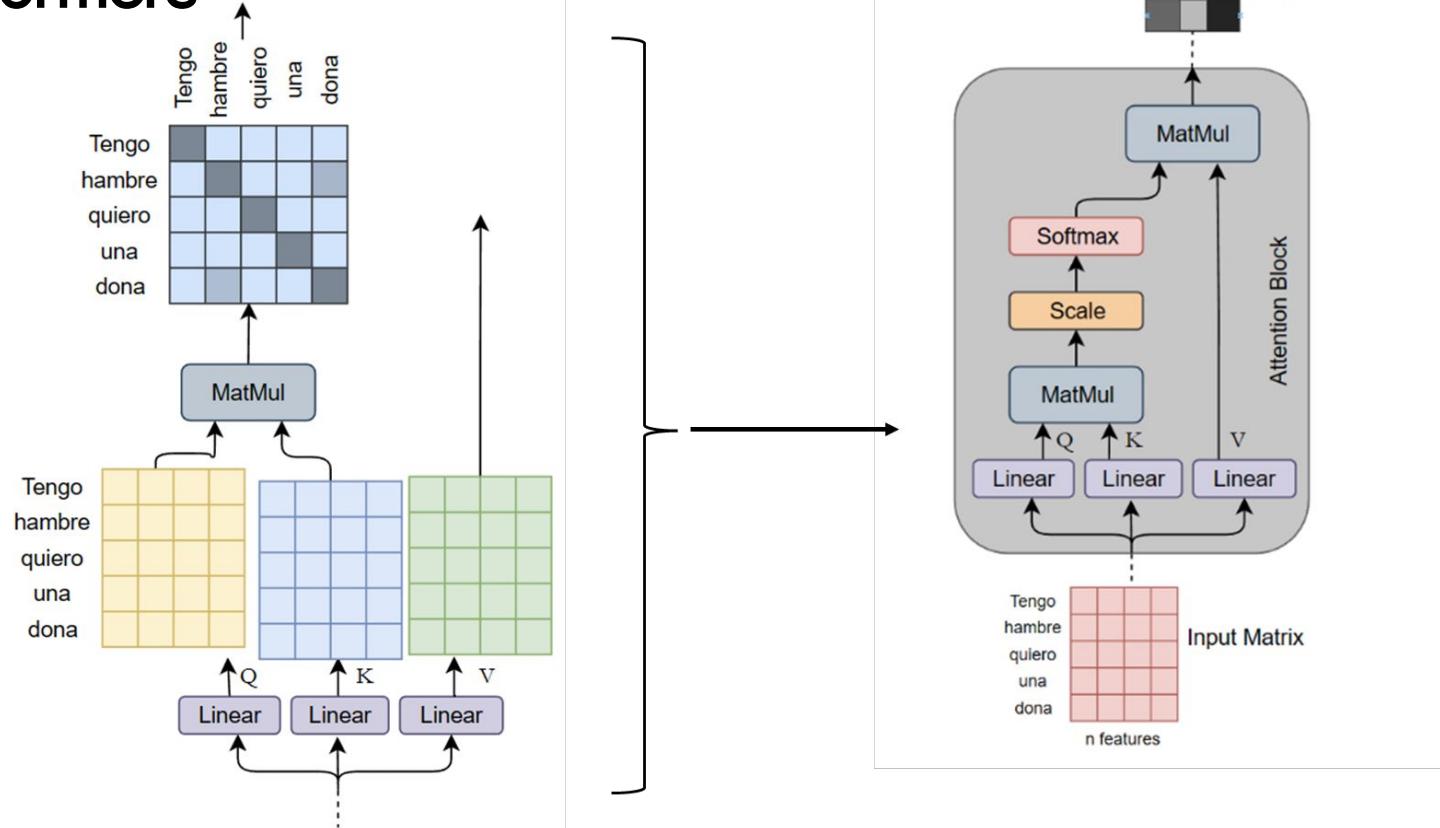
“Tengo hambre, quiero una dona”



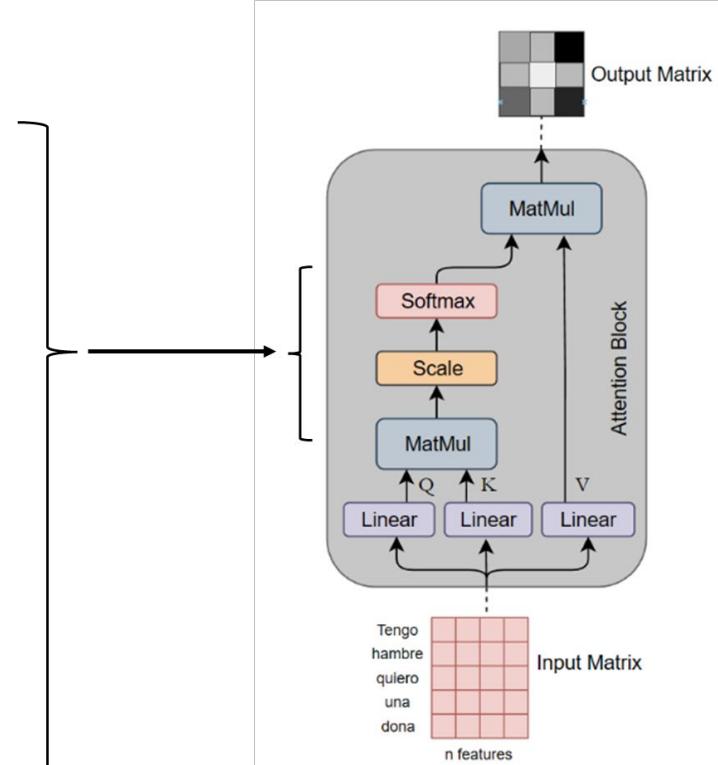
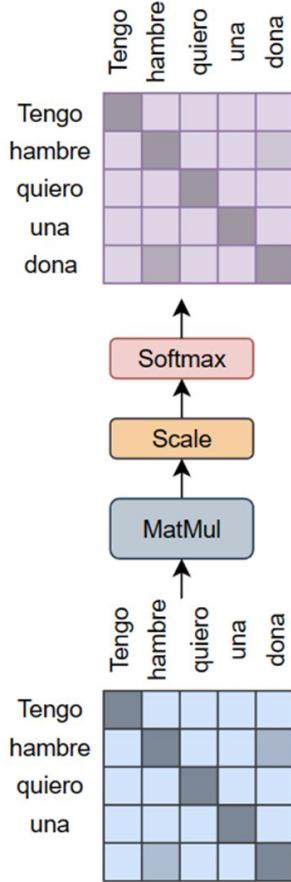
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



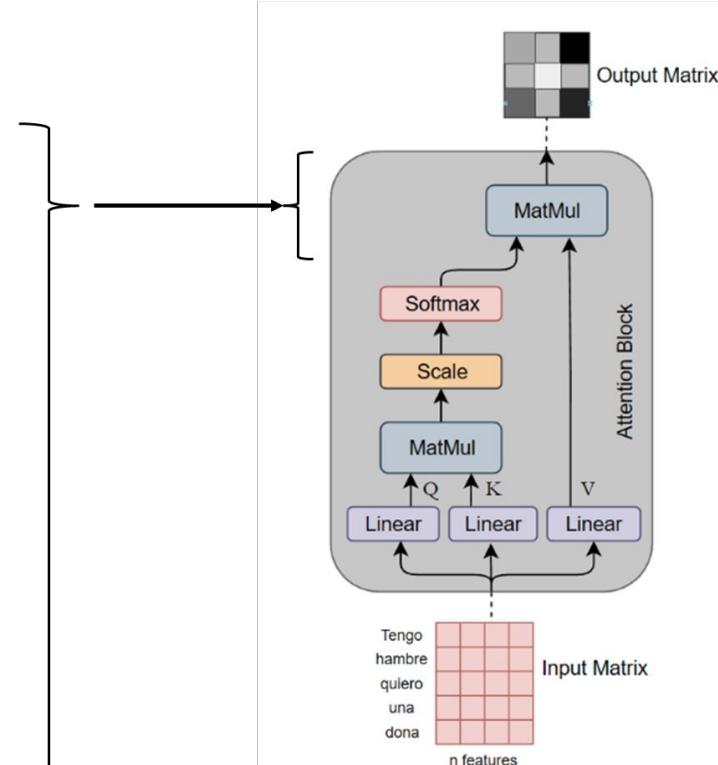
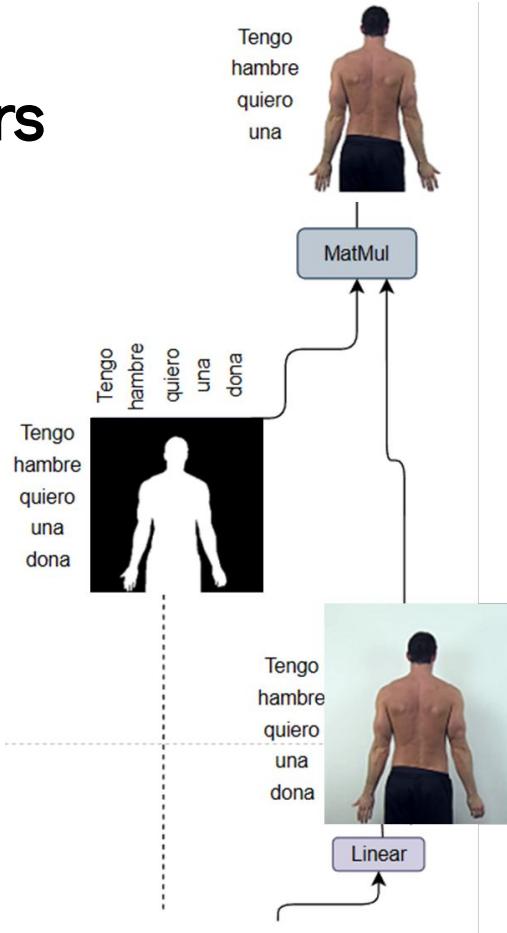
# Transformers



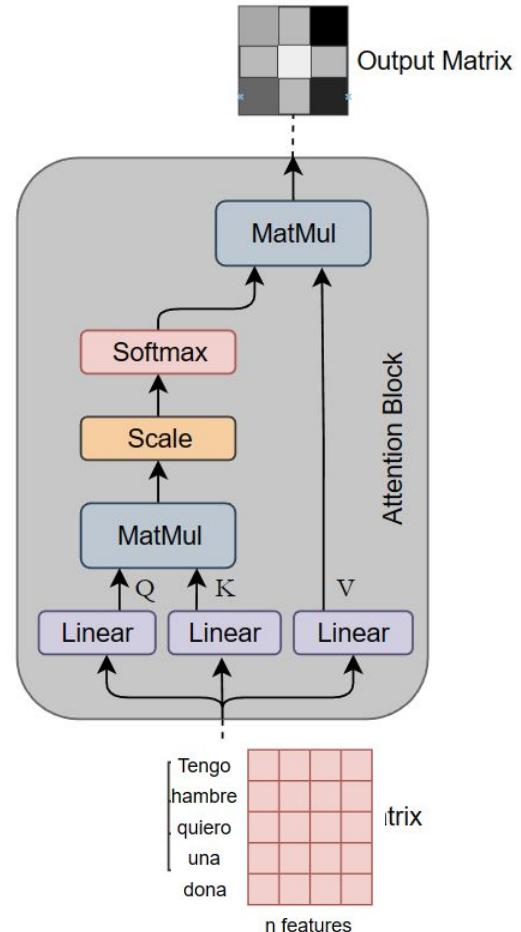
# Transformers



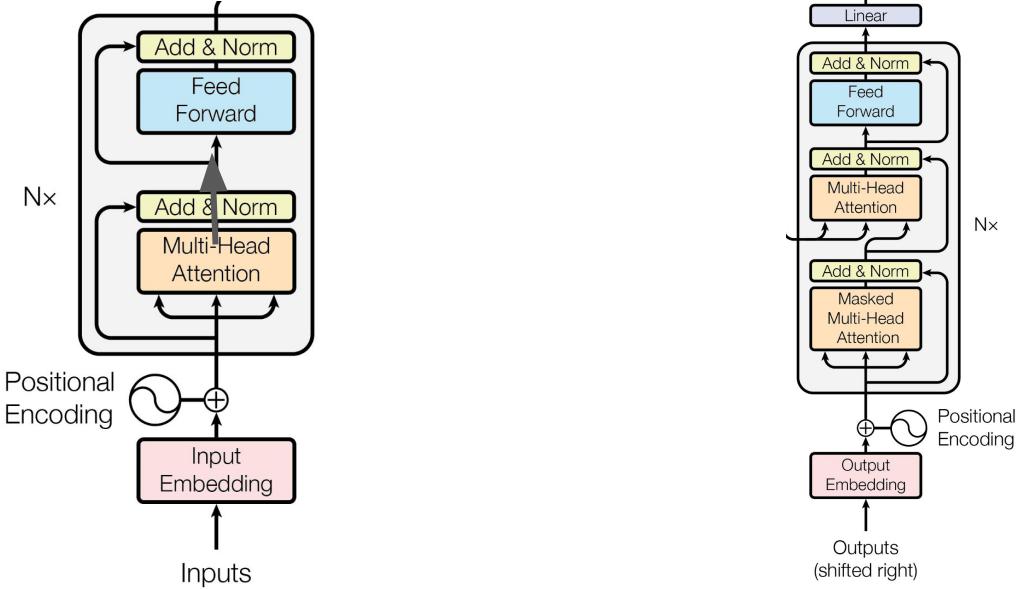
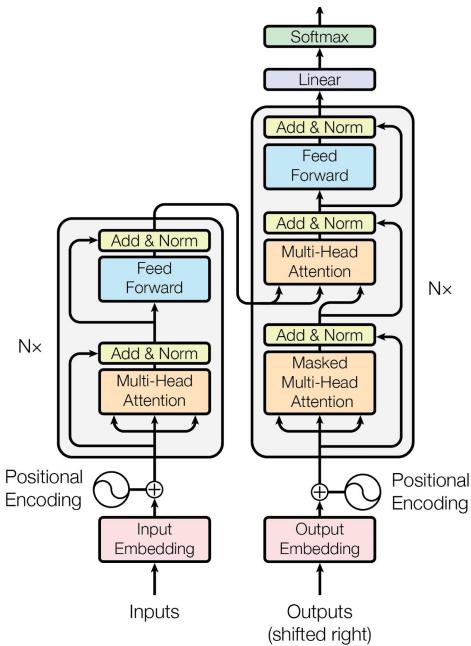
# Transformers



# Practice - The Attention Mechanism



# Architectures



## Encoder

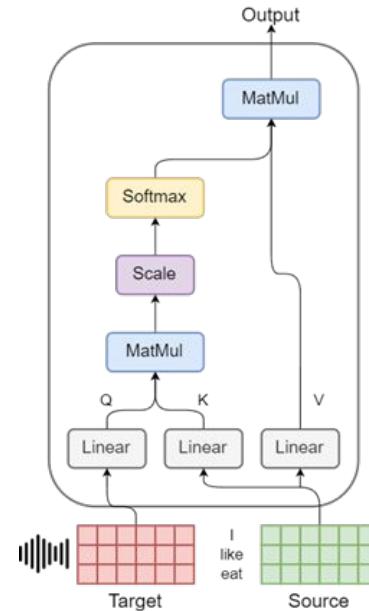
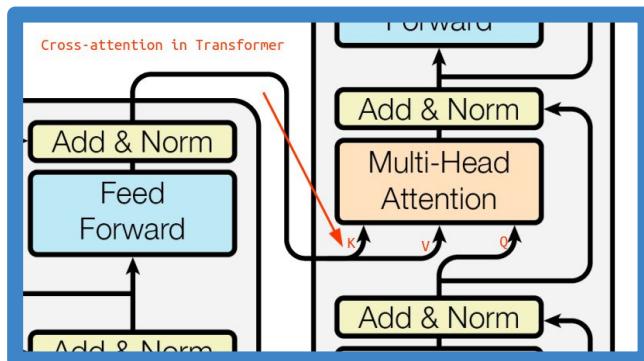
- *Positional encoding*
- *Usual output: embeddings*
- *MHA, AN+FF*

## Decoder

- *Autoregressive*
- *Usual output: token probability*
- *MHA, AN+FF, PE*

# Transformers

**Cross attention:** how relevant are the *words* in a sentence to words from another sentence ( $K$  and  $V$  from a sentence,  $Q$  from the other one)



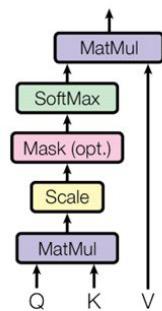
Ashish Vaswani et al. Attention is all you need. NeurIPS 2017

<http://nlp.seas.harvard.edu/2018/04/03/attention.html>  
<https://www.youtube.com/watch?v=8wUVJRm8DX0>

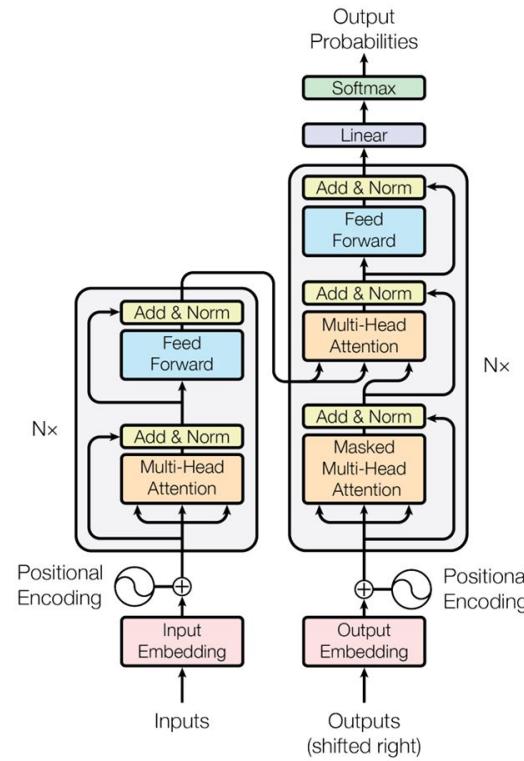
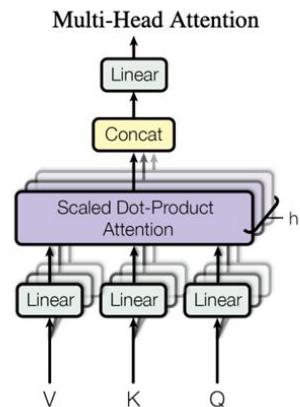
# Transformers

## Multi-head attention: parallel attention layers

Scaled Dot-Product Attention



Multi-Head Attention



Ashish Vaswani et al. Attention is all you need. NeurIPS 2017

<http://nlp.seas.harvard.edu/2018/04/03/attention.html>

<https://www.youtube.com/watch?v=8wUVJRm8DX0>



<https://tinyurl.com/3vnz3kmu>

# Transformers

**Tokenization:** a way to split sequential data into primitives (tokens)

**Positioning encoding:** The only explicit mechanism to account for sequential information

$e_0$	$p_0$	$e_1$	$p_1$	$e_2$	$p_2$	$e_3$	$p_3$
0.42	0	0.87	1	0.02	2	0.02	3
0.31	0	-0.64	1	0.01	2	0.01	3
0.73	0	0.81	1	-0.24	2	-0.24	3
0.36	0	0.26	1	-0.07	2	-0.07	3
0.99	0	-0.35	1	0.00	2	0.00	3

## ORIGINAL TOKENIZATION

Leicestershire

beat

Somerset

by

an

innings

## WORDPIECES

Leicester

##shire

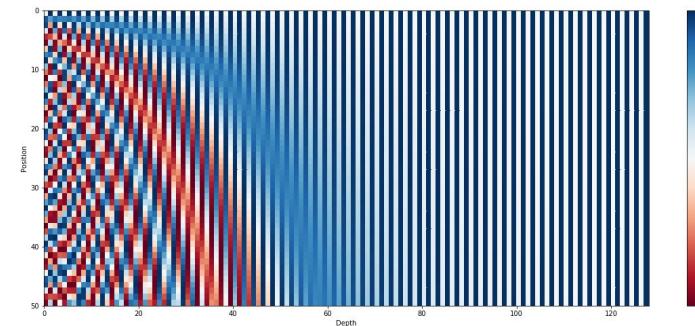
beat

Somerset

by

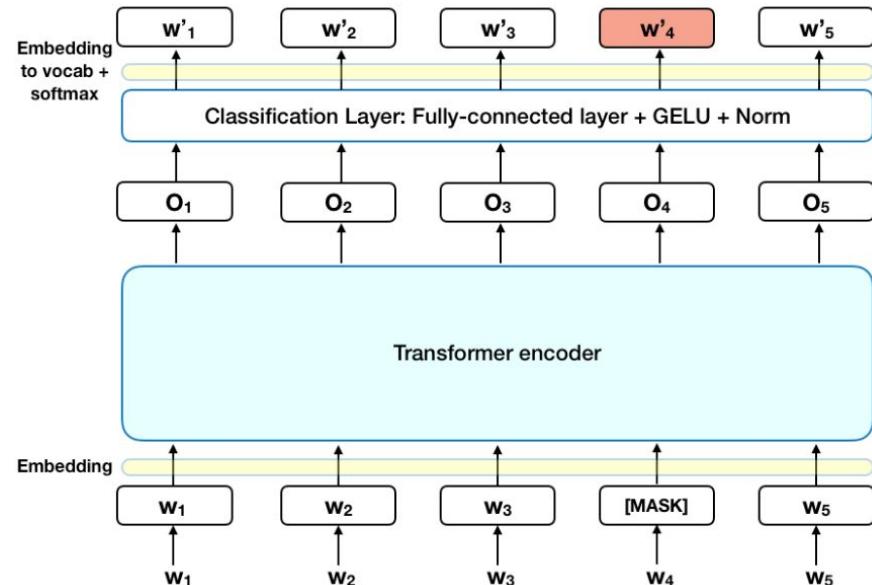
an

innings



# Popular transformers: BERT

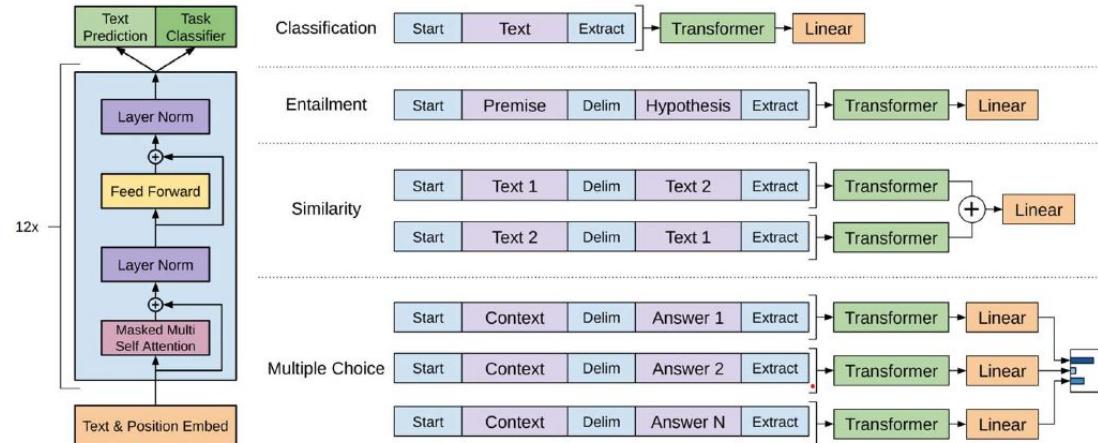
- Bidirectional Encoder Representations from Transformers
  - Large scale pretraining
  - Pretext tasks: MLM & NSP
  - Bidirectional encoder model
  - Special tokens



# Popular Transformers: GPT

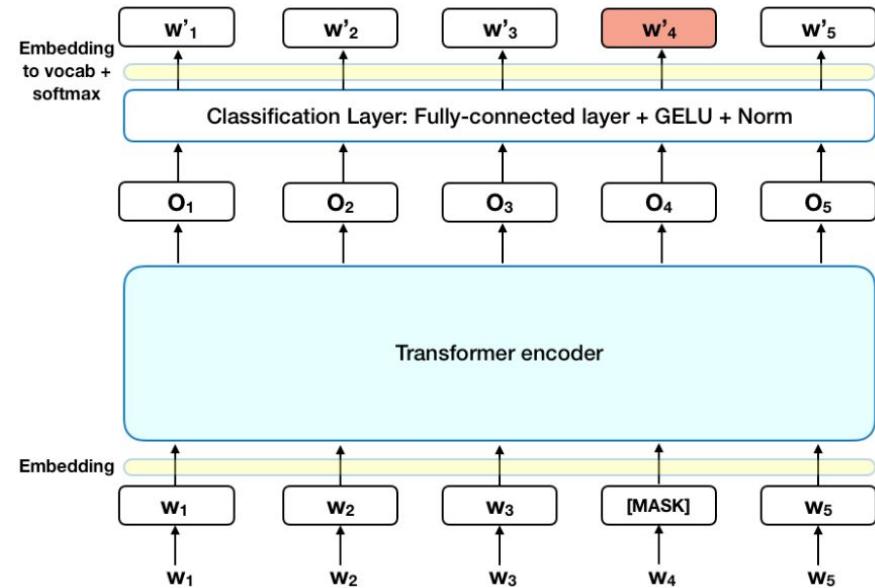
- Generative pretraining transformer
  - Large scale pretraining
  - Autoregressive decoder language model (left to right)
  - Fine tuning for solving downstream tasks for the task at hand

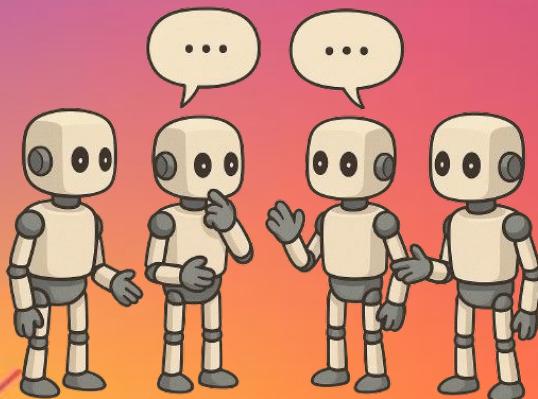
$$L_1(\mathcal{U}) = \sum \log P(u_i|u_{i-k}, \dots, u_{i-1}; \Theta)$$



A. Radford et al. Improving Language Understanding by Generative Pre-Training. 2018  
[https://cdn.openai.com/research-covers/language-unsupervised/language\\_understanding\\_paper.pdf](https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf)

# Practice - Using BERT to encode a sentence





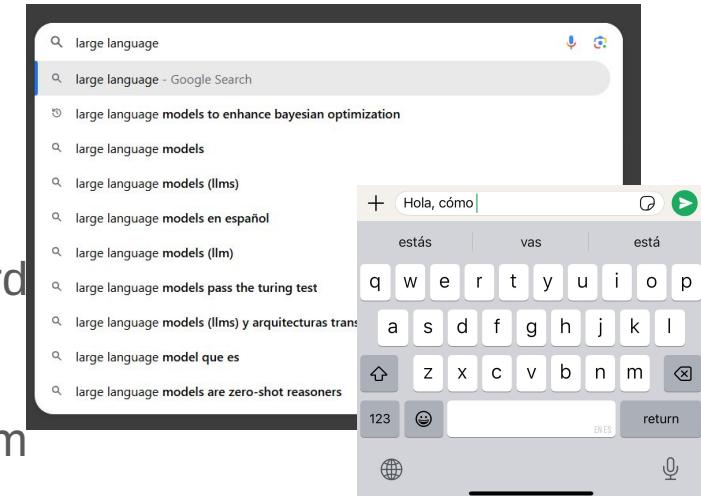
## Content

- Language model definition
- Architecture visualization
- Quantization
- Prompting techniques

# Language Models

# Language Models

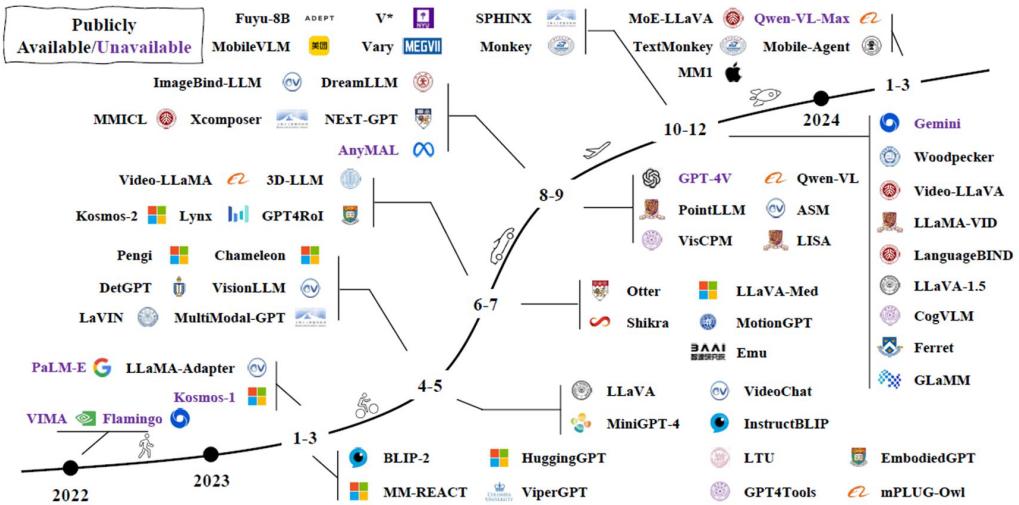
- A language model learns to predict the next word in a sentence
  - It aims to “*understand*” and generate language according to patterns learned from a large corpus of documents



- Maximize:  $L_1(\mathcal{U}) = \sum \log P(u_i | u_{i-k}, \dots, u_{i-1}; \Theta)$
- LLMs are LMs trained on massive amounts of data

# Large Language Models

- How many LLMs are out there?
- Are LLMs restricted to text?
- Which one is better?
- What are the criteria used to evaluate LLMs?

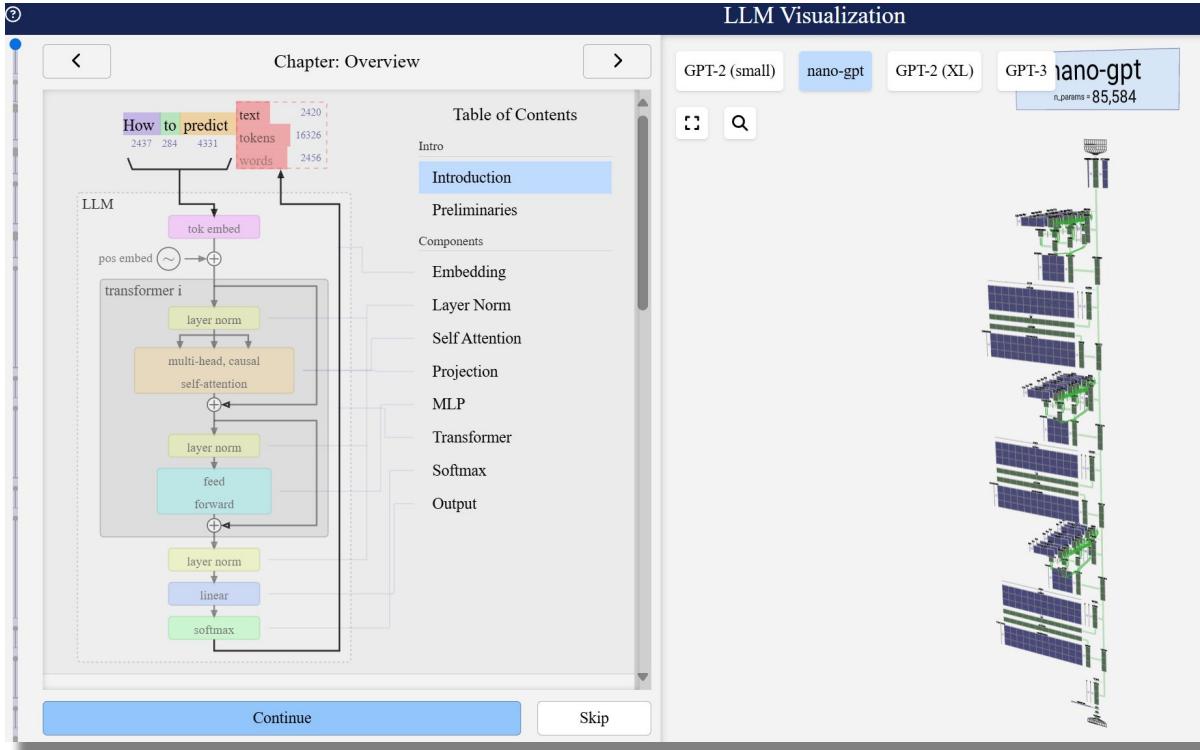


[https://en.wikipedia.org/wiki/List\\_of\\_large\\_language\\_models](https://en.wikipedia.org/wiki/List_of_large_language_models)  
[https://www.vellum.ai/llm-leaderboard?utm\\_source=google&utm\\_medium=organic](https://www.vellum.ai/llm-leaderboard?utm_source=google&utm_medium=organic)  
<https://huggingface.co/spaces/lmarena-ai/lmarena-leaderboard>

S. Yin et al. A Survey on Multimodal Large Language Models. TPAMI, 2025,  
<https://github.com/BradyFU/Awesome-Multimodal-Large-Language-Models>

Rank (UB)	Model	Score	Votes
1	G gemini-2.5-pro	1451	54,087
1	AI claudie-opus-4-1-20250805-thi...	1447	21,306
1	AI claudie-sonnet-4-5-20250929-t...	1445	6,287
1	G gpt-4.5-preview-2025-02-27	1441	14,644
2	AI chatgpt-4-0-latest-20250326	1440	40,013
2	AI o3-2025-04-16	1440	51,293
2	AI claudie-sonnet-4-5-20250929	1438	6,144
2	G gpt-5-high	1437	23,580
2	AI claudie-opus-4-1-20250805	1437	33,298
3	AI qwen3-max-preview	1434	18,078

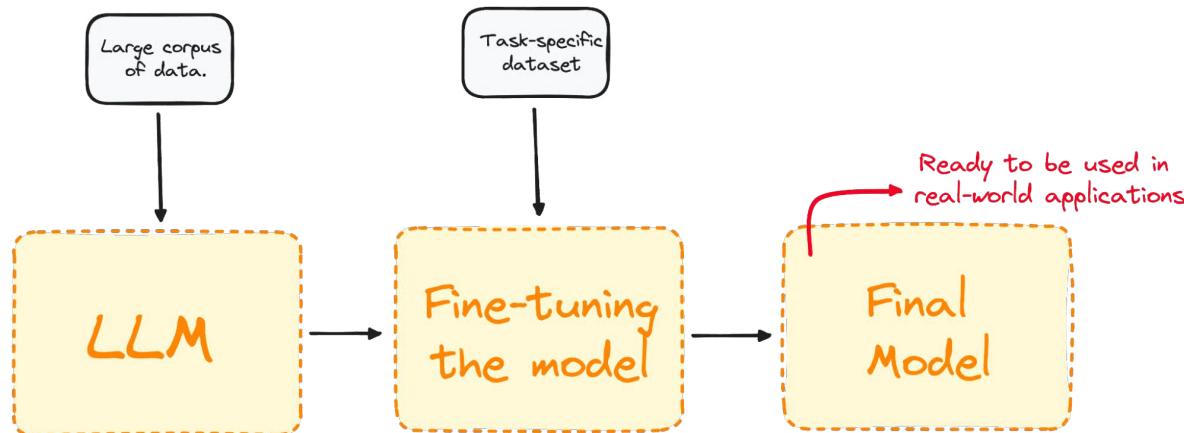
# Visualization of LLMs



Brendan Bycroft: <https://bbycroft.net/llm>

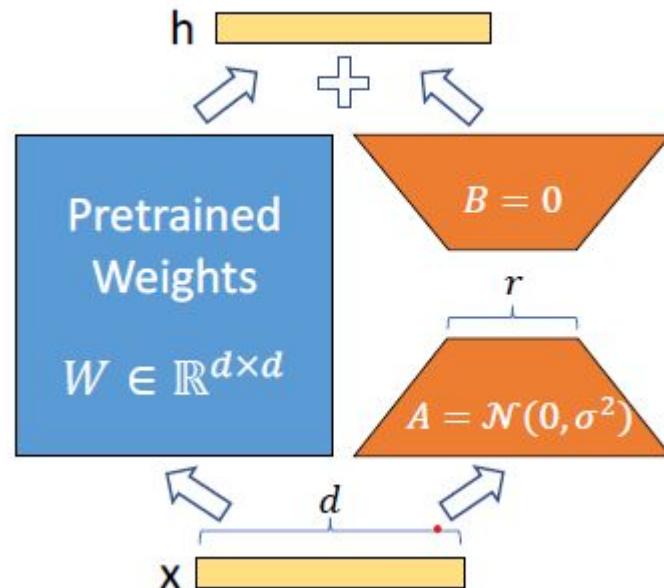
# LLM Fine tuning

- Given a (generic) trained LLM, adjust the model for a specific domain by running additional learning rounds over a task-specific dataset



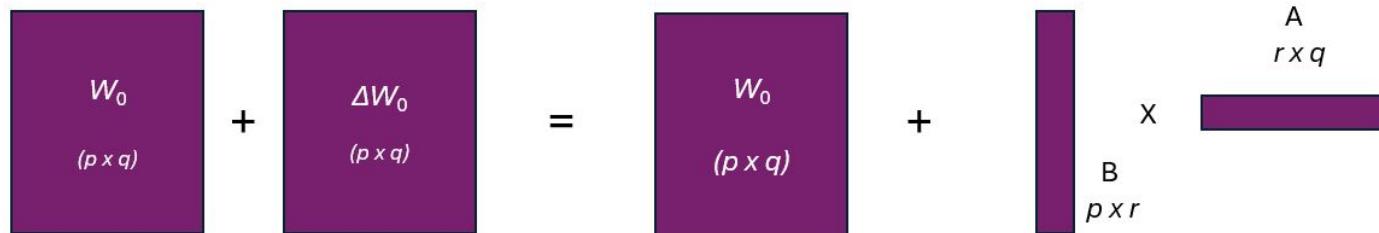
# Fine tuning with LoRA

- Efficient fine tuning by freezing the initial LLM and learning (small) decomposition matrix
- These matrices learn the changes in weights as tuning occurs



# Fine tuning with LoRA

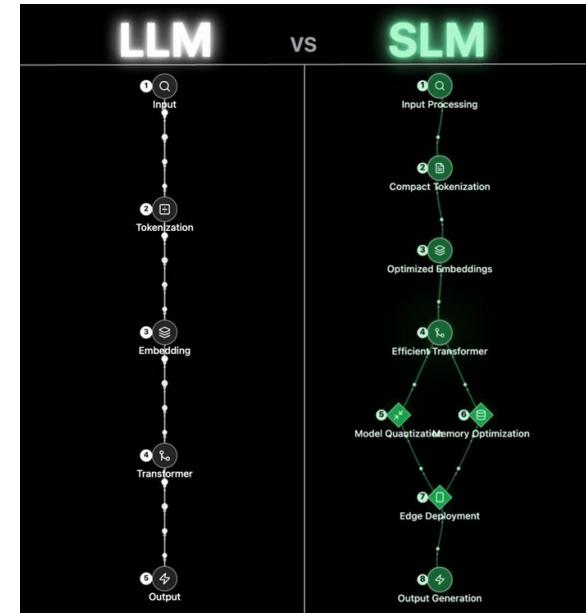
$$W_0 + \Delta W = W_0 + BA$$



# Small Language Models

“... Small Language Models (SLMs) are lightweight versions of traditional language models designed to operate efficiently on resource-constrained environments such as smartphones, embedded systems, or low-power computers... ”

Dimension	LLMs	SLMs
Parameters	Billions	Millions– few B
Speed	Slower	Faster
Cost	High	Low
Deployment	Cloud	Edge / Local
Examples	GPT-5, Claude 3 DeepSeek-V3.2 Llama 4, Gemini 2.5	Phi-3.5-Mini-3.8B SmolLM2-1.7B DeepSeeek-R1-1.5B



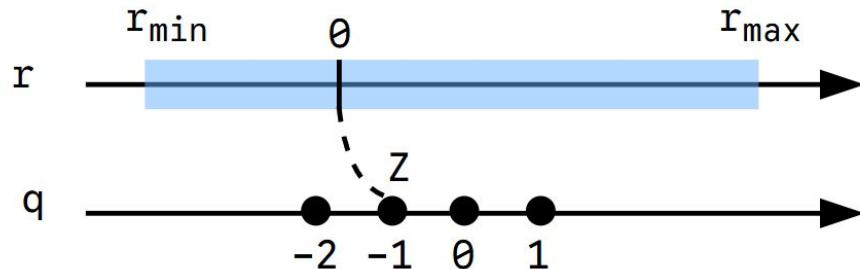
Credit: Manthan Patel

# Model Quantization

How to go from a Large Language Model to a *Small Language Model?*

- Asymmetric Linear Quantization

Represent real values with a finite number of states



32-bit float			2-bit signed int		
2.09	-0.98	1.48	1	-2	0
0.05	-0.14	-1.08	-1	-1	-2
-0.91	1.92	0	-2	1	-1

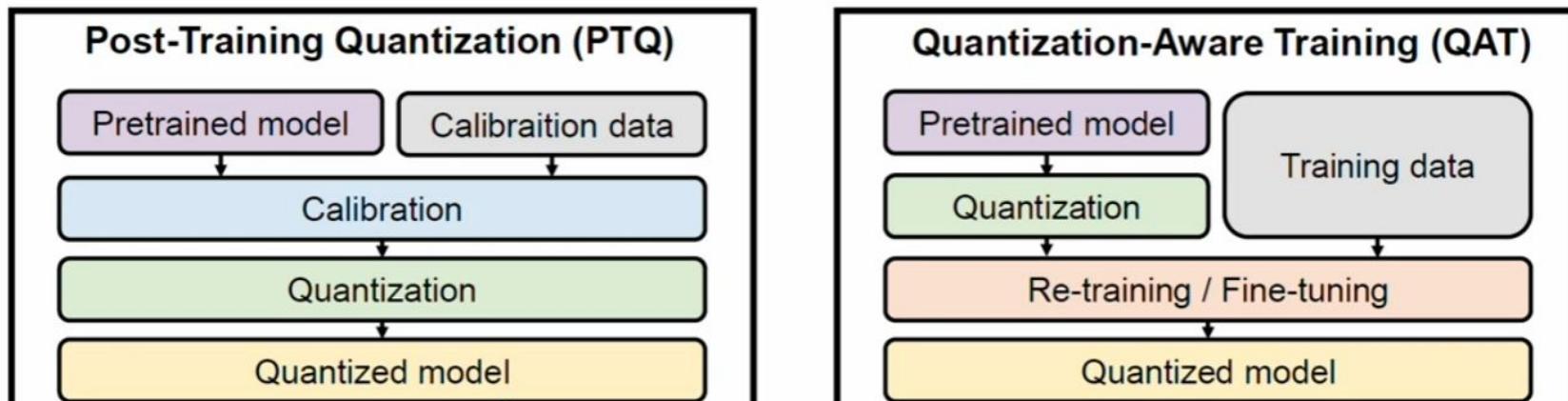
Must approximate from the encoded states to the original value...

$$S = \frac{(r_{\max} - r_{\min})}{(q_{\max} - q_{\min})}$$
$$r = (q - z) \times S$$

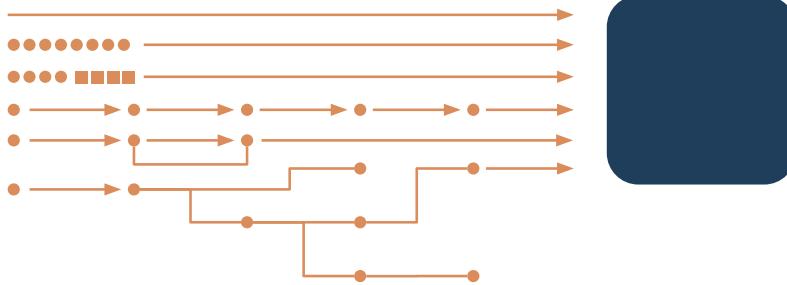
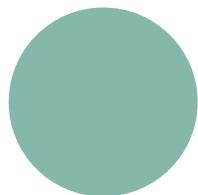
# Model Quantization

There are two common quantization techniques

- Two types of DNN quantization: QAT and PTQ
  - Post-Training Quantization (PTQ) directly converts parameters and activations to reduced-precision
  - Quantization-Aware Training (QAT) exploits re-training to compensate quantization error



# Prompting Techniques



- Zero-Shot
- Few-Shot
- Chain-of-Thought Prompting
- Self-Consistency
- Generated Knowledge
- ToT
- RAG

# Zero-Shot Prompting



**PROMPT:** \_\_\_\_\_

Classify the text into neutral, negative or positive.

Text: I think the vacation is okay.  
Sentiment:

**OUTPUT:** \_\_\_\_\_

# Zero-Shot Prompting



**PROMPT:** \_\_\_\_\_

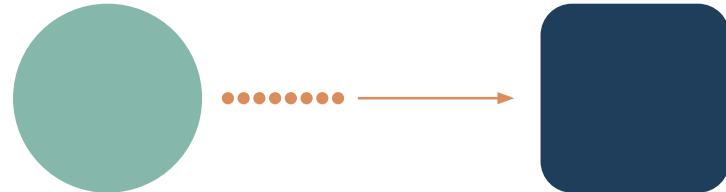
Classify the text into neutral, negative or positive.

Text: I think the vacation is okay.  
Sentiment:

**OUTPUT:** \_\_\_\_\_

neutral

# Few-Shot Prompting



PROMPT: \_\_\_\_\_

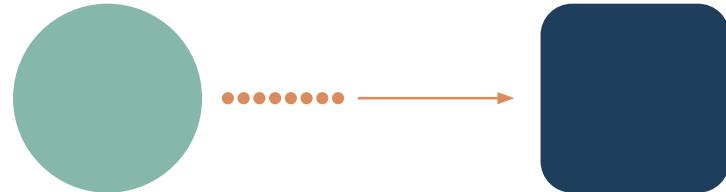
A "whatpu" is a small, furry animal native to Tanzania. An example of a sentence that uses the word whatpu is:

We were traveling in Africa and we saw these very cute whatpus.

To do a "farduddle" means to jump up and down really fast. An example of a sentence that uses the word farduddle is:

OUTPUT: \_\_\_\_\_

# Few-Shot Prompting



PROMPT: \_\_\_\_\_

A "whatpu" is a small, furry animal native to Tanzania. An example of a sentence that uses the word whatpu is:

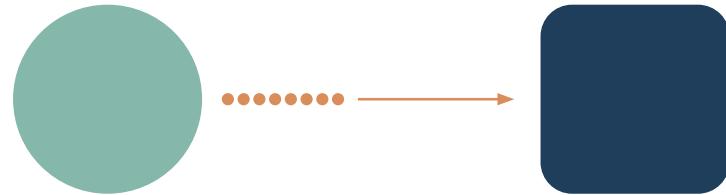
We were traveling in Africa and we saw these very cute whatpus.

To do a "farduddle" means to jump up and down really fast. An example of a sentence that uses the word farduddle is:

OUTPUT: \_\_\_\_\_

When we won the game, we all started to farduddle in celebration.

# Few-Shot Prompting (2)



PROMPT: \_\_\_\_\_

This is awesome! // Negative

This is bad! // Positive

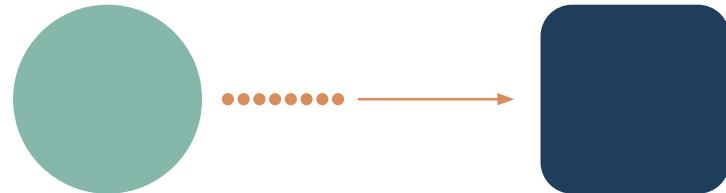
Wow that movie was rad! // Positive

What a horrible show! //

OUTPUT: \_\_\_\_\_

Negative

# Few-Shot Prompting (Limitations)



PROMPT: \_\_\_\_\_

The odd numbers in this group add up to an even number: 15, 32, 5, 13, 82, 7, 1.

A:

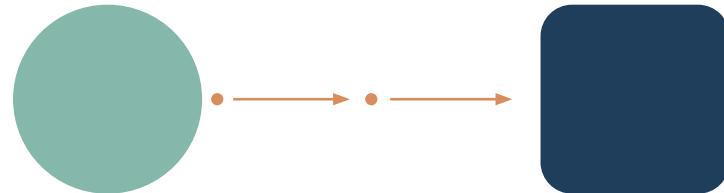
OUTPUT: \_\_\_\_\_

Yes, the odd numbers in this group add up to 107, which is an even number.

OUTPUT: \_\_\_\_\_

False

# Chain-of-Thought Prompting



PROMPT:

The odd numbers in this group add up to an even number: 4, 8, 9, 15, 12, 2, 1.

A: Adding all the odd numbers (9, 15, 1) gives 25. The answer is False.

The odd numbers in this group add up to an even number: 15, 32, 5, 13, 82, 7, 1.

A:

OUTPUT:

Adding all the odd numbers (15, 5, 13, 7, 1) gives 41. The answer is False.

# Self-Consistency Prompting

**PROMPT:** \_\_\_\_\_

Q: Michael had 58 golf balls. On tuesday, he lost 23 golf balls. On wednesday, he lost 2 more. How many golf balls did he have at the end of wednesday?

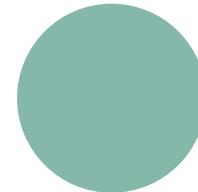
A: Michael initially had 58 balls. He lost 23 on Tuesday, so after that he has  $58 - 23 = 35$  balls. On Wednesday he lost 2 more so now he has  $35 - 2 = 33$  balls. The answer is 33.

Q: Olivia has \$23. She bought five bagels for \$3 each. How much money does she have left?

A: She bought 5 bagels for \$3 each. This means she spent \$15. She has \$8 left.

Q: When I was 6 my sister was half my age. Now I'm 70 how old is my sister?

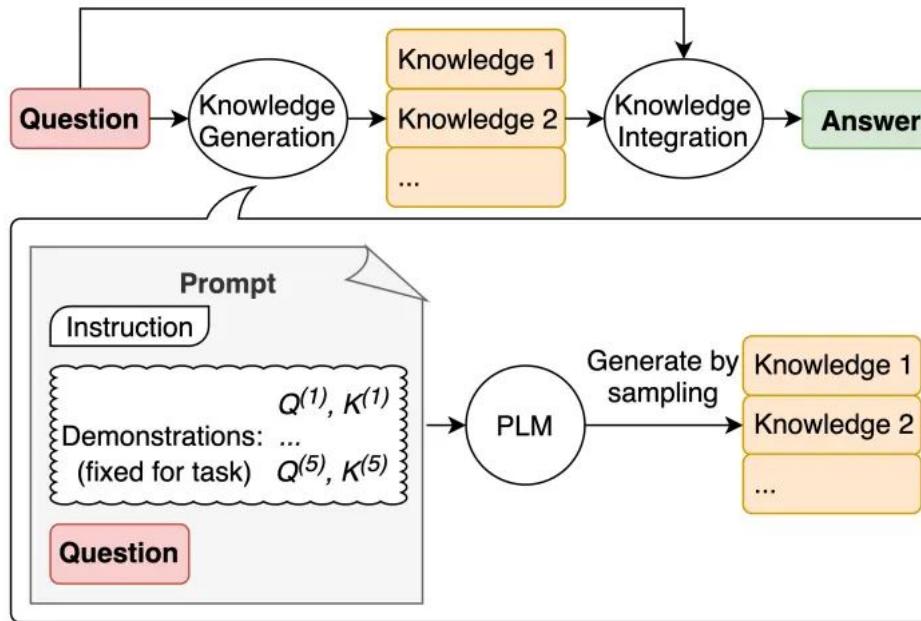
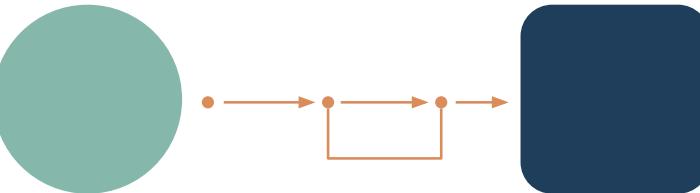
A:



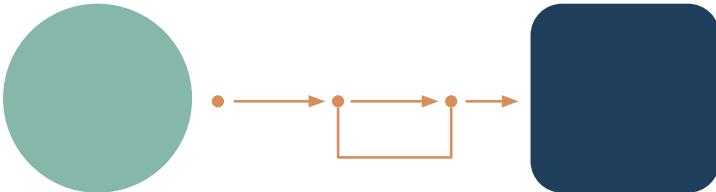
**OUTPUT:** \_\_\_\_\_

When I was 6 my sister was half my age, so she was 3. Now I am 70, so she is  $70 - 3 = 67$ . The answer is 67.

# Generated Knowledge Prompting



# Generated Knowledge Prompting

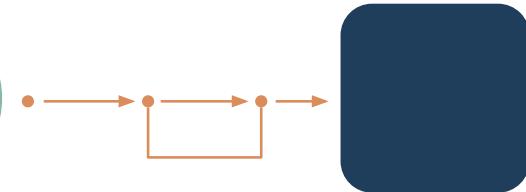
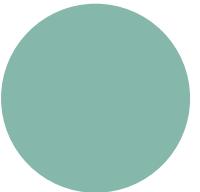


PROMPT: \_\_\_\_\_

Part of golf is trying to get a higher point total than others. Yes or No?

OUTPUT: \_\_\_\_\_

# Generated Knowledge Prompting



PROMPT: \_\_\_\_\_

**Input:** A rock is the same size as a pebble.

**Knowledge:** A pebble is a clast of rock with a particle size of 4 to 64 millimetres based on the Udden-Wentworth scale of sedimentology. Pebbles are generally considered larger than granules (2 to 4 millimetres diameter) and smaller than cobbles (64 to 256 millimetres diameter).

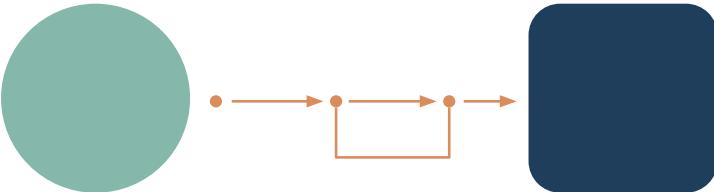
**Input:** Part of golf is trying to get a higher point total than others.

**Knowledge:**

OUTPUT: \_\_\_\_\_

The objective of golf is to play a set of holes in the least number of strokes. A round of golf typically consists of 18 holes. Each hole is played once in the round on a standard golf course. Each stroke is counted as one point, and the total number of strokes is used to determine the winner of the game.

# Generated Knowledge Prompting



PROMPT: \_\_\_\_\_

Part of golf is trying to get a higher point total than others. Yes or No?

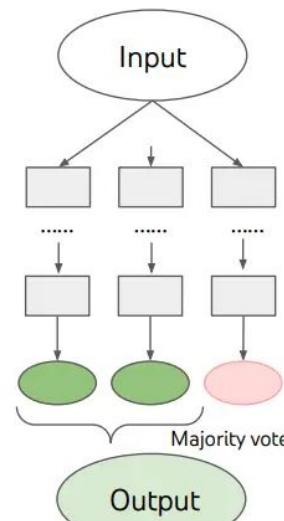
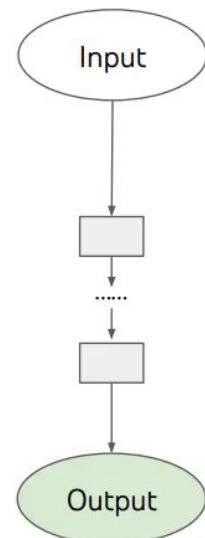
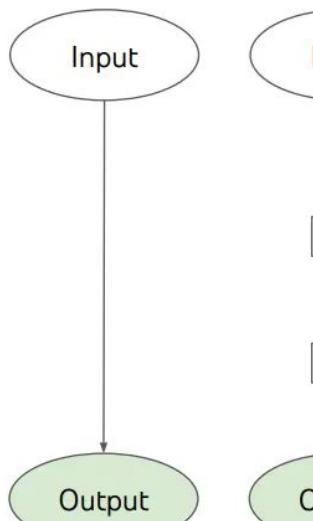
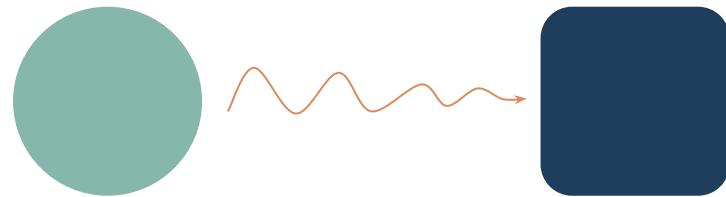
Knowledge: The objective of golf is to play a set of holes in the least number of strokes. A round of golf typically consists of 18 holes. Each hole is played once in the round on a standard golf course. Each stroke is counted as one point, and the total number of strokes is used to determine the winner of the game.

Explain and Answer:

OUTPUT: \_\_\_\_\_

No, the objective of golf is not to get a higher point total than others. Rather, the objective is to play a set of holes in the least number of strokes. The total number of strokes is used to determine the winner of the game, not the total number of points.

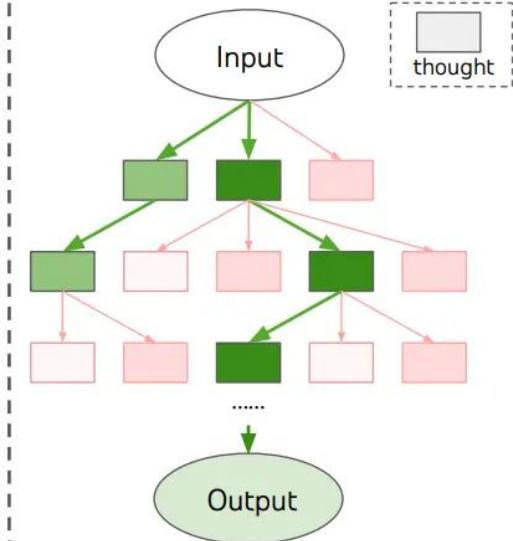
# Tree of Thoughts (ToT)



(a) Input-Output  
Prompting (IO)

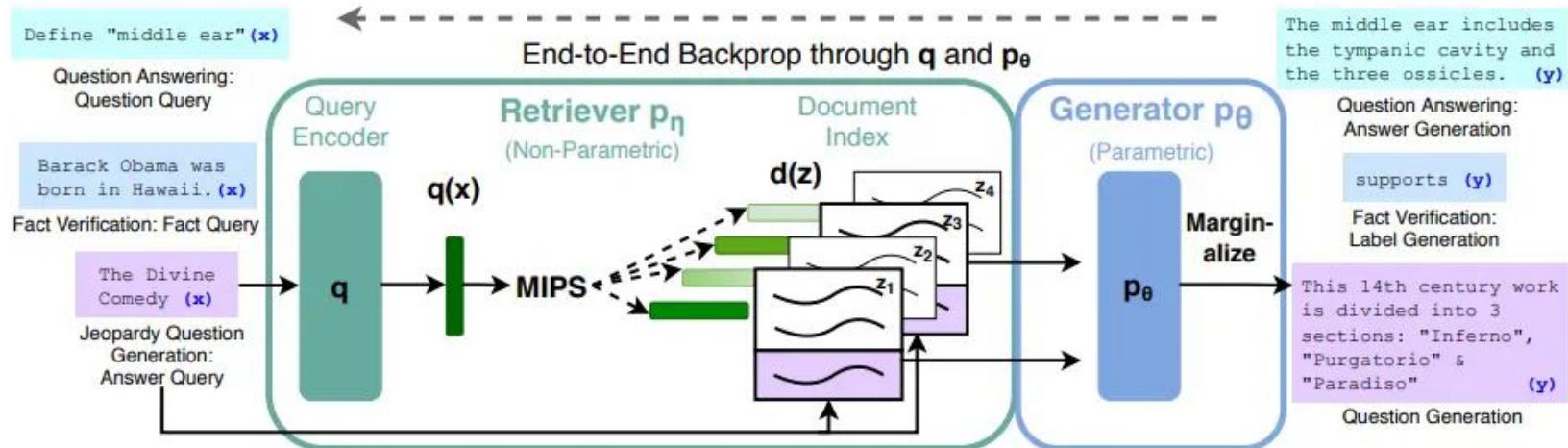
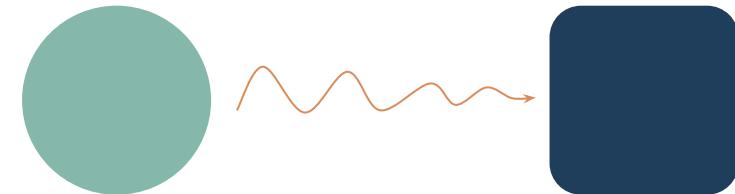
(c) Chain of Thought  
Prompting (CoT)

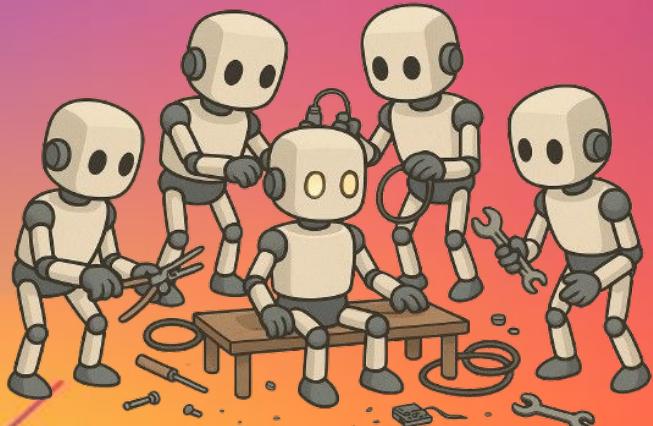
(c) Self Consistency  
with CoT (CoT-SC)



(d) Tree of Thoughts (ToT)

# Retrieval Augmented Generation (RAG)





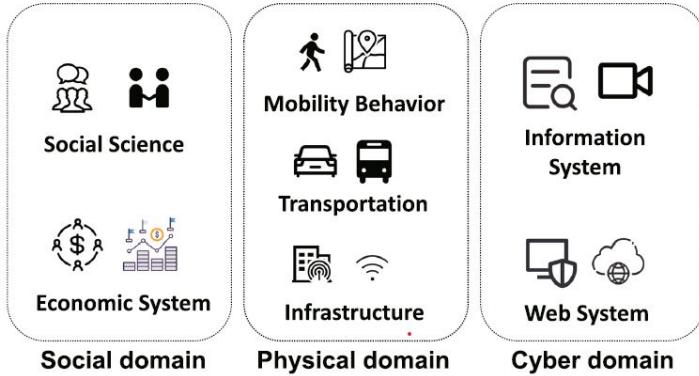
## Content

- From LLMs to AI Agents
- Anatomy of an AI Agent
- Function calling & MCPs
- Building an AI Agent

# AI Agents

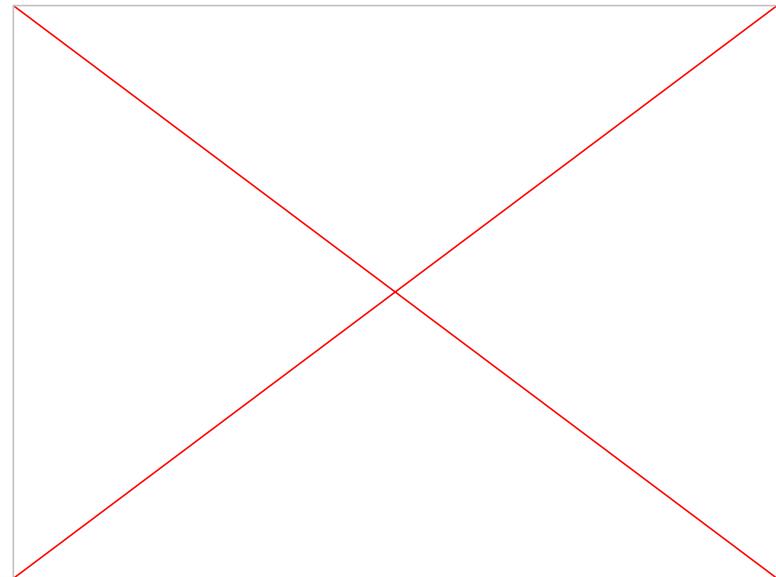
# What can one do with “Agents”?

- Simulation of individual behavior
- Simulation of collective behavior



Gao, C., Lan, X., Li, N. et al. Large language models empowered agent-based modeling and simulation: a survey and perspectives. Humanit Soc Sci Commun 11, 1259 (2024). <https://doi.org/10.1057/s41599-024-03611-3>

<https://www.youtube.com/watch?v=ZdoU9vI2yCg>



Joon Sung Park et al. Generative Agents: Interactive Simulacra of Human Behavior. arXiv:2304.03442v2

# What can one do with “AI Agents”?

Trending AI Agent applications...



Coding



Custom Agentic Workflows



Note Taking

# What can one do with “AI Agents”?



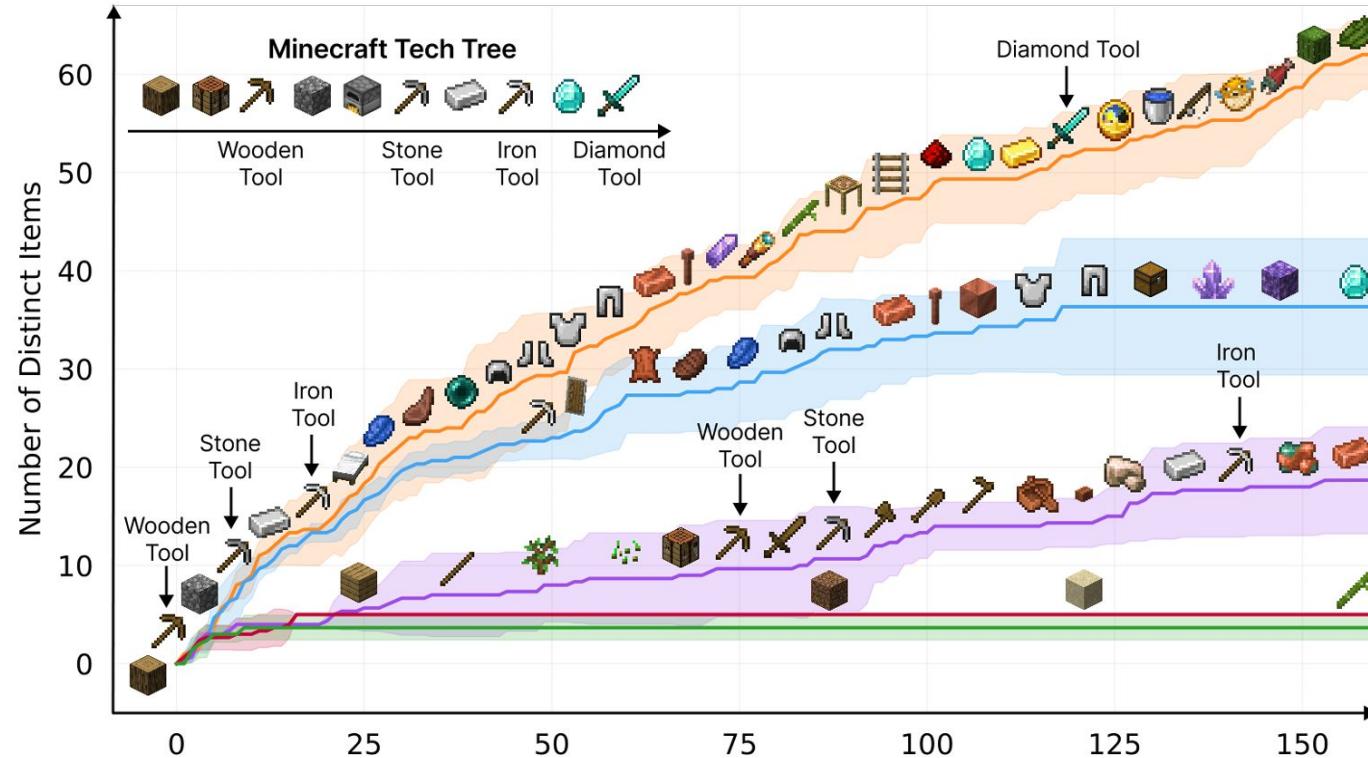
<https://www.youtube.com/watch?v=ZdoU9vl2yCg>

# What can one do with “AI Agents”?



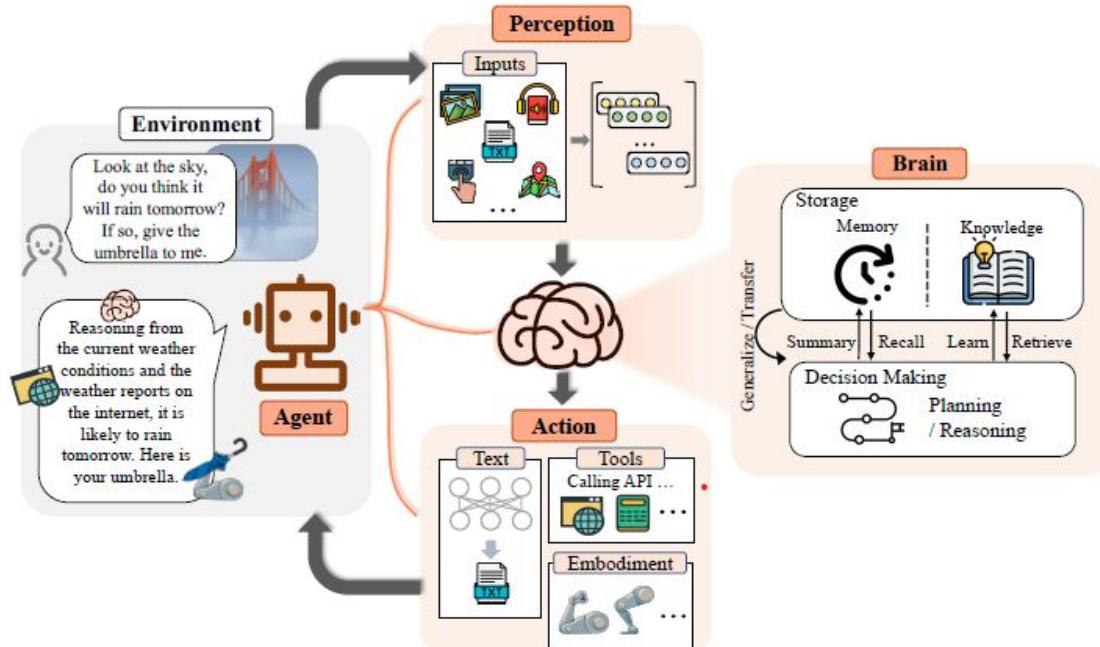
# What can one do with “AI Agents”?

Guanzhi Wang, et al. Voyager: An Open-Ended Embodied Agent with Large Language Models  
<https://arxiv.org/abs/2305.16291>



# Anatomy of an AI Agent (in a nutshell)

1. An agent interacts with an environment
2. It perceives the environment through available modalities
3. Processes information and plans the next step
4. Takes action by executing tools



# From LLMs to AI Agents



## WebGPT: Browser-assisted question-answering with human feedback

Reiichiro Nakano

### REACT: SYNERGIZING REASONING AND ACTING IN LANGUAGE MODELS

Shunyu Yao<sup>\*1</sup>, Jeffrey Zhao<sup>2</sup>, Dian Yu<sup>2</sup>, Nan Du<sup>2</sup>, Izhak Shafran<sup>2</sup>, Karthik Narasimhan<sup>1</sup>, Yuan Cao<sup>2</sup>

<sup>1</sup>Department of Computer Science, Princeton University

<sup>2</sup>Google Research, Brain team

<sup>2</sup>{ jeffreyz}



### Toolformer: Language Models Can Teach Themselves to Use Tools

Timo Schick   Jane Dwivedi-Yu   Roberto Dessì<sup>†</sup>   Roberta Raileanu  
Maria Lomeli   Luke Zettlemoyer   Nicola Cancedda   Thomas Scialom

Meta AI Research   <sup>†</sup>Universitat Pompeu Fabra

# Function Calling

How do we make the LLM to generate **valid** function calls?



What is the weather today in  
Guanajuato?

get\_weather\_api(date="2025-11-03",  
city="Guanajuato")

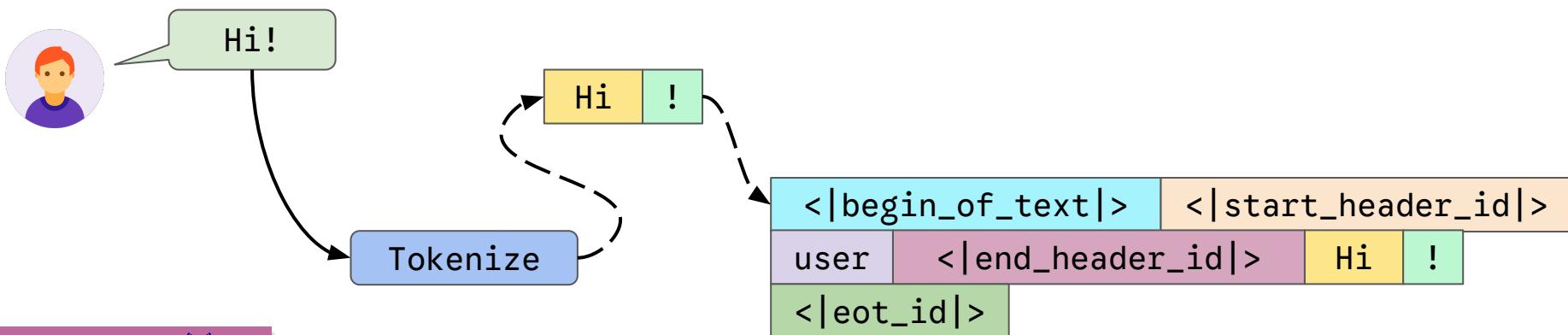


- 1) Teach the language model how to respond by **finetuning** (SFT)
- 2) Constrain the model at decoding-time... grammars, regex, and state machines are techniques used for **controlled generation**

# Function Calling

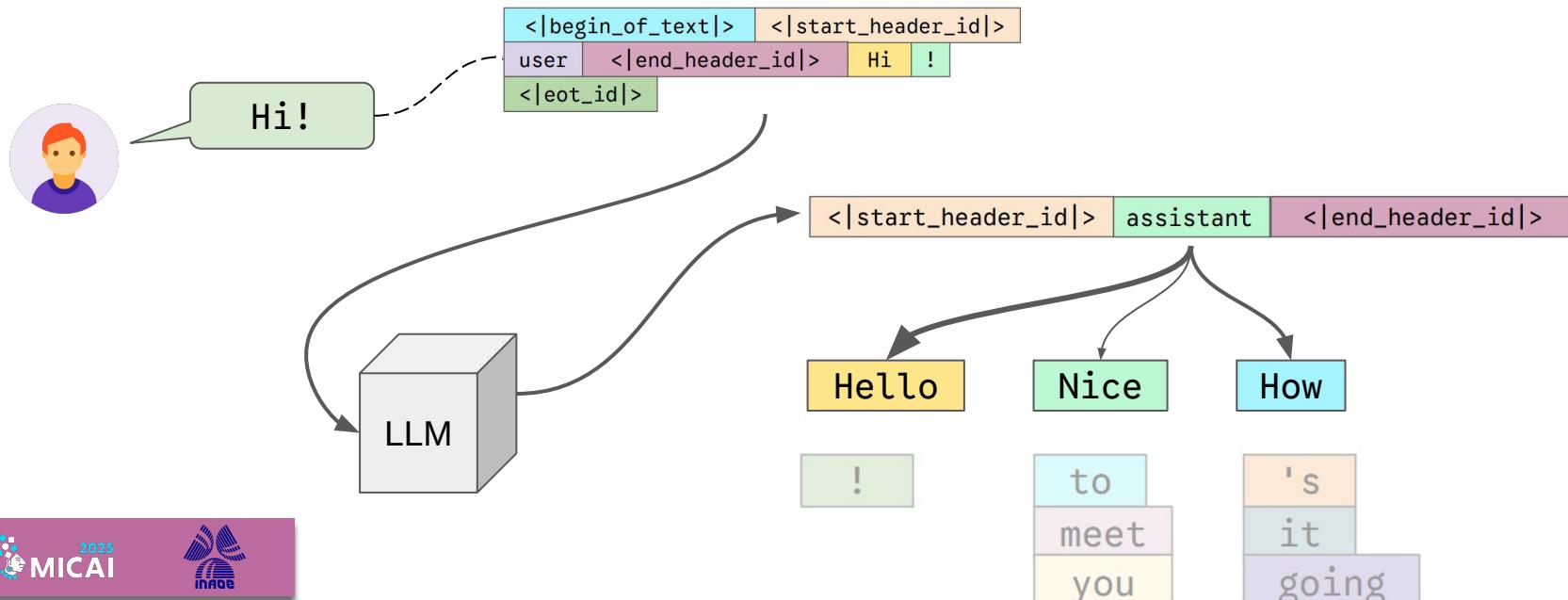
Let's recall a language model processes raw input by first tokenizing.

- Each model has its own vocabulary, and some special tokens are added depending on the use case



# Function Calling

If using **greedy sampling**, the model would choose the most probable token.



# Function Calling

Going back to function calling, **how do we prepare a dataset for SFT?**

1. Leverage system prompt in training dataset to instruct for function calling



You are a helpful assistant with tool calling capabilities.

2. Pass the function JSON



You have access to the following functions...

3. Introduce a new ID when finetuning the model for function calling

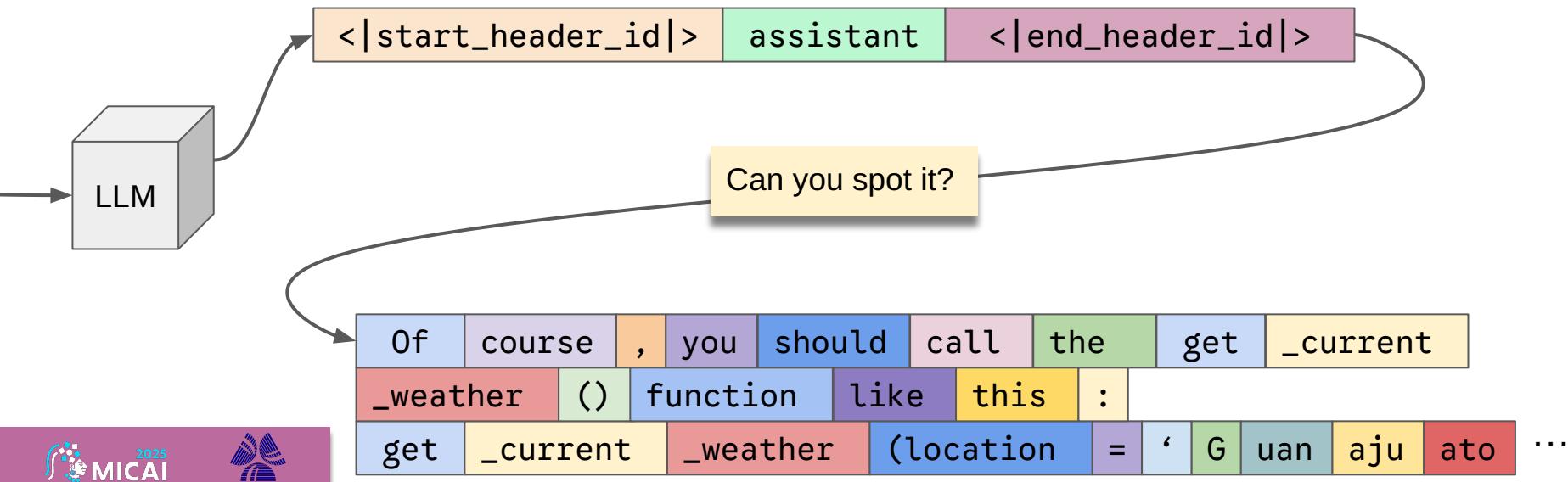
<|start\_header\_id|> tool <|end\_header\_id|>

# Function Calling

```
1  <|begin_of_text|><|start_header_id|>system<|end_header_id|>
2
3  Cutting Knowledge Date: December 2023
4  Today Date: 23 July 2024
5
6  When you receive a tool call response, use the output to format an answer to
   the orginal user question.
7
8  You are a helpful assistant with tool calling capabilities.<|eot_id|>
   <|start_header_id|>user<|end_header_id|>
9
10 Given the following functions, please respond with a JSON for a function call
    with its proper arguments that best answers the given prompt.
11
12 Respond in the format {"name": function name, "parameters": dictionary of
    argument name and its value}. Do not use variables.
13
14 {
15     "type": "function",
16     "function": {
17         "name": "get_current_conditions",
18         "description": "Get the current weather conditions for a specific
           location",
19         "parameters": {
20             "type": "object",
```

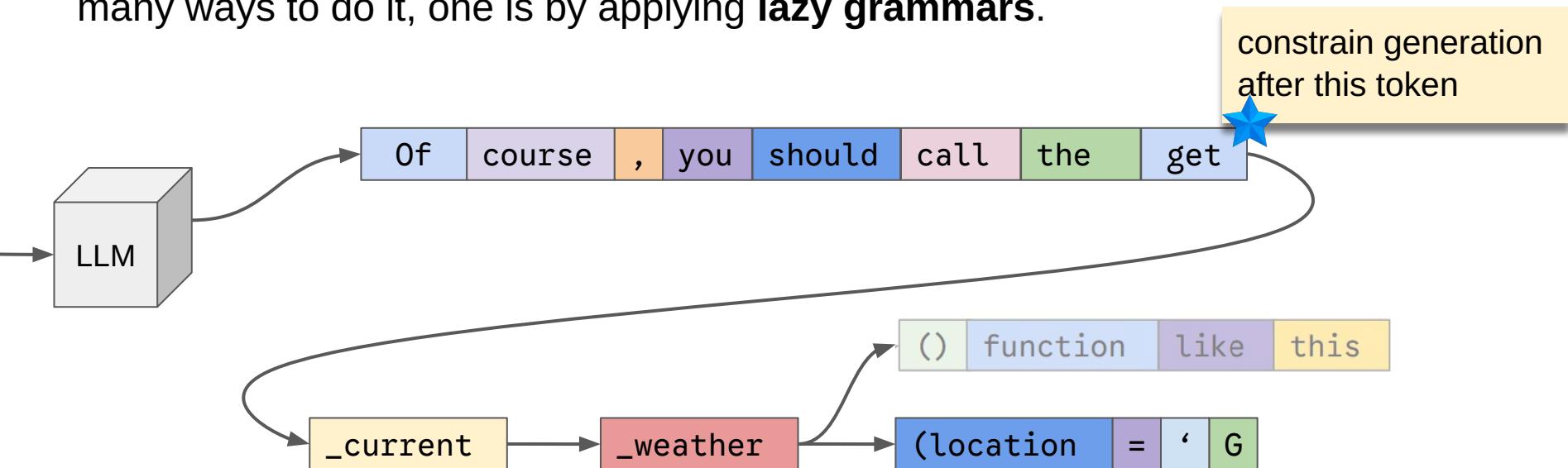
# Function Calling

First approaches for function calling let the model generate freely, then **parsed** the function call from the model's output. However, generation errors make it hard to parse the function.



# Function Calling

One way to avoid this issue is to **control the generation process**. There are many ways to do it, one is by applying **lazy grammars**.



# Function Calling



*Let's try function calling with the OpenAI Python completions module.*

# Model Context Protocol (MCP)

- Define of MCP
- Fundamentals of MCP
- Build MCP Server (later)



# What is MCP?

“An open protocol that standardizes how your LLM applications connect to and work with your tools and data sources” (ANTHROP\C)

# What is MCP (Model Context Protocol)?

## REST APIs

Standardize how **web applications** interact with the **backend**

## LSP

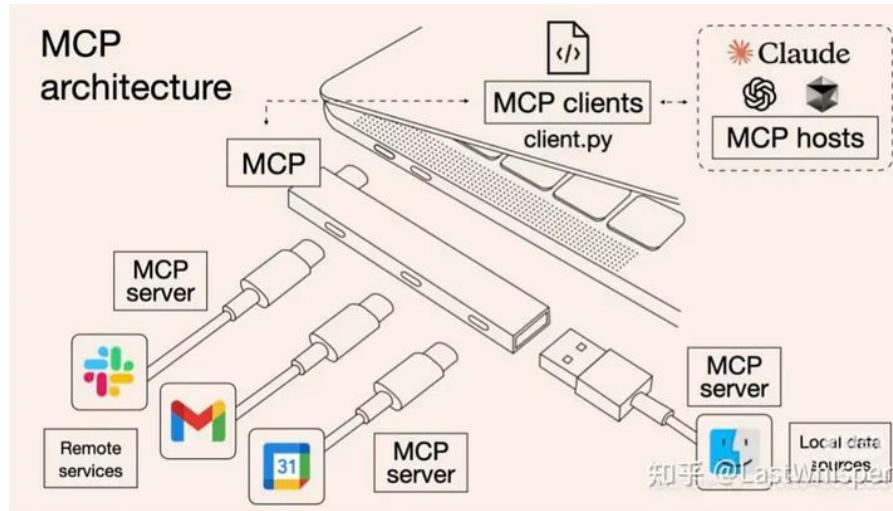
Standardizes how **IDEs** interact with **language-specific tools**

## MCP

Standardizes how **AI applications** interact with **external systems**

# What is MCP (Model Context Protocol)?

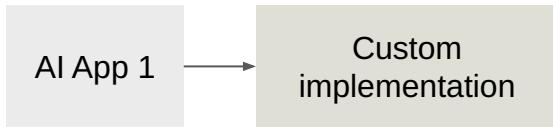
“Think of MCP like a USB-C port for AI applications. Just as USB-C provides a standardized way to connect electronic devices, MCP provides a standardized way to connect AI applications to external systems.”



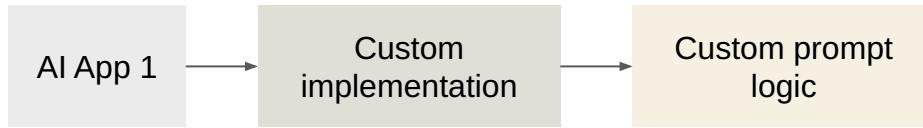
# Without MCP: Fragmented AI Development

AI App 1

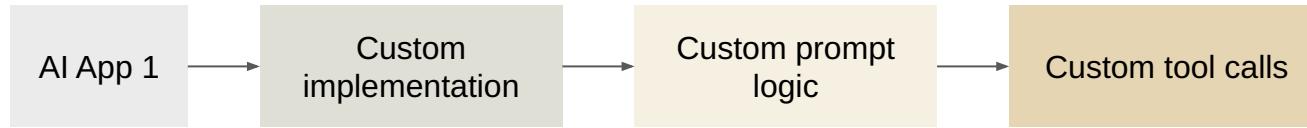
# Without MCP: Fragmented AI Development



# Without MCP: Fragmented AI Development



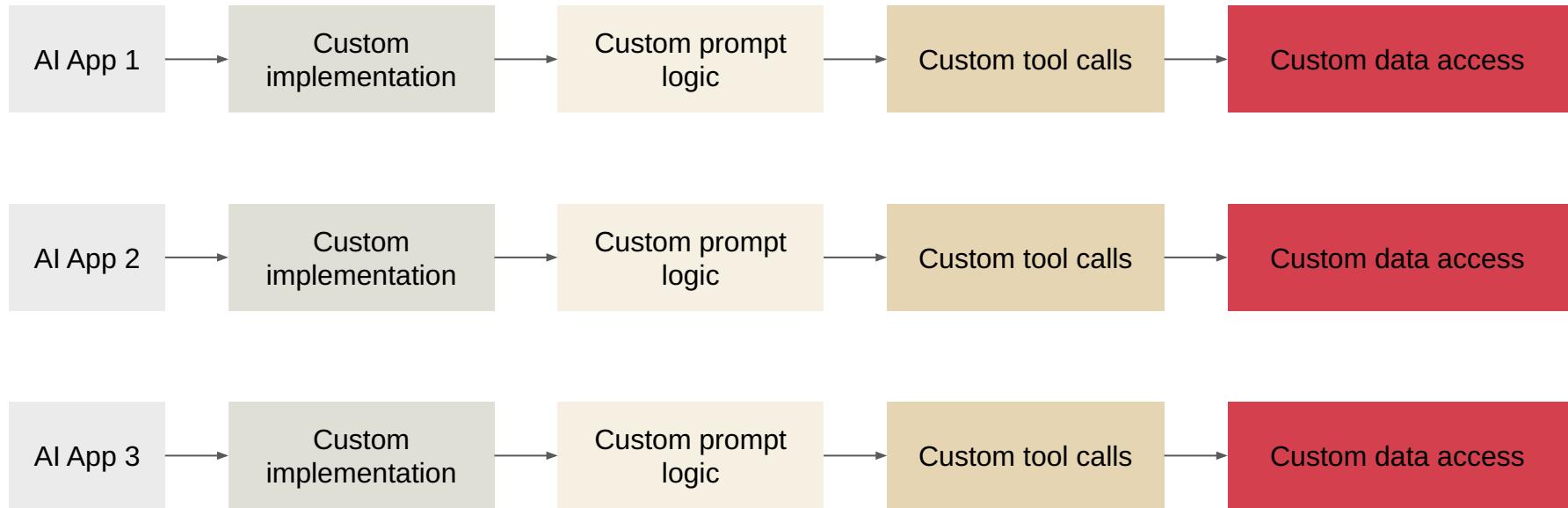
# Without MCP: Fragmented AI Development



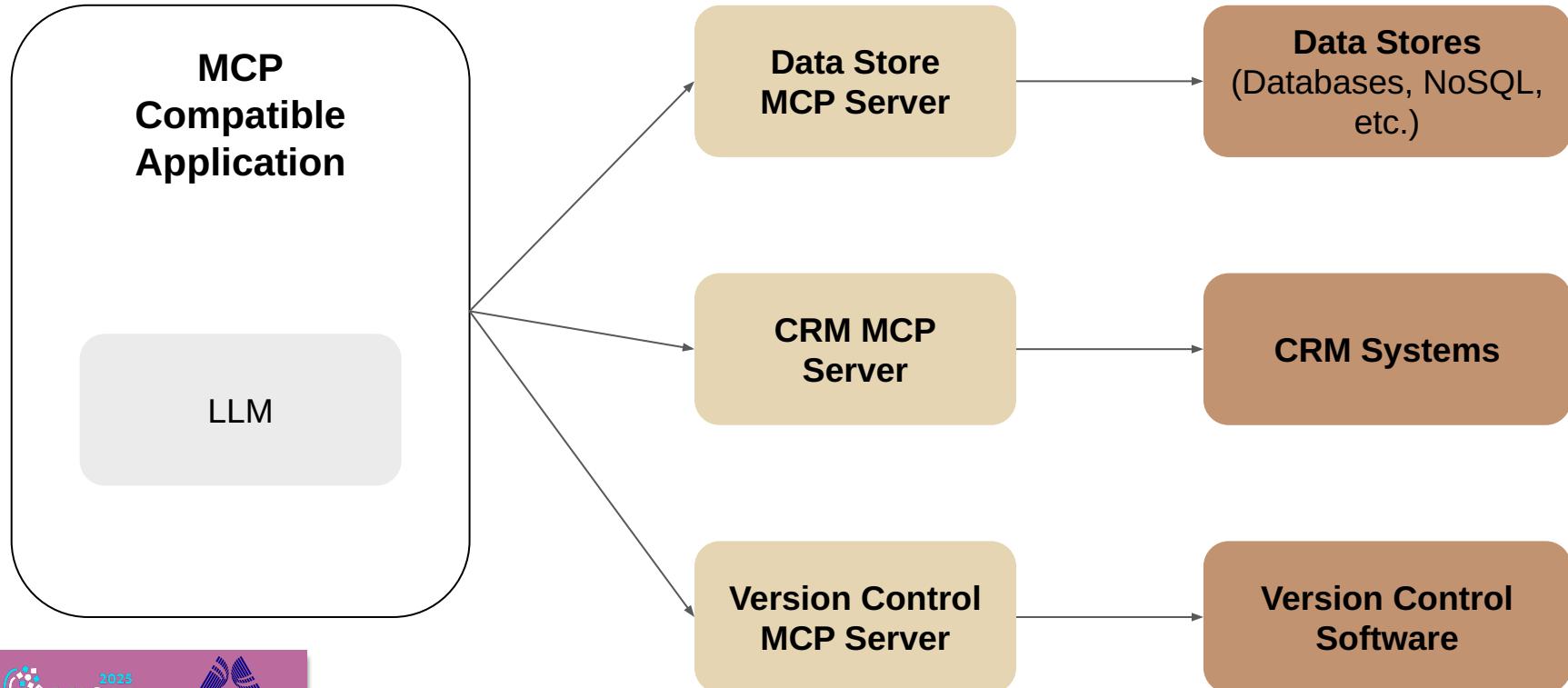
# Without MCP: Fragmented AI Development



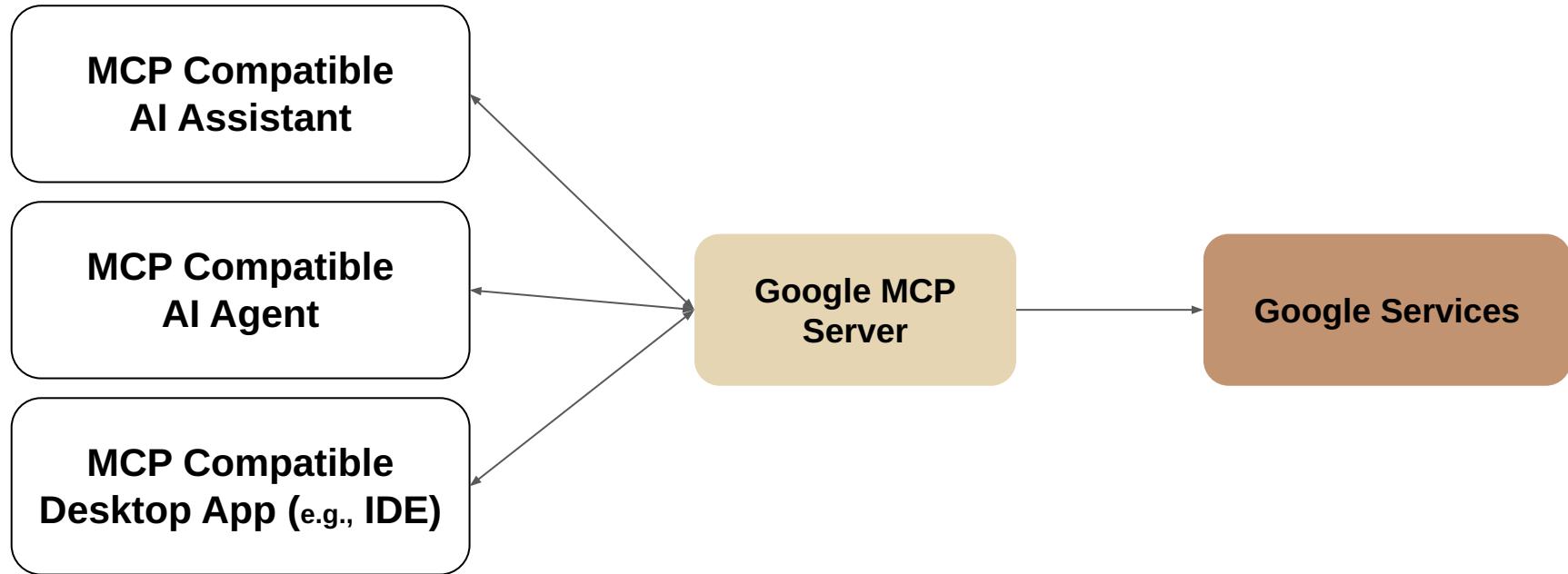
# Without MCP: Fragmented AI Development



# With MCP : Standardized AI Development



# With MCP : MCP Server are reusable by AI-Apps



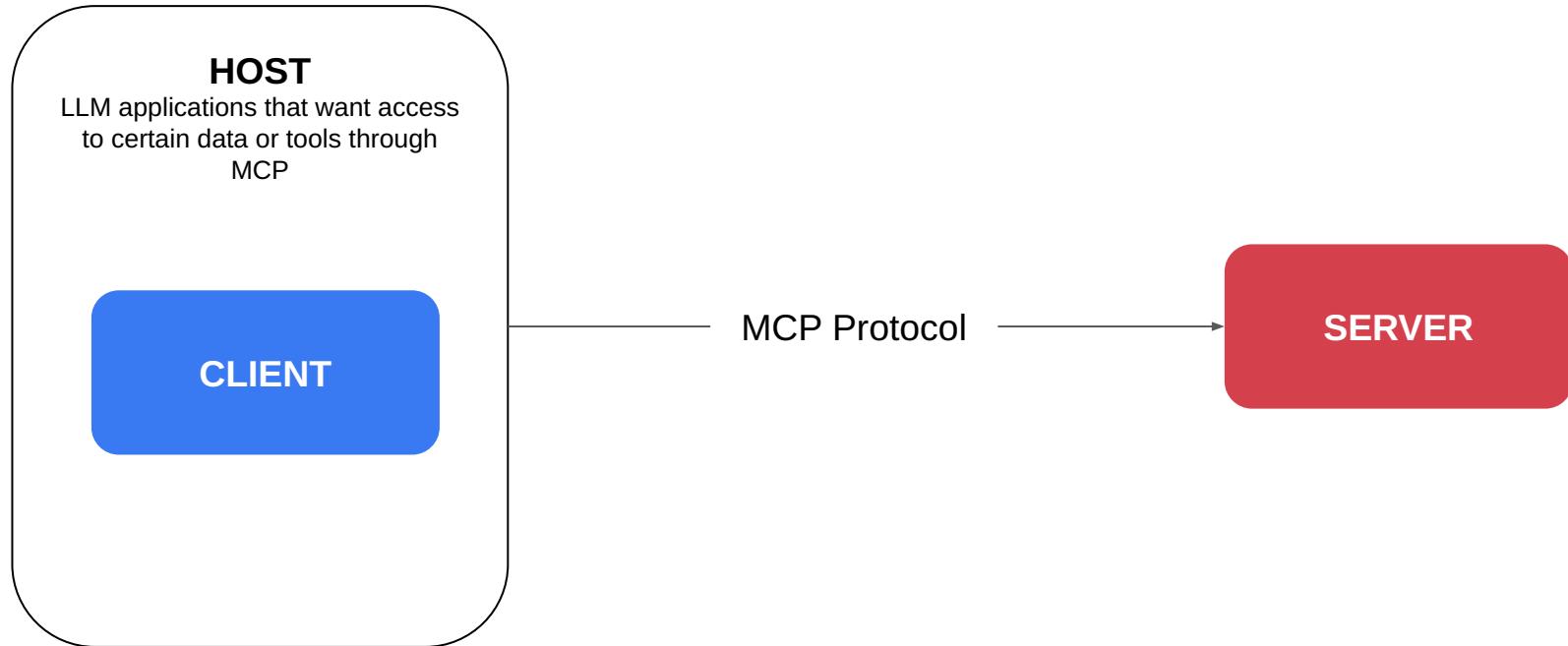
# Client-Server Architecture

HOST

CLIENT

SERVER

# Client-Server Architecture



# Components inside an MCP Server

TOOLS

RESOURCES

PROMPT  
TEMPLATES

# Components inside an MCP Server

## TOOLS

Functions and tools  
that can be invoked  
by the client

## RESOURCES

## PROMPT TEMPLATES

# Components inside an MCP Server

TOOLS

RESOURCES

PROMPT  
TEMPLATES

Read-only data that  
is exposed by the  
server

# Components inside an MCP Server

TOOLS

RESOURCES

PROMPT  
TEMPLATES

Structured prompt  
blueprints

# Communication Life Cycle

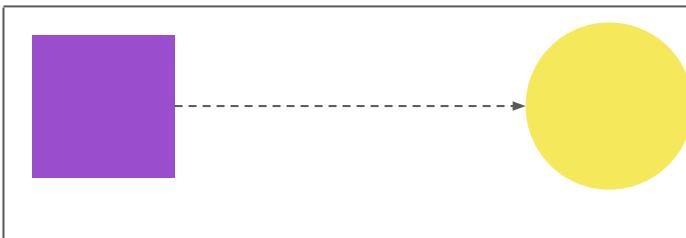
## Initialization



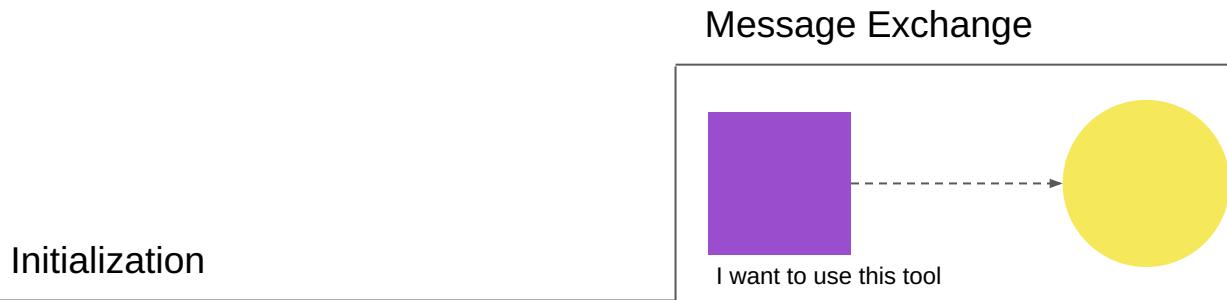
Hey, I would like things from u now

# Communication Life Cycle

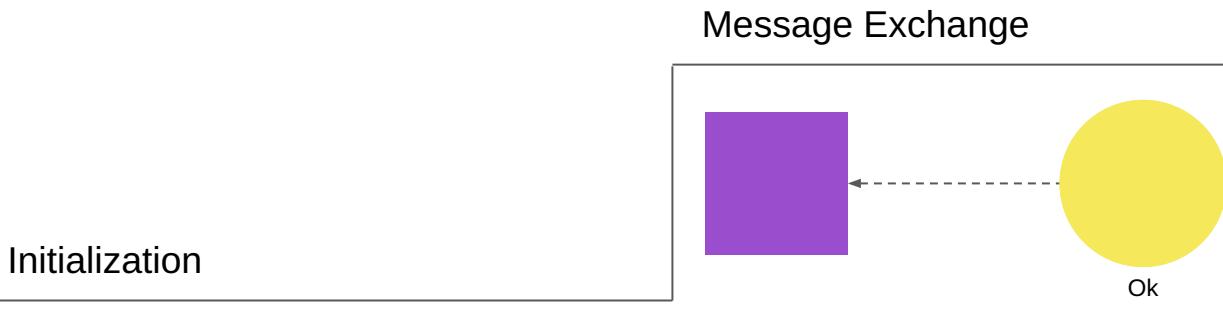
Initialization



# Communication Life Cycle



# Communication Life Cycle



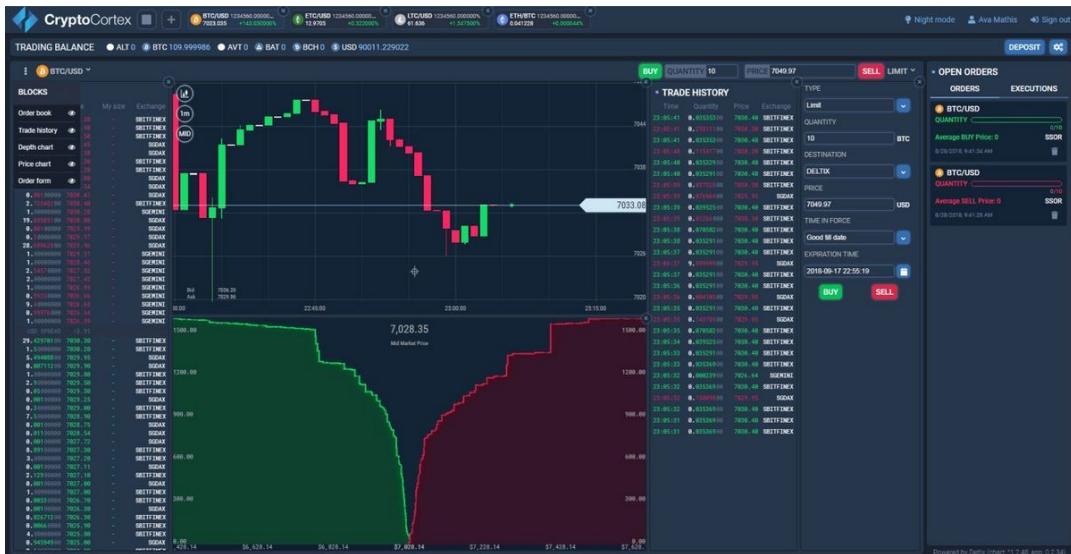
# Communication Life Cycle



# Practice - MCP Servers



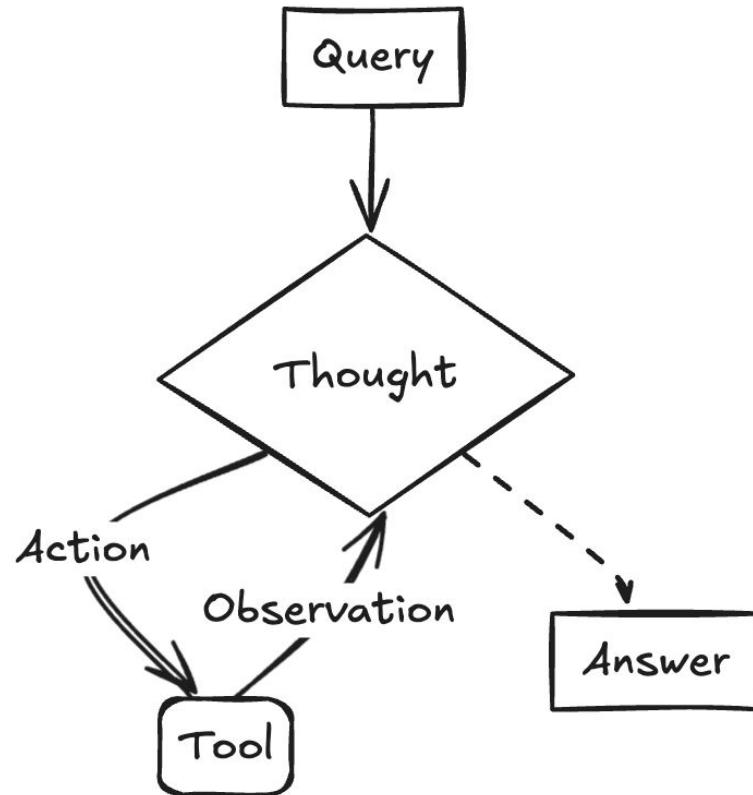
# GitHub



# LLM Agent Architecture

Currently, there are multiple LLM agent architectures...

- The **ReAct** paper is a first reference on how to prompt language models for synergizing *reasoning and acting* general task solving

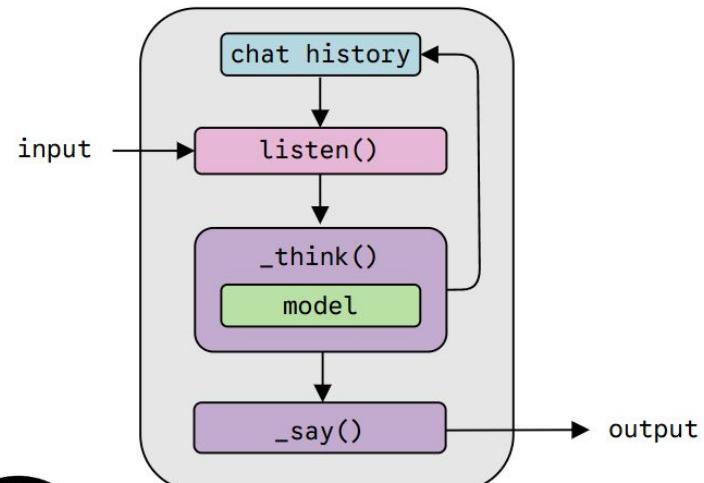


# LLM Agent Architecture



Let's set the foundations of an **agentic loop**

- Set `listen()` method to query the agent
- Generate a completion with the `_think()` method
- Respond back to user with `say()` method
- Store chat history in a list of messages
- `listen()` just packs what's in the chat history and sends it to the `_think()` method

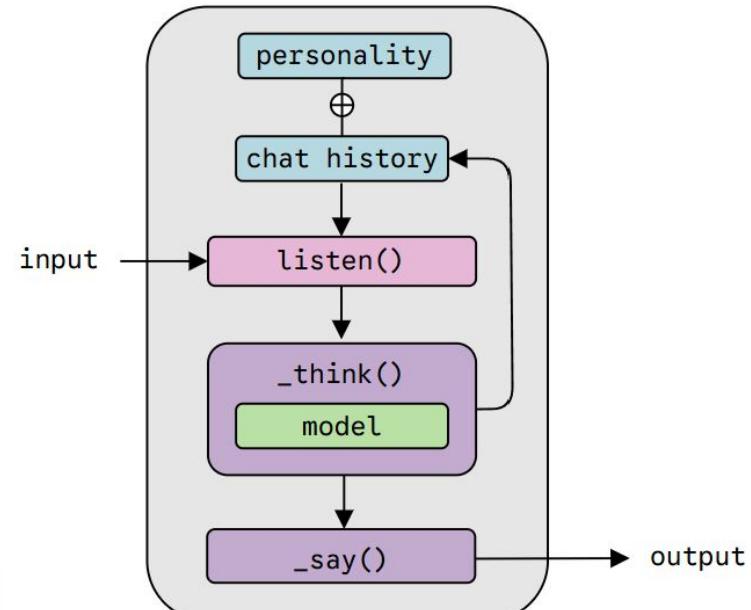
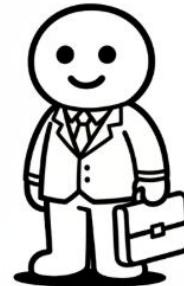


# LLM Agent Architecture

Then, we can use the system prompt to define:

- Personality - for better task specific generations
- Generation guidelines - for setting the objective and indicating how it should respond on every call

Notice how the agent loop stops on every `listen()` call.

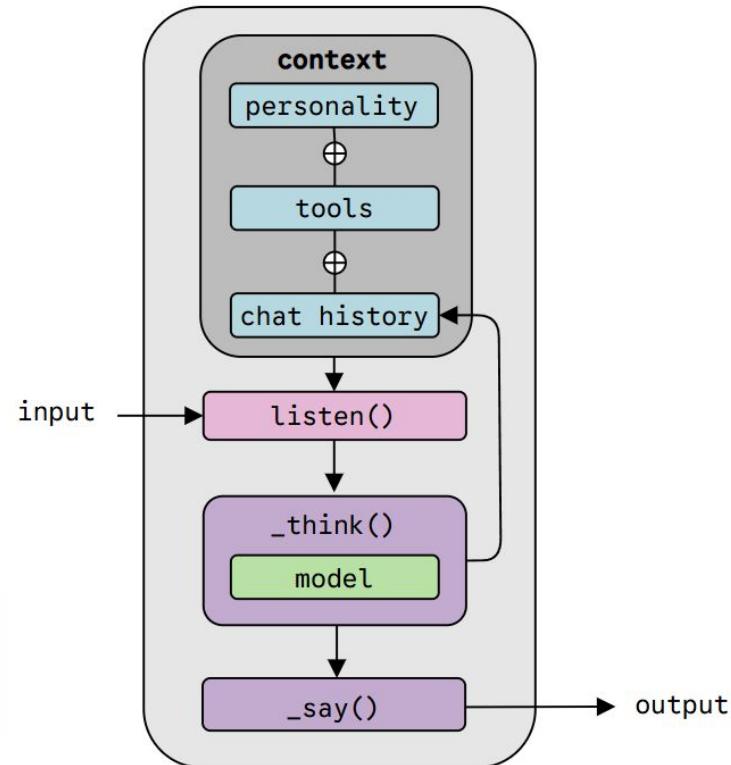


# LLM Agent Architecture

Continue by adding tools for the agent to select, but keep in mind that

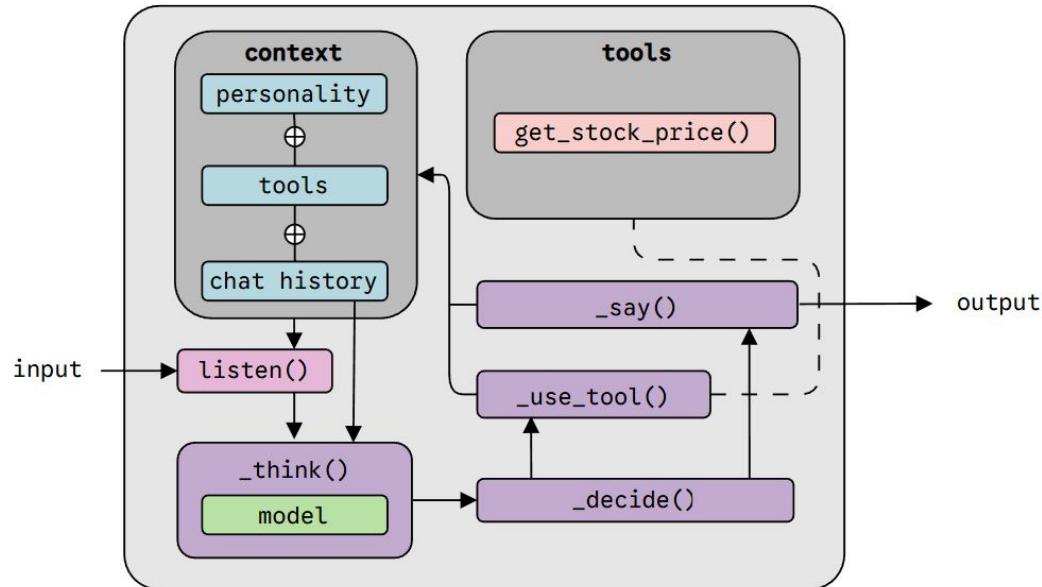
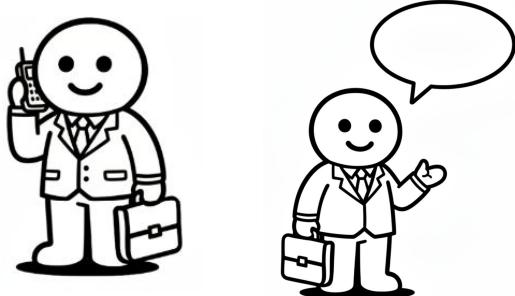
- The tool definition must be included as part of the context
- The model can generate a normal response, or a tool call, or both

We must handle this options by introducing a new block...



# LLM Agent Architecture

- The `_decide()` block allows us to **handle** the next step after a generation





## Content

- From LLMs to AI Agents
- AI Agents vs Agentic AI
- Popular applications

# From LLMs to Agents

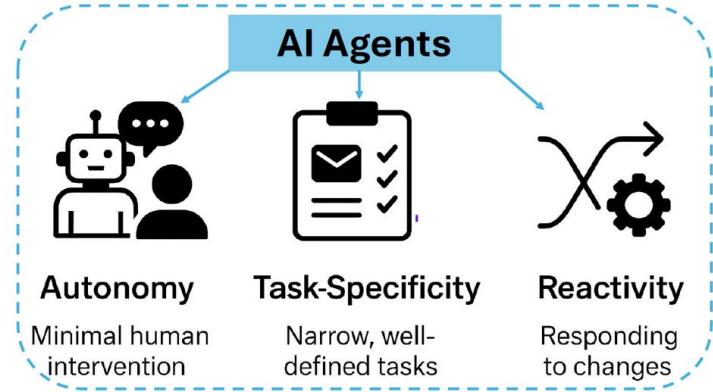
- What LLMs/SLMs can do?
  - “*Understanding, reasoning, generating, etc.*”
- Limitation
  - They are passive, in the sense that they cannot take action by themselves
- AI Agents and Agentic AI
  - Combine LLMs with tools (APIs, memory, reasoning loops, etc.)
  - Agents can perceive, plan and act in the “*world*”
  - Examples: AutoGPT, Devin, ChatGPT with tools



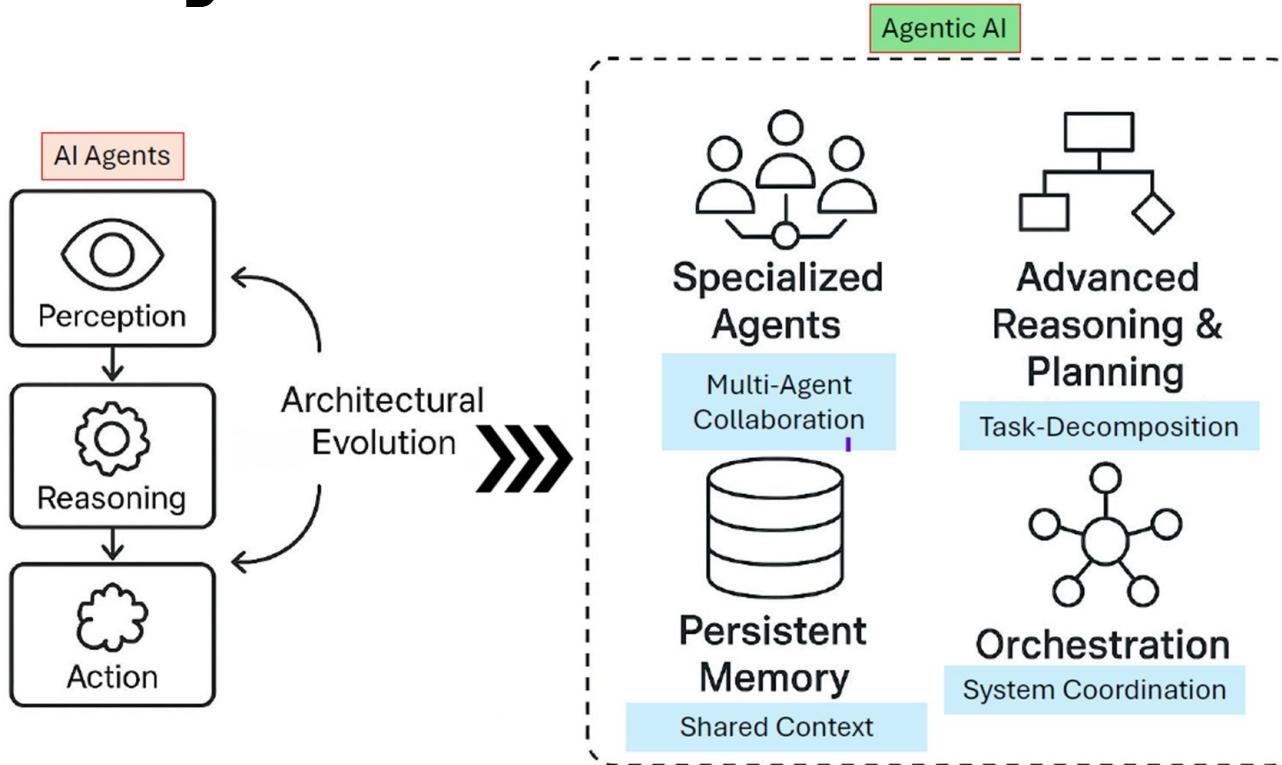
<https://talent500.com/blog/large-language-models-enterprise-ai-adoption/>

# AI Agents vs Agentic AI

- **AIA:** Autonomous (*software*) entities for goal-directed task execution within bounded digital environments
  - Show reactive intelligence and adaptability
- **AAI:** Systems of multiple AI agents collaborating to achieve complex goals through structured communication, shared memory and dynamic role assignment

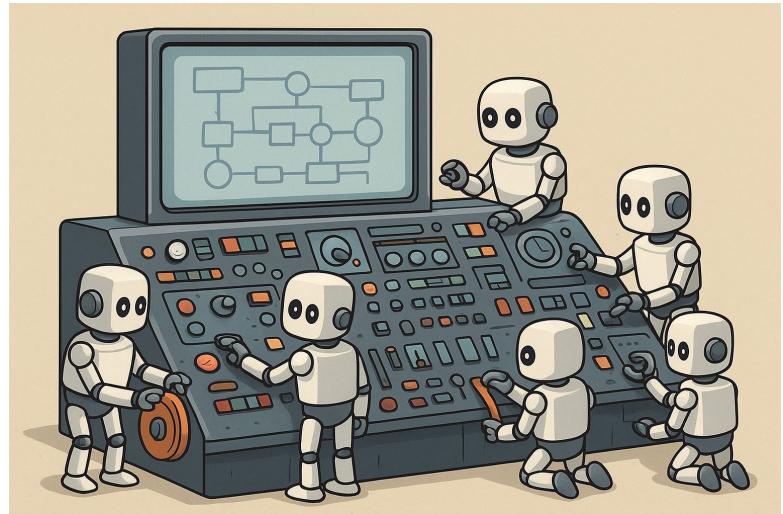


# AI Agents vs Agentic AI



# Applications

- AI Agents
  - Customer support automation
  - Personalized content recommendation
  - Scheduling assistants
- Agentic AI
  - Research assistants
  - Intelligent robotics coordination
  - Multi-agents in gaming

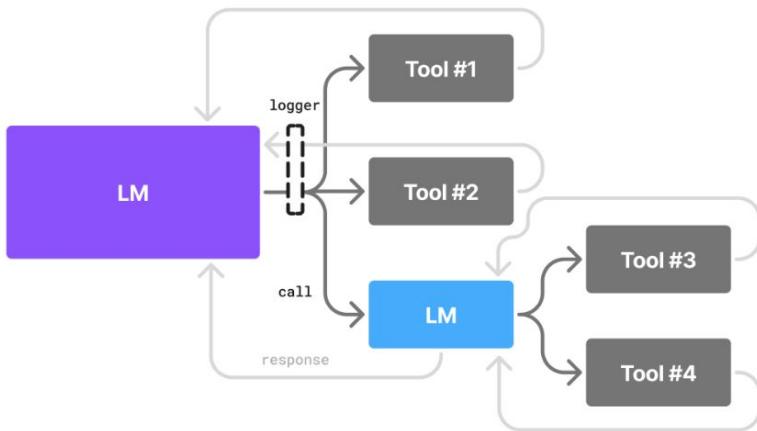


# Are SLMs appropriate for AIA & AAI?

- LMs in AIA and AAI are used for repetitive, simple, scoped and non-conversational tasks (using LLMs could be a waste of LM resources)
- **Advantages of SLMs for AIA and AAI:** lower latency, reduced memory and computational requirements, and significantly lower operational costs
- **Further benefits:** reduced economical cost, democratization of AI



# Building Agents with SLMs



Example Control Flow:



# What can one do with AIA & AAI?



**Idea Generation**

## Content



- Supervisor Architecture
- Hierarchical Architecture
- Network Architecture

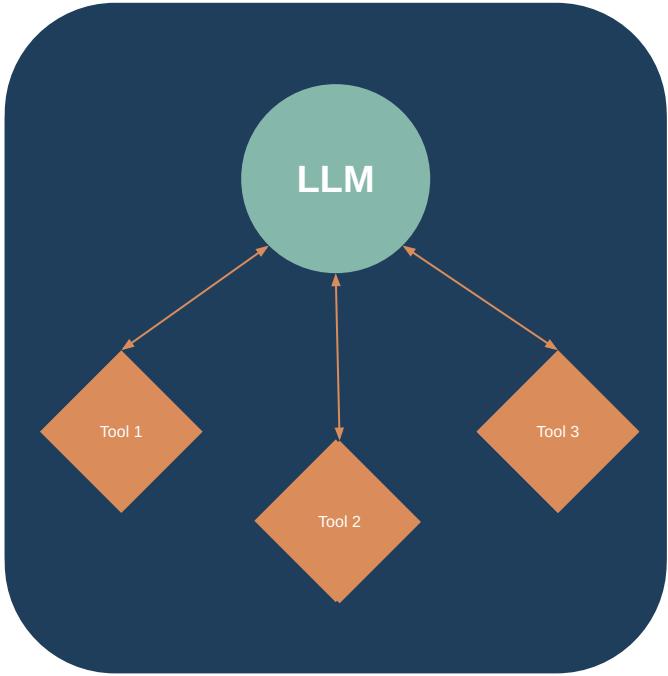
# Multi-Agent Systems

# Multi-Agent Systems (MAS)



- Single Agents
- Supervision
- Hierarchical
- Network
- Custom

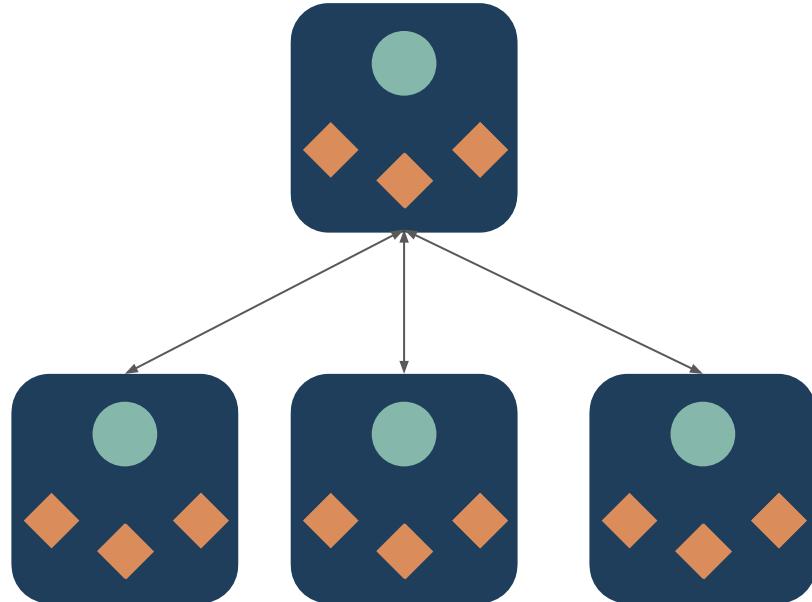
# Single Agents



We only have one agent with a host of tools and connected to an LLM.

The agent then takes in a query and makes decisions on what tool to call using the LLM for reasoning.

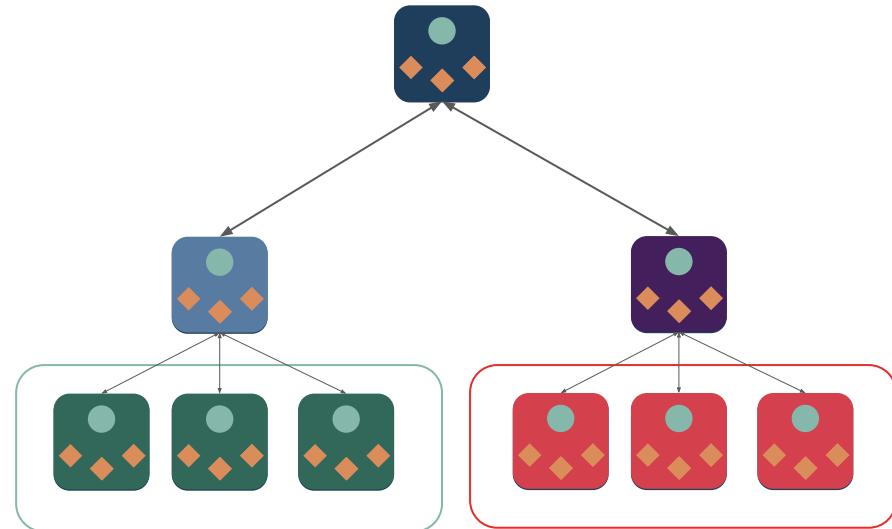
# Supervisor Agent Architecture



The supervisor, basically, you can think of it as the orchestrator in a church choir.

It takes in queries, uses an LLM to reason over it and decide on what agents to delegate a task to or a sub-task to.

# Hierarchical Agent Architecture

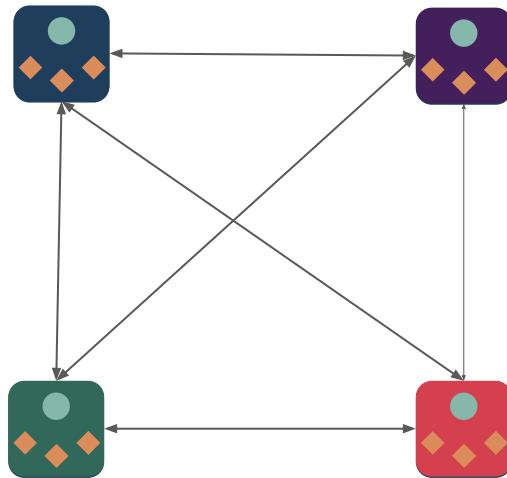


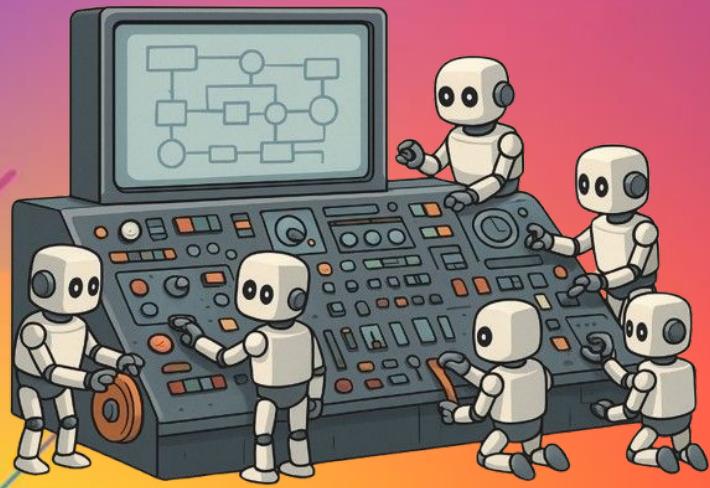
One common approach is to create teams, where each team has its own supervisor, that supervisor is specialized in managing its own team of sub-agents.

Basically breaking the graph into subgraphs, each sub-graph has its supervisor and sub-agents under that supervisor.

# Network Agent Architecture

Each agent can talk to any other agent in the whole collection.



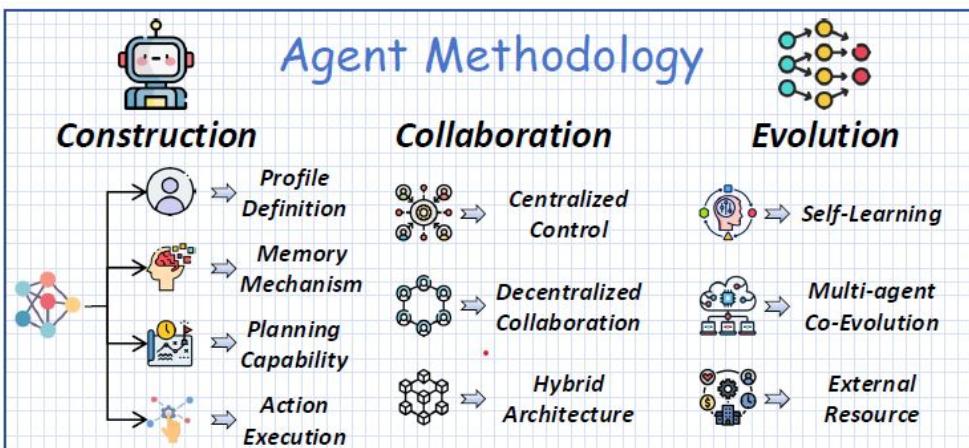


# AI Agents using Small Language Models



# Agent Brain Decomposition

1. An agent interacts with an environment
2. **Perceiving** through any modality
3. Processes information and **plans** the next step
4. Takes **action** by executing tools



Junyu Luo et al. Large Language Model Agent: A Survey on Methodology, Applications and Challenges.  
<https://arxiv.org/abs/2503.21460>