

Introduction to Large Language Models (LLMs)

Training and Adaptations

- Priyam Ghosh
(Upcoming AI Research Intern
@Nvidia)
(Ex JPM @NASA)

HISTORY OF ARTIFICIAL INTELLIGENCE



Alan Turing



John McCarthy

How LLMs Learn??



Behind the magic, LLMs are just glorified matrix multiplication machines.

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \end{bmatrix} = \begin{bmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{bmatrix}$$

$$c_{11} = a_{11} * b_{11} + a_{12} * b_{21} + a_{13} * b_{31}$$

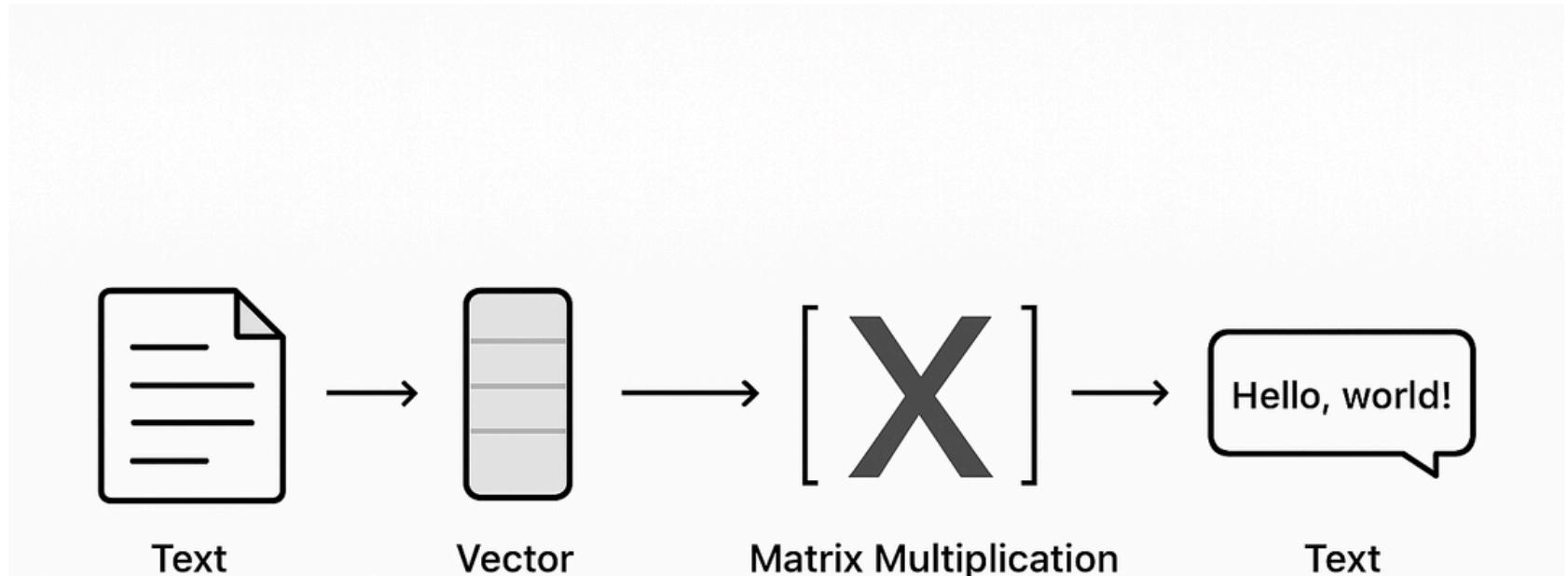
$$c_{12} = a_{11} * b_{12} + a_{12} * b_{22} + a_{13} * b_{32}$$

$$c_{21} = a_{21} * b_{11} + a_{22} * b_{21} + a_{23} * b_{31}$$

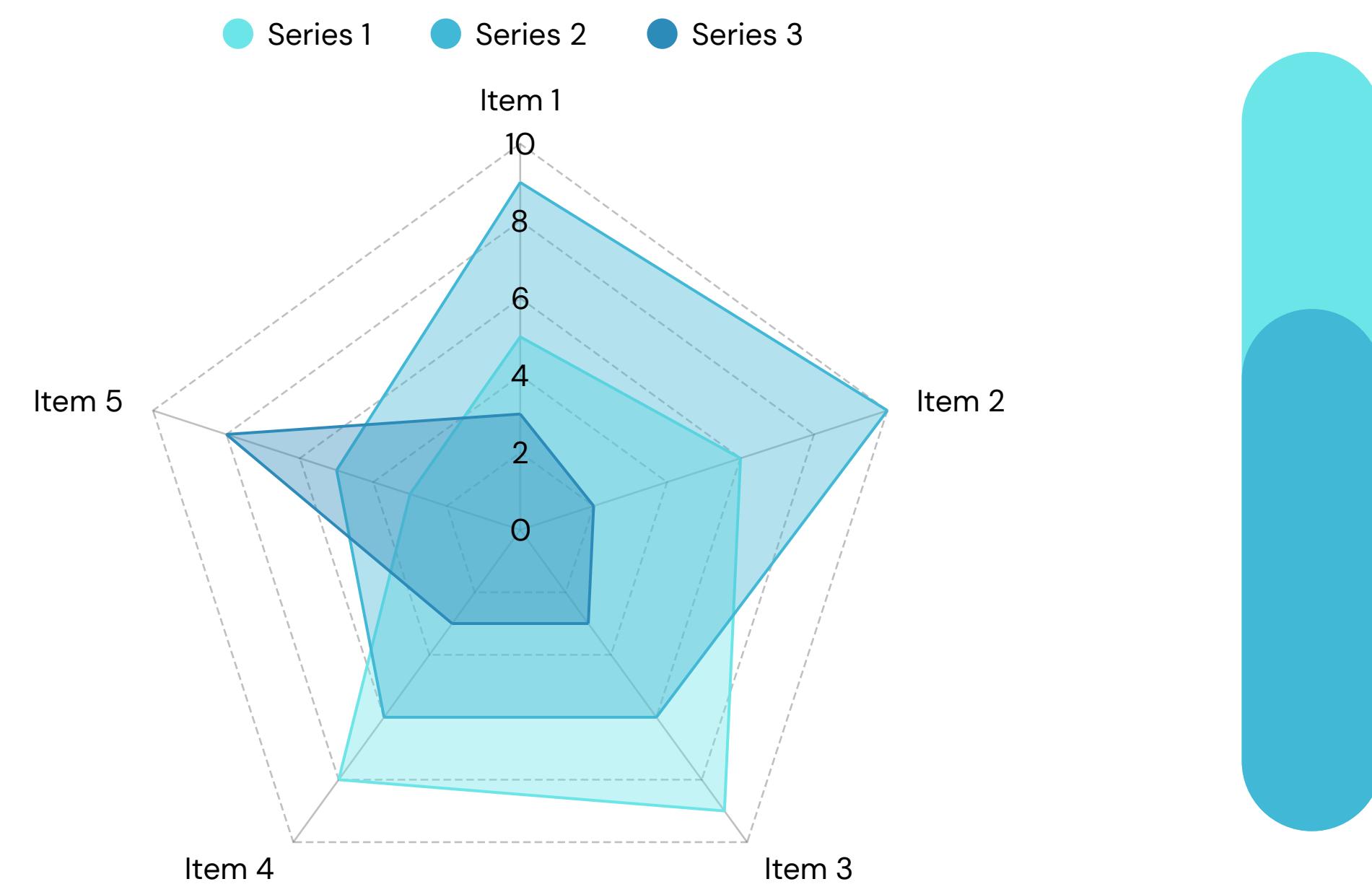
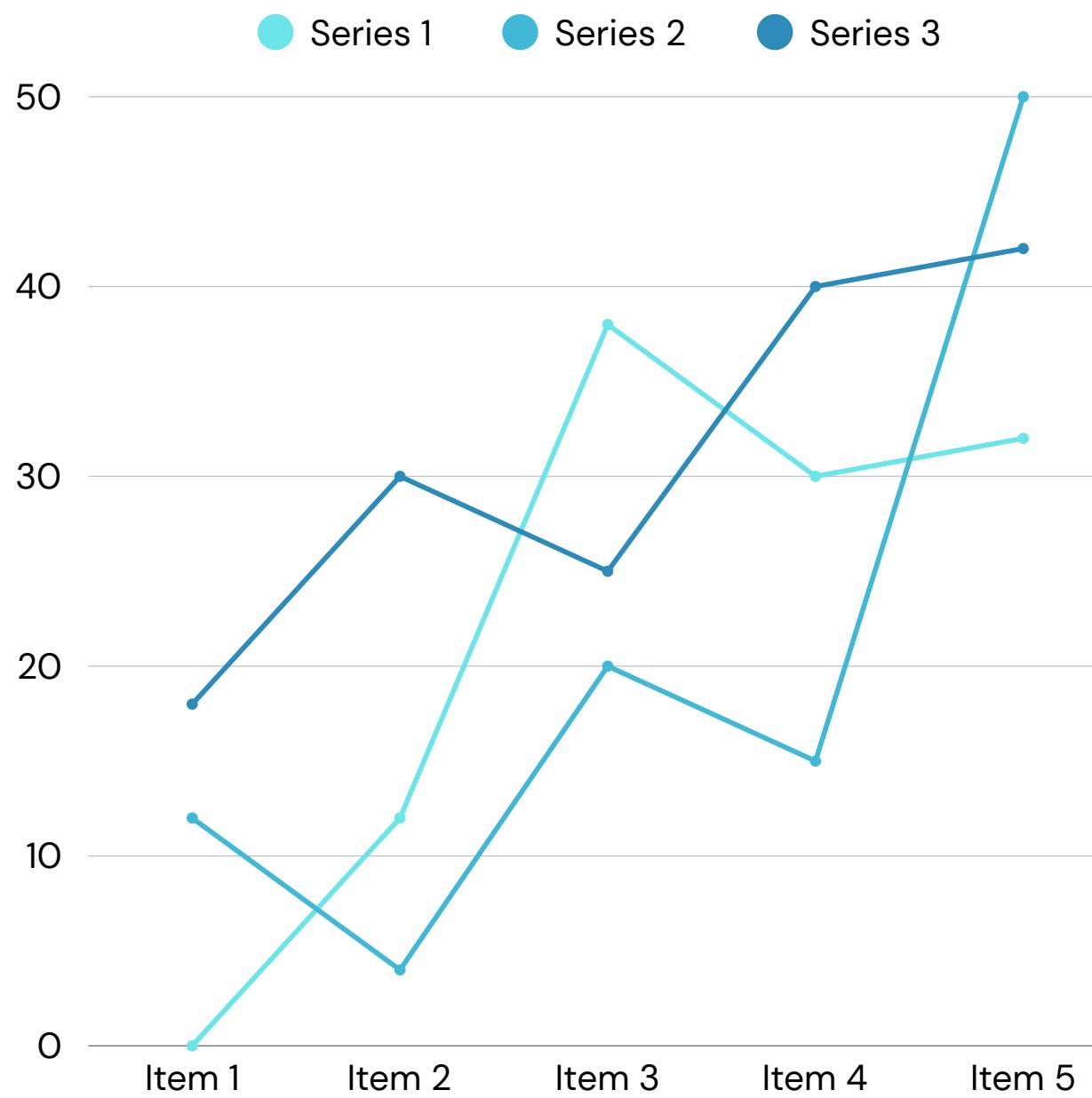
$$c_{22} = a_{21} * b_{12} + a_{22} * b_{22} + a_{23} * b_{32}$$

"Input Embedding" → "Weight Matrices" → "Output Embedding"

- Every word, sentence, and meaning is converted into **numbers (vectors)**.
- These numbers are **multiplied** by massive **matrices** (learned during training).
- The output is new numbers, which are **decoded** back into **words**.
- At their core, LLMs are huge chains of **matrix multiplications — billions** of times!



Data: The Best Friend of Large Language Models (LLMs)



Training Process: Predicting the Next Word

1

Core Objective

LLMs learn by predicting the next word in a sequence, such as guessing "jumps" after "The quick brown fox."

2

Input and Output

The model receives a sequence of words as input and outputs probabilities over a large vocabulary representing possible next words.

3

Learning Through Prediction

By continually refining predictions against actual next words during training, the model develops deep contextual understanding of language.



The "Tireless Musician" Analogy

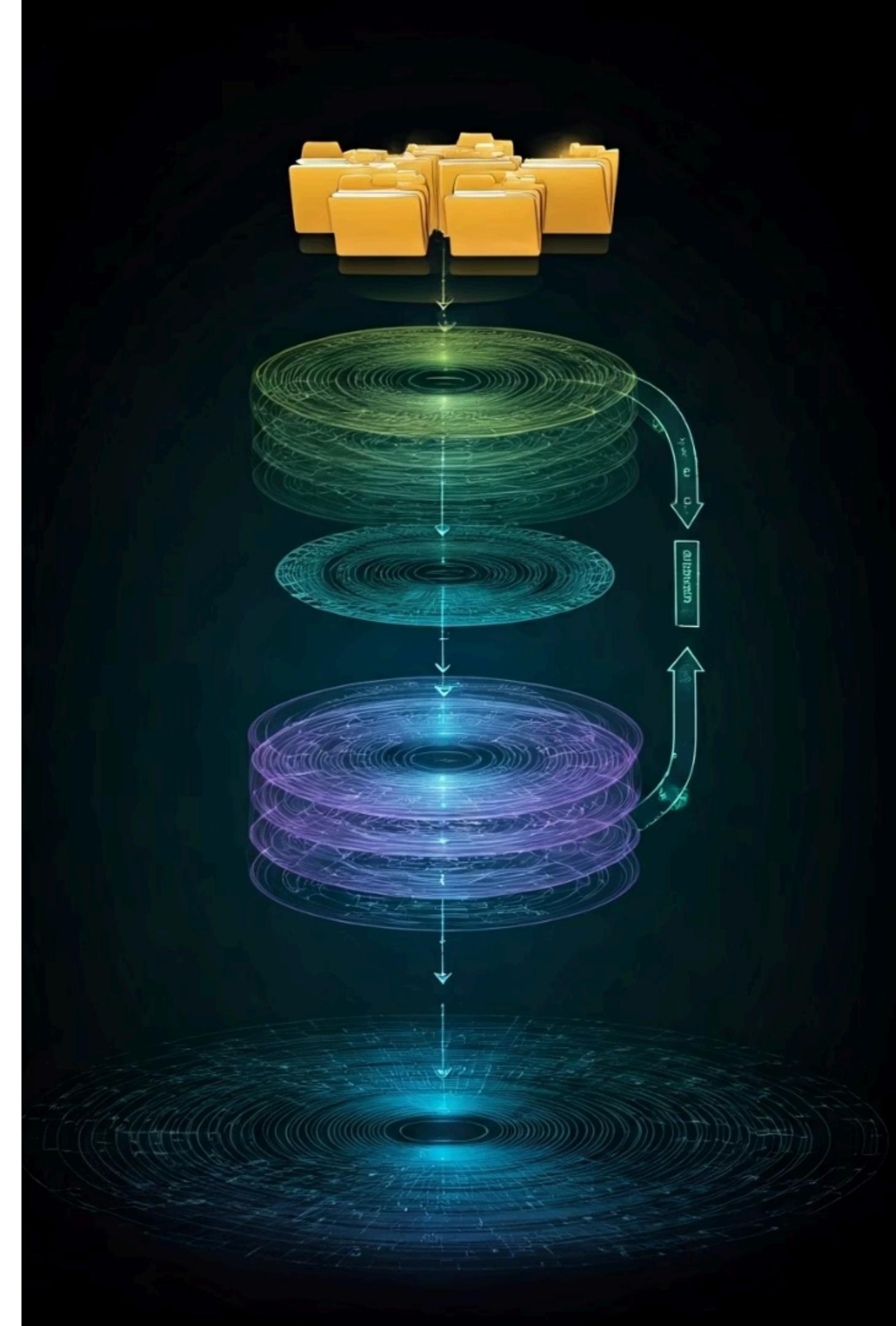


Ujjwal is ____?

1. sleeping
2. dead
3. ਲੈਂਡਿਆਬਾਜ਼

- 1. Data Feasting :** student reading the book and understanding stuffs
- 2. Self-supervised Predictions** trying to solve a puzzle game with predictions and internal error solving
- 3. Model Architecture & Forward Pass** like an assembly line going from one layer to other trying to create the best possible output
- 4. Learning Loop** redoing each step and in every step it learns more and more and repeats it again.
- 5. Scale and Repeat**

The Training Loop: Forward Pass, Loss, Backpropagation



Fine-Tuning LLMs: What It Is and Why It's Needed



Pre trained Model

Domain Specific Fine tuning

- Updated Accuracy
- Customization
- Data Privacy
- Handling Rare Scenarios

LLM Tuning Methods

deci.

Prompt
Engineering

PEFT

RAG

Full
Fine-tuning



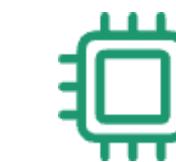
Fine-tuning and Adaptation



Full Fine-Tuning
Complete retraining
of all model
parameters on
specific task data to
tailor performance.



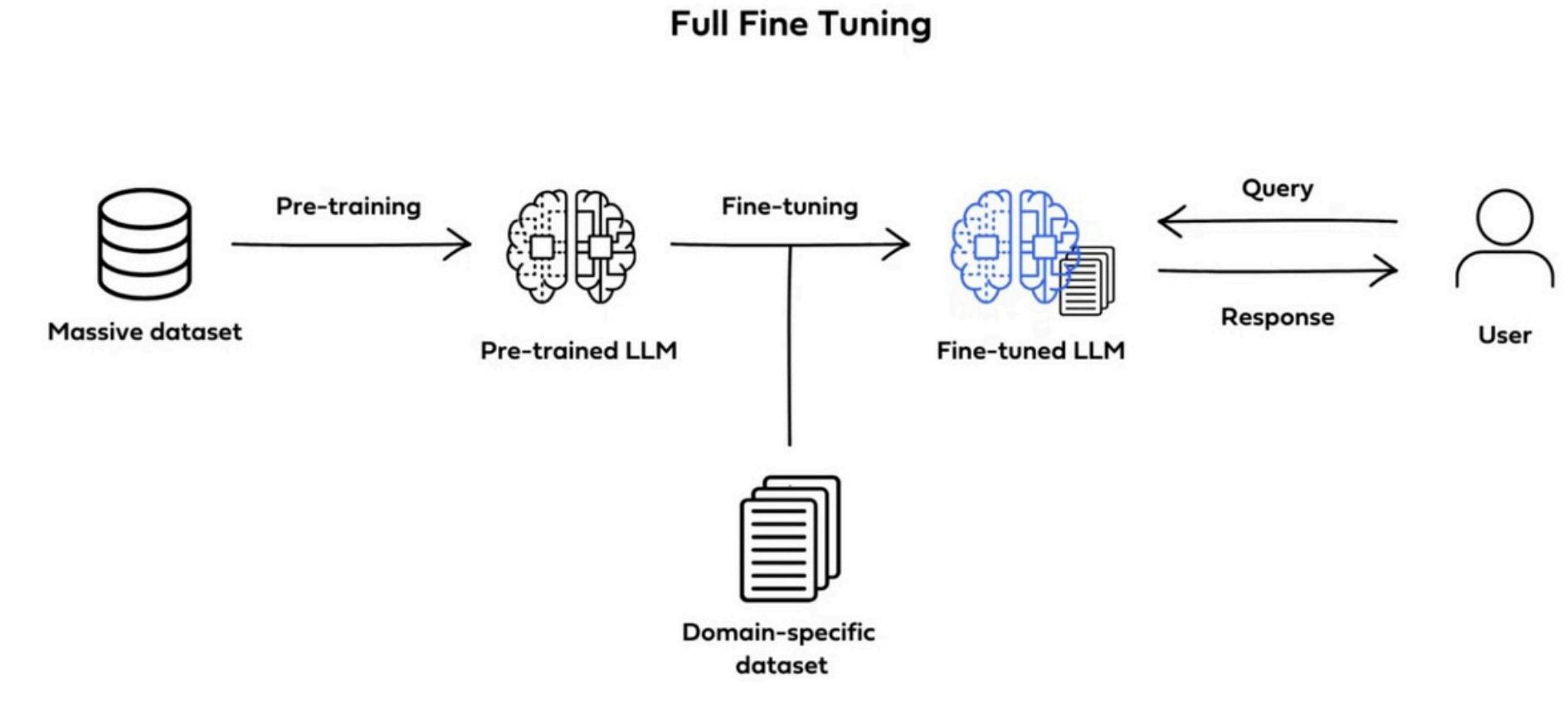
**Prompting /
Few-Shot
Learning**
Using clever input
prompts with few
examples to steer
model outputs
without retraining.



**Parameter-
Efficient Fine-
Tuning (LoRA)**
Updating only a small
subset of parameters
to adapt models
efficiently with lower
resource
requirements.

Full Fine tuning

- Updates all model parameters for deep adaptation
- Needs large datasets and high computational resources
- Delivers maximum accuracy for specialized tasks
- Risk of overfitting and losing general knowledge
- Best for mission-critical, highly customized applications



Model parameters are the internal variables of a machine learning model that are learned from data during training.

deci.^o

Prompt Tuning

- Trains small, learnable prompt embeddings
- Model weights stay frozen
- Persistent effect (task-based)
- Efficient, reusable task adaptation

Soft Prompts: These are artificial, learnable tokens (not human-readable) added to the beginning of the input sequence. They are optimized during training to guide the model toward better performance on a specific task

Analogy: Intelligence Machine with Passcode

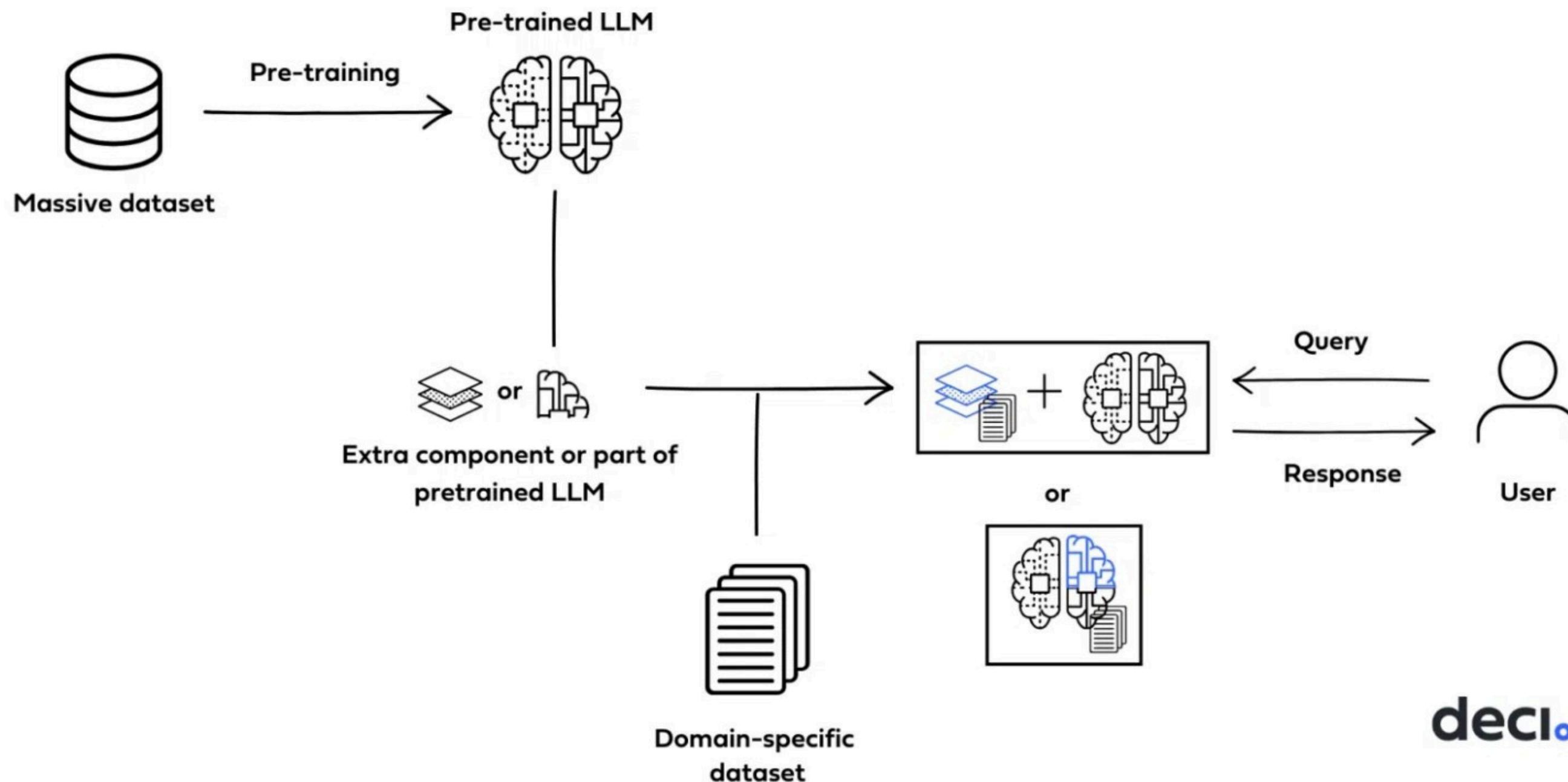


Parameter Efficient Fine Tuning (PEFT)

Low Rank Adaptation (LoRA)



Parameter-Efficient Fine-Tuning

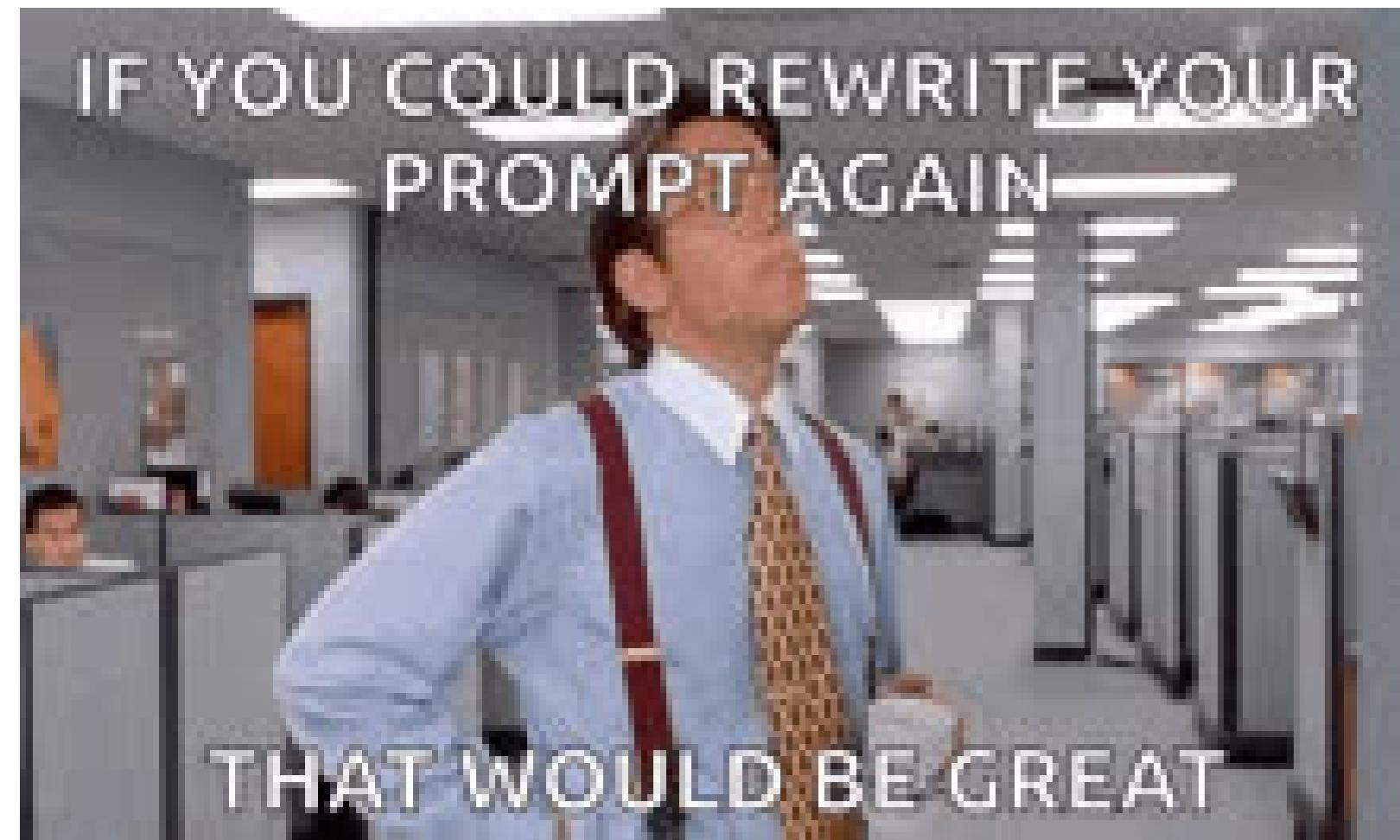


deci.

- Adds small trainable matrices to the model (adapters).
- Keeps original model weights frozen.
- Dramatically reduces memory and compute needs.
- Achieves performance close to full fine-tuning.
- Allows quick, cheap adaptation to new tasks.
- Easy to swap between different task adapters.
- No extra inference cost after merging adapters.
- Great for customizing large models efficiently.

Aspect	Full Fine-Tuning	Prompt/Few-Shot	LoRA Fine-Tuning
		Learning	
What is it?	Retrain all model parameters on new data	Guide the model with examples in the prompt	Update a small subset of parameters using adapter layers
Training Needed?	Yes	No	Yes
Resource Usage	Very high (needs lots of compute & time)	None (just prompt crafting)	Low (efficient, less compute needed)
Accuracy Potential	Highest (best for deep customization)	Limited by model's existing knowledge	High (close to full fine-tuning)
Risk of Forgetting	Yes (may lose general knowledge)	None	Low
Speed to Implement	Weeks	Minutes	Days
Best For	Critical, highly specialized tasks	Quick tests, simple or generic tasks	Most domain-specific applications
Example Use Case	Medical diagnostics, legal analysis	Email formatting, sentiment analysis	Customer support bots, industry chatbots

Method	Analogy	Resource	Flexibility	Best For
		Cost		
Full Tuning	Relearning everything	Very High	Low (task-specific)	Critical, high-accuracy tasks
LoRA	Targeted cheat sheets	Low	High (swap adapters)	Efficient multi-task learning
Prompt Tuning	Secret exam hints	Minimal	Moderate (prompt design)	Quick, low-resource tasks





GPT and Modern LLMs

- Large language model that generates **human-like text**
- **Pre-trained** on massive datasets (books, websites, articles)
- **Transformer architecture** uses self-attention for deep context understanding
- Works by predicting the next word in a **sequence**, one at a time
- Can be fine-tuned for specific tasks (e.g., medical, legal, coding)
- Versatile: Used for chatbots, content creation, translation, summarization, and more
- Produces natural, relevant, and coherent text
- Scalable: **Handles huge amounts of data quickly**



Limitations and Challenges of LLMs

Bias and Hallucination

LLMs inherit societal biases from data and sometimes generate false or misleading information, which is a critical challenge.

Context Limitations

LLMs struggle with maintaining coherence over very long text due to fixed context window sizes.

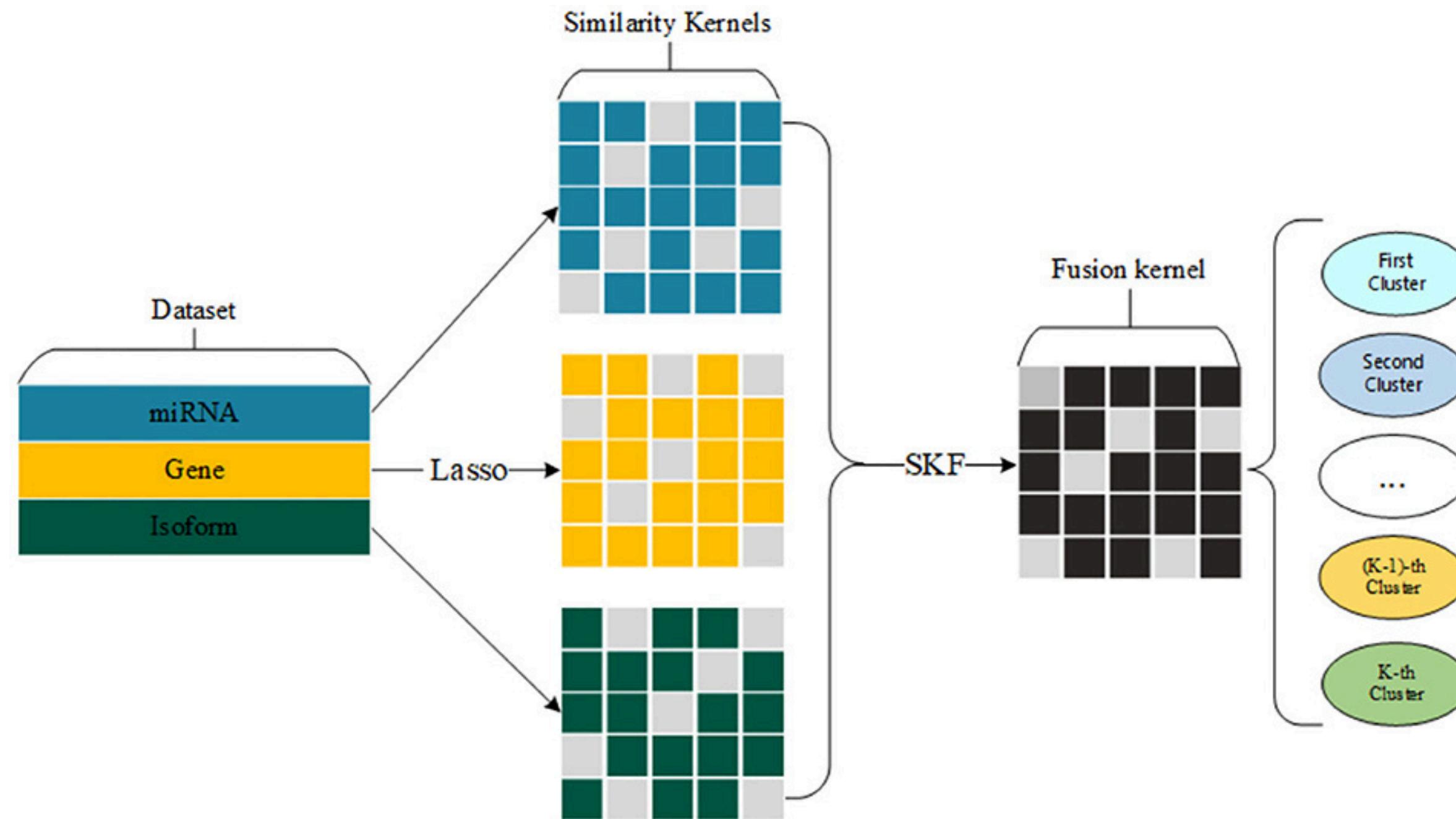
High Resource Costs

Training and running LLMs demand massive computational power and memory, posing environmental and accessibility concerns.

Lack of Real Understanding

These models generate plausible text without true comprehension or real-world awareness, limiting reasoning.

The Shift of LLMs from GPU to CPU-Based Inference



- Normally, each GPU kernel (chef) does its job and stores intermediate results in slow global memory (the shelf)
- With kernel fusion, multiple operations are combined into one kernel (both chefs at the counter), so intermediate results are passed directly using fast registers or shared memory (no shelf)
- This reduces memory traffic, speeds up computation, and makes the whole system more efficient

 SMALL PACKAGES

Microsoft's “1-bit” AI model runs on a CPU only, while matching larger systems

Future AI might not need supercomputers thanks to models like BitNet b1.58 2B4T.

KYLE ORLAND – APR 19, 2025 1:16 AM | 105



Can big AI models get by with "smaller" weights? Credit: Getty Images

Real-World Applications of

Content Generation

LLMs create articles, marketing copy, and creative writing efficiently, supporting media and advertising industries.

Chatbots and Assistants

They power conversational agents for customer support and personal virtual assistants improving user interactions.

Translation and Coding

LLMs offer accurate language translation services and assist programmers with code completion and bug detection.

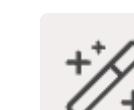


Conclusion: The Future of LLMs



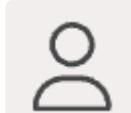
Ongoing Innovation

Research continues to address current limitations, improving accuracy, efficiency, and understanding in LLMs.



Transformative Potential

LLMs are poised to impact industries from education to healthcare, amplifying human capabilities through AI.



Ethics and Responsibility

Responsible development and awareness of ethical implications remain crucial to harnessing the power of LLMs safely.

Dhanyavad