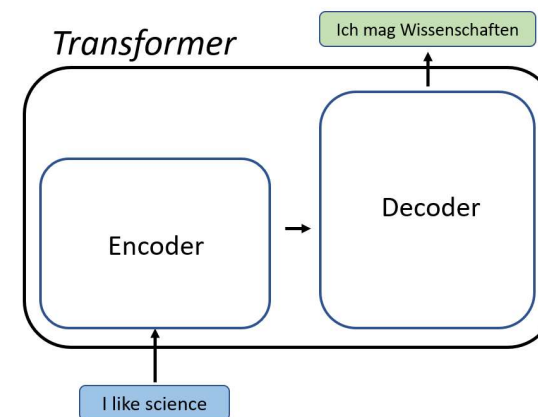
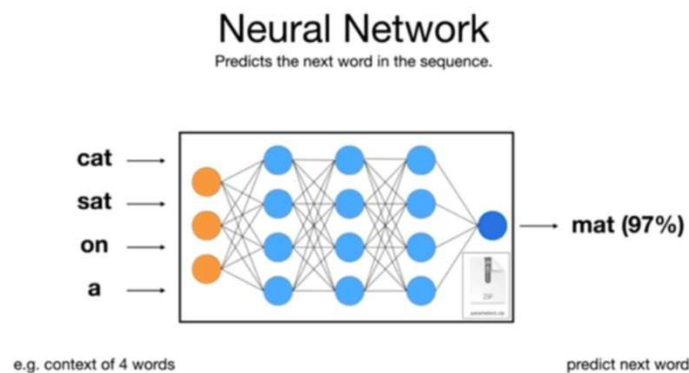
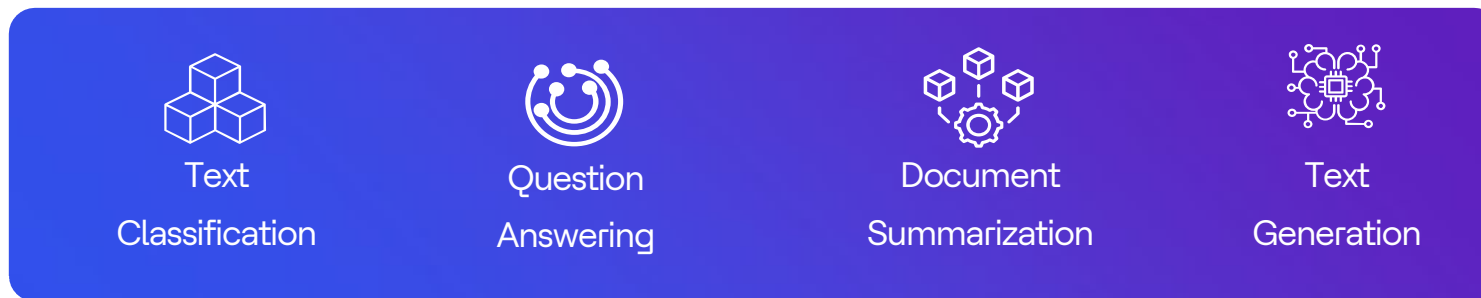


Cloud Native & AI Labs

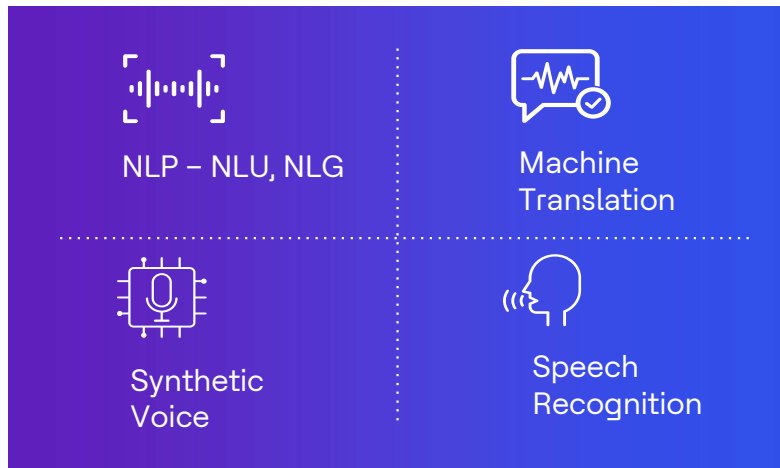
LLMs and Generative AI

Large Language Models

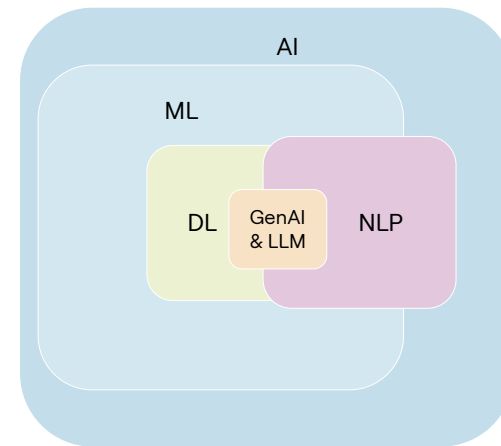
LLMs are trained to solve language problems like...



Generative AI



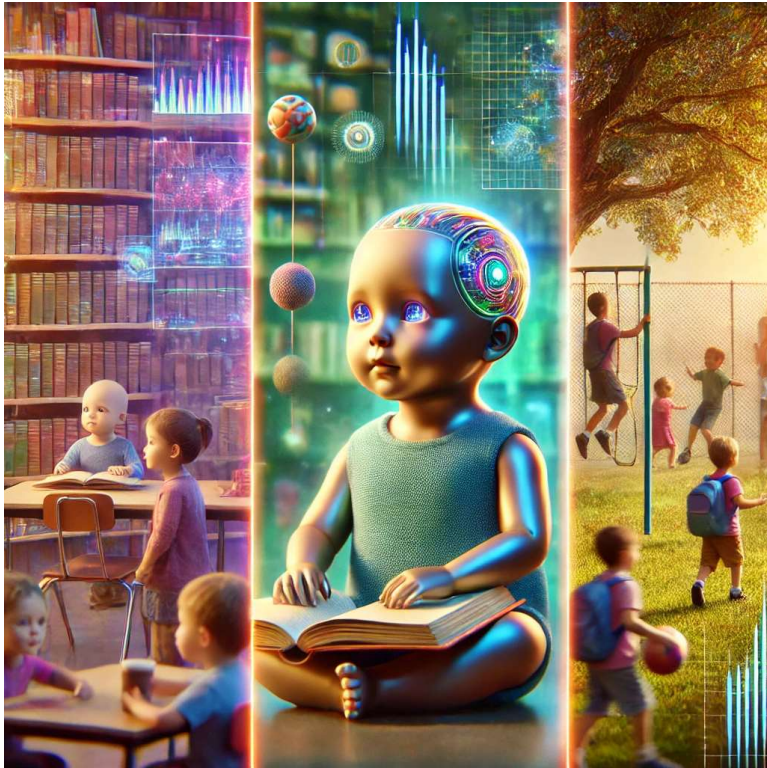
Artificial Intelligence



Domains of Generative AI



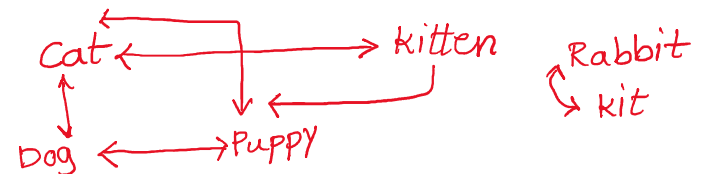
The Philosophy of Learning



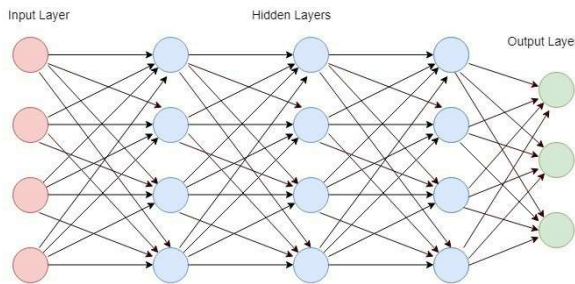
1. The world is chaotic, and our brains are optimized to learn and adapt to chaos.
2. Imperative → Declarative Paradigms → AI/Machine Learning
3. Types of Learning:
 - Supervised
 - Unsupervised
 - Reinforcement
 - Self-supervised

1. Word ordering
2. Image jigsaw
3. Denoising (Diffusion)
4. Game of compression

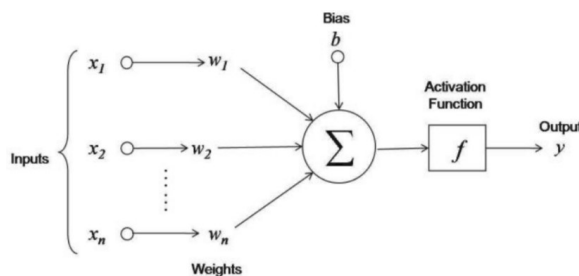
We are looking at new concepts by looking at the relationships between different concepts. We also reason in Language.



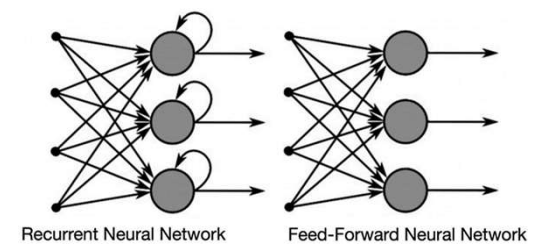
Neural Networks



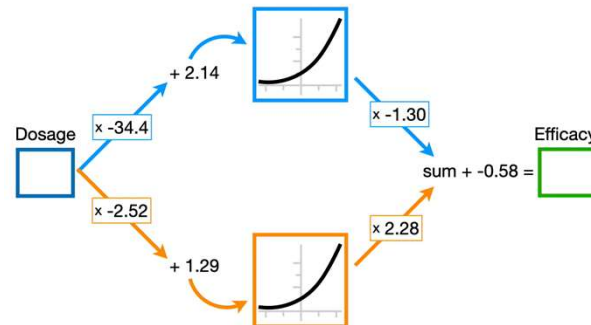
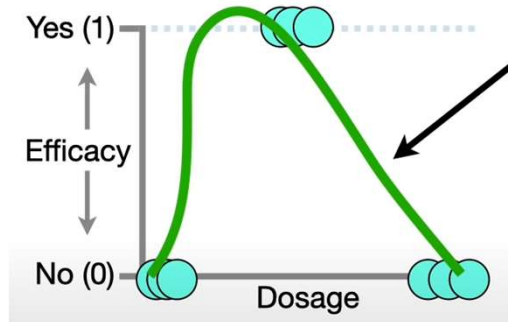
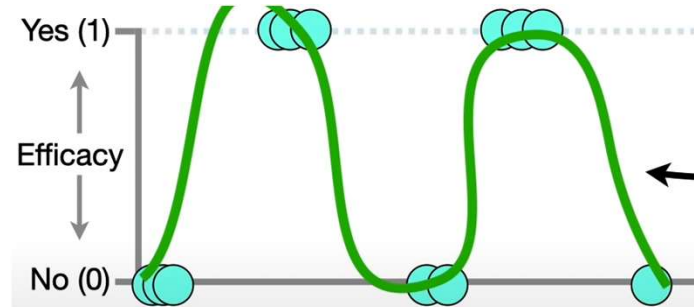
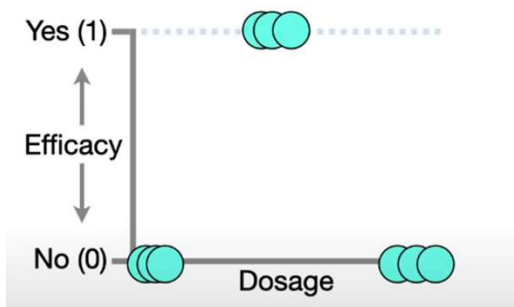
- Neurons - Fundamental unit of processing in a neural network.
- Neural network – a network of interconnected nodes or neurons in a layered structure that resembles the human brain. Used for machine learning process called deep learning
- 3 layers in a neural network – Input, Hidden, Output



- Types of neural network:
 - Feedforward Neural Networks
 - Recurrent Neural Networks
 - Convolutional neural networks
 - Generative Adversarial Networks
 - Etc.

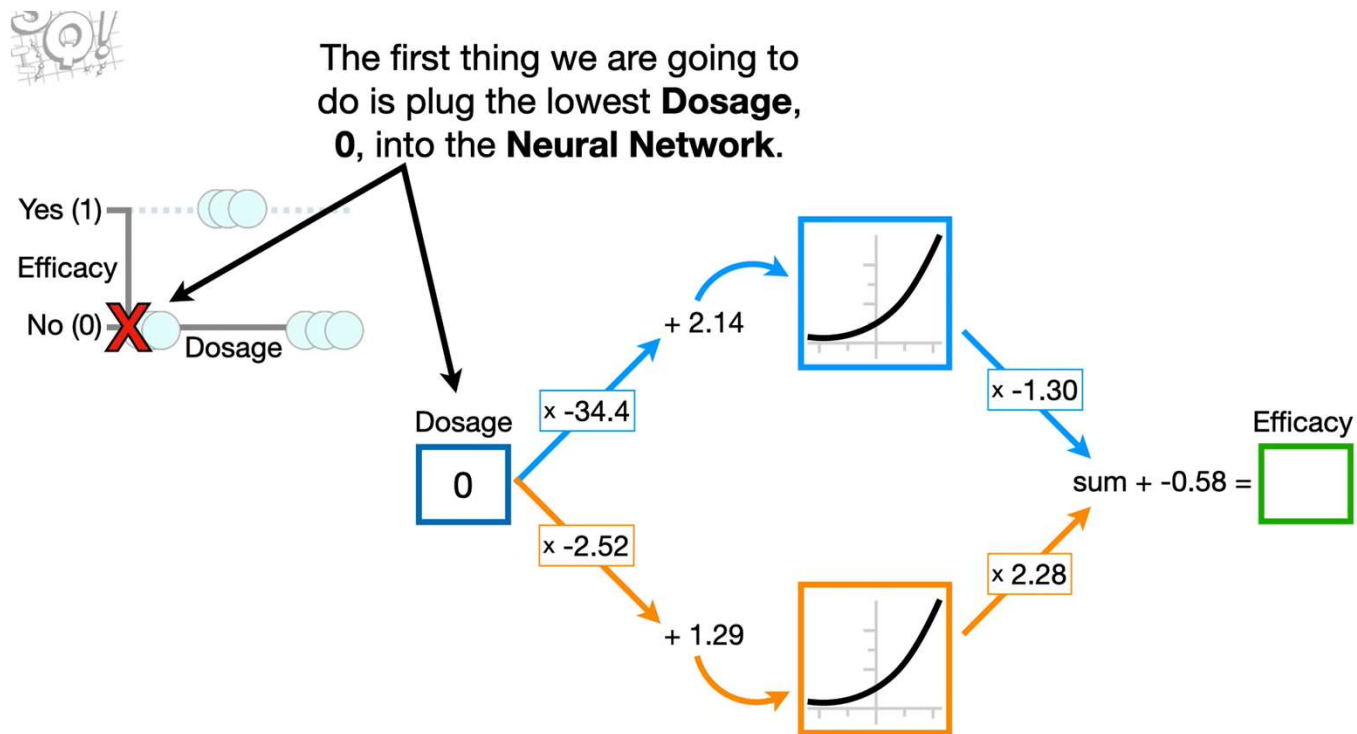


Neural Networks

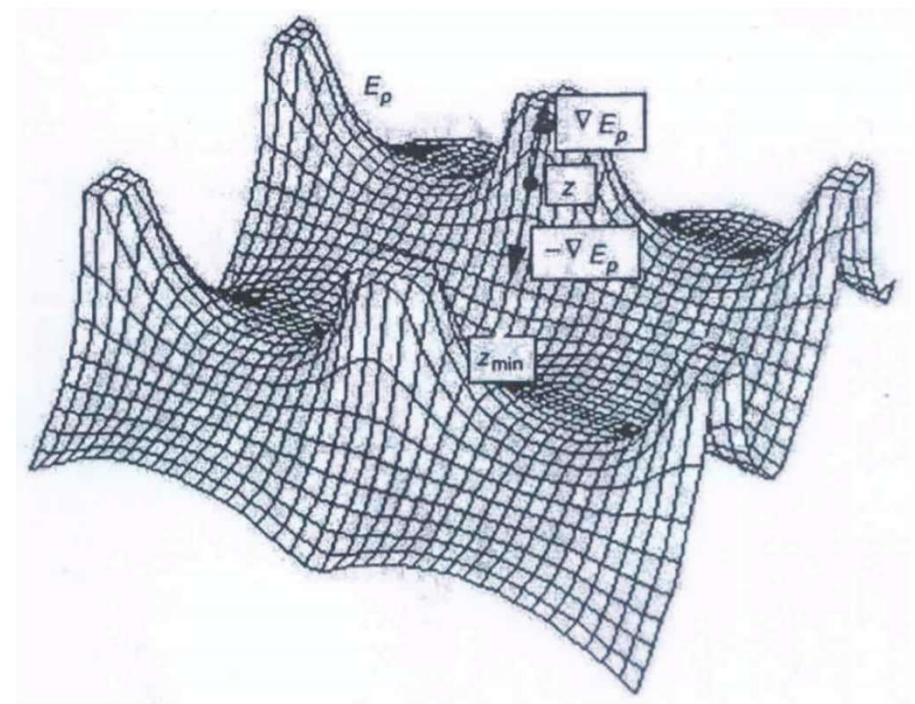
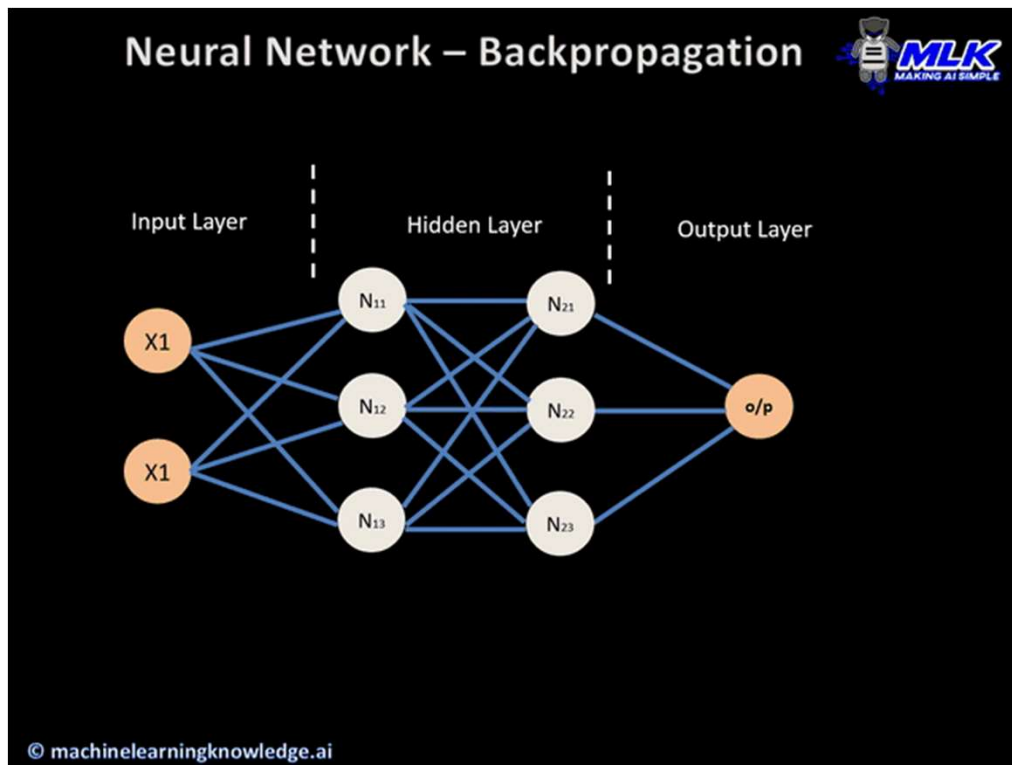


<https://deepnote.com/app/martin-alvarez-59be/Neural-Networks-1c9e6060-64bc-412e-9539-aeb20ef76214>

Neural Networks



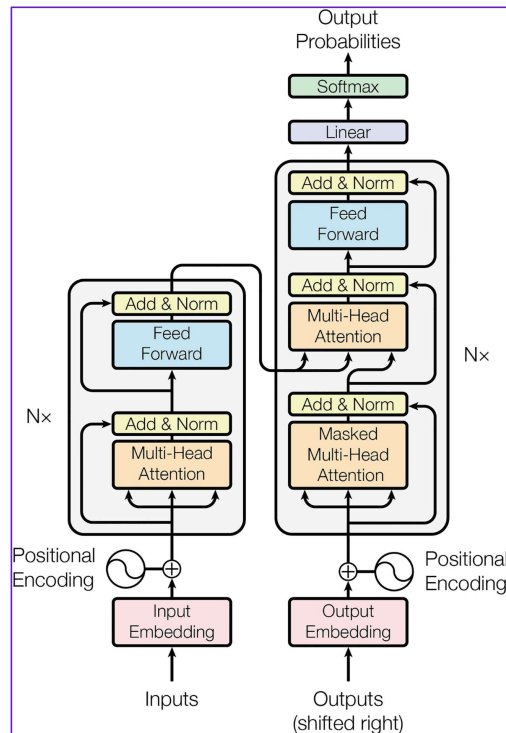
Backpropagation



Transformers

Seq2Seq Neural Network

Transformers



- **Self-Attention Mechanism** - allows the model to assign weights to different words in a sentence based on their importance and then compute a weighted sum of values
- **Multi-Head Attention** - To capture the relationships between the words in a sentence, the transformer uses multiple sets of self-attention mechanisms, known as heads. Each head learns different aspects of the relationships, and their results are concatenated and linearly transformed.
- **Positional Encoding** - Unlike RNNs, transformers don't have an inherent sense of order due to their parallel processing behavior. Positional encodings are added to the input embeddings to provide information about the position of each token in the sequence.
- **Encoder-Decoder Architecture** - The original transformer architecture consists of an encoder and a decoder. In machine translation, the encoder processes the source language sentence, and the decoder generates the target language translation.
- **Pretraining and Fine-Tuning** - Transformers are often pretrained on large set of data using self-supervised tasks. After pretraining, the models can be fine-tuned for specific tasks, adapting their knowledge to the target application.

{ LLM is two files = ① code ② parameters }

for Llama 2 70B parameters \times 2 Bytes 140 GB \leftarrow NN weights

140 GB is self-contained !

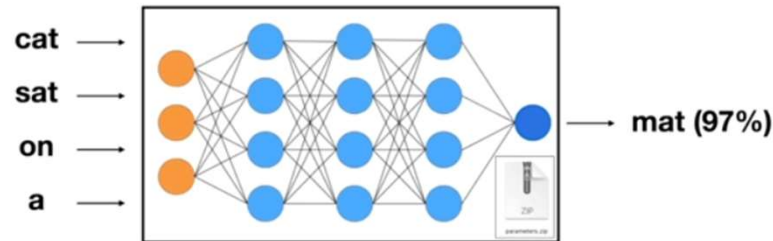
Training (compression) \sim 10TB \rightarrow $\frac{6000 \text{ GPUS}}{12 \text{ days}}$ \rightarrow \sim 140 GB

* Lossy compression

\rightarrow Llama 3
 \rightarrow Phi 3

Neural Network

Predicts the next word in the sequence.



e.g. context of 4 words

predict next word

! LLM $\frac{140 \text{ GB}}{\text{parameters}}$ $\frac{500 \text{ lines}}{\text{run.c}}$ llama-2-70b

$\sim 10 \text{ TB} = 6000 \text{ GPUs for 12 days} = \$2\text{M} \rightarrow 140 \text{ GB (lossy compression)}$

Mohsen is a person \leftarrow transformer architecture
 \downarrow Probability dreams Internet

it is an odd way of compression

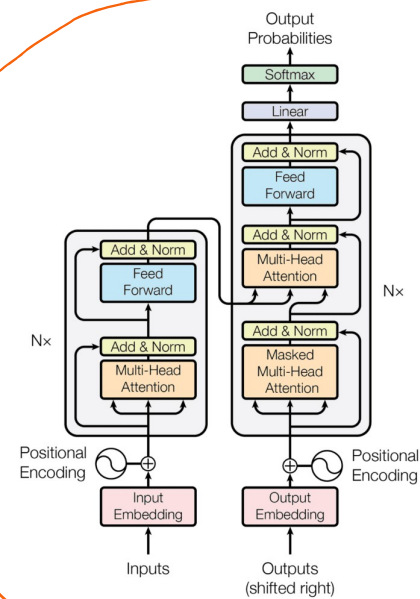
① Base model
Transformer

PreTraining
Tons of Internet
Every Year

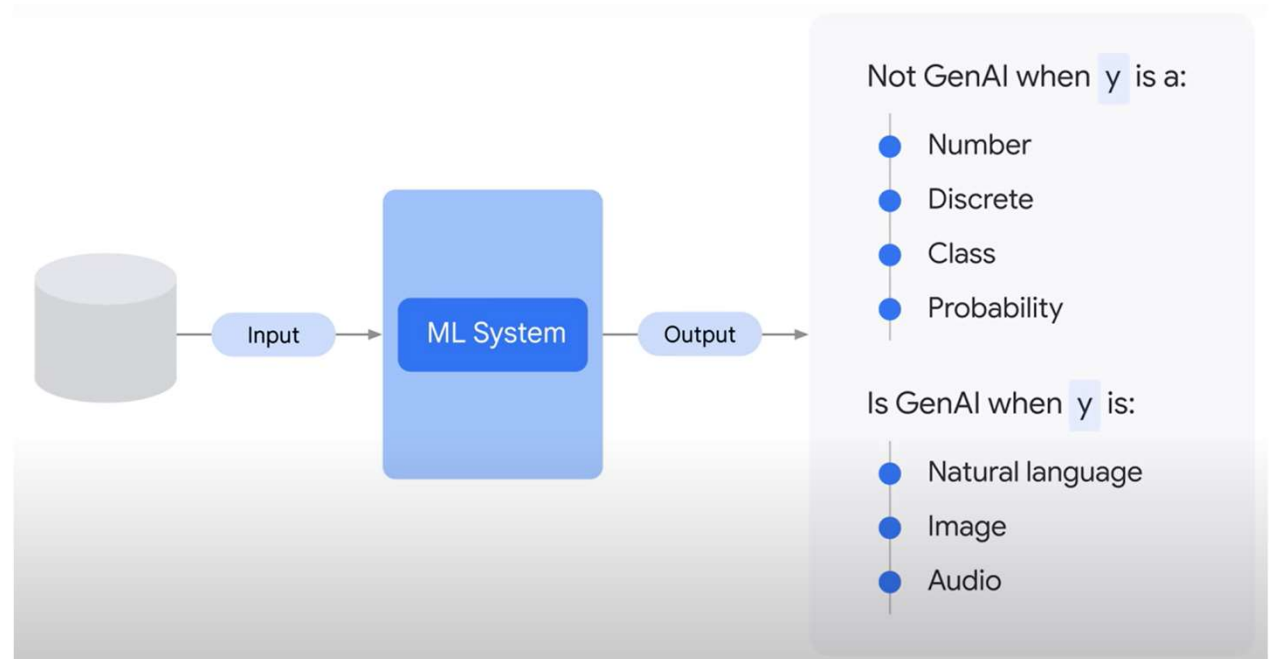
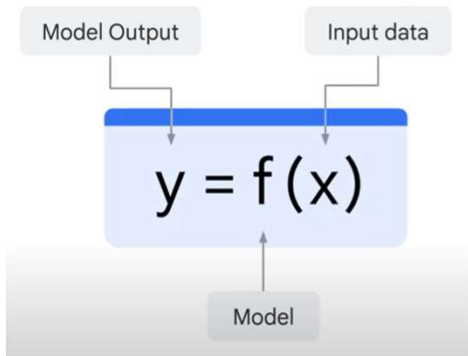
② assistant Model
1000k answers
and convo

"Fine Tuning"
helpful assistant
Every week

③ Comparison
stage
RLHF / classification



Gen AI vs Not Gen AI



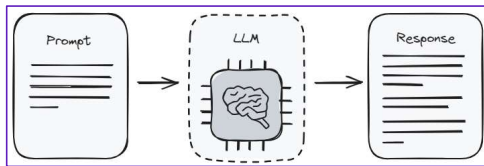
Generative AI – Fetching Data from LLMs

Prompting

Retrieval
Augmented
Generation

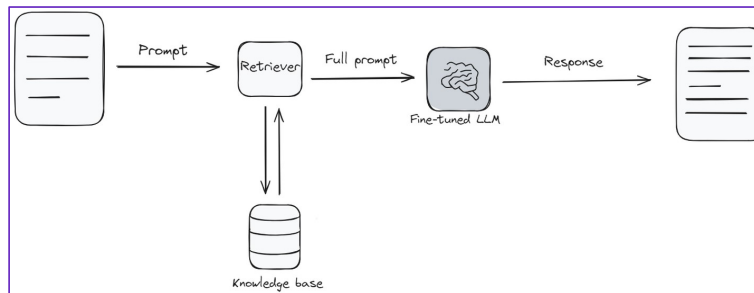
Fine Tuning

Model Training



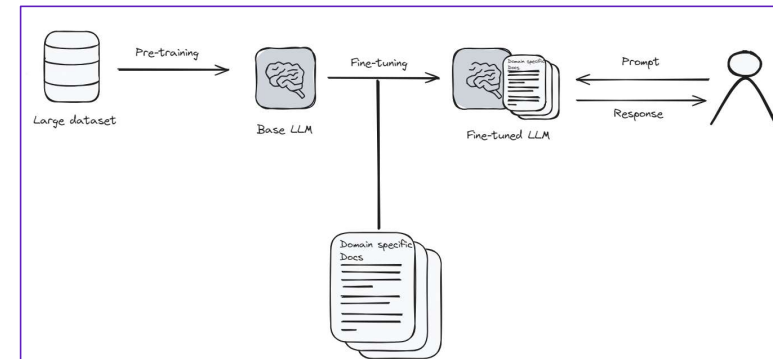
Prompt Engineering

- Easy to use
- Cost effective
- Flexible
- Inconsistent
- Limited customization
- Dependency on model's knowledge



RAG

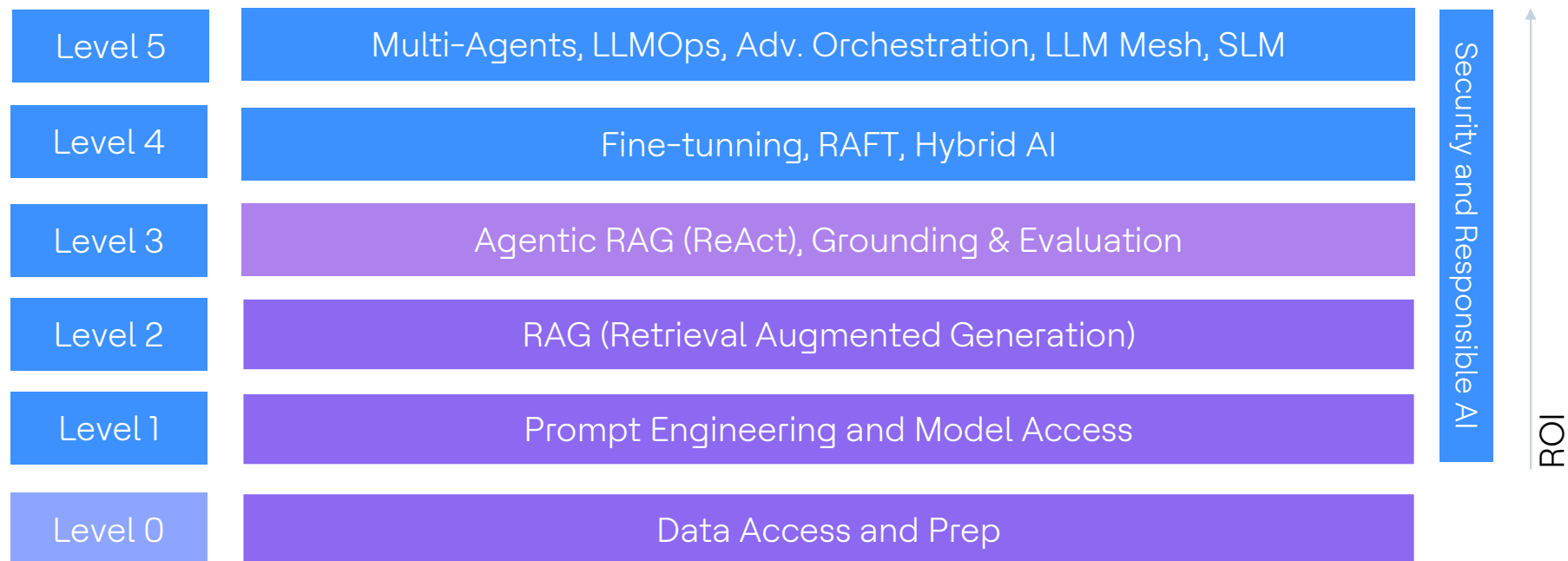
- Dynamic information
- Cost Balance
- Contextual data
- Complexity
- Resource intensive
- Data dependency
- Performance dependency



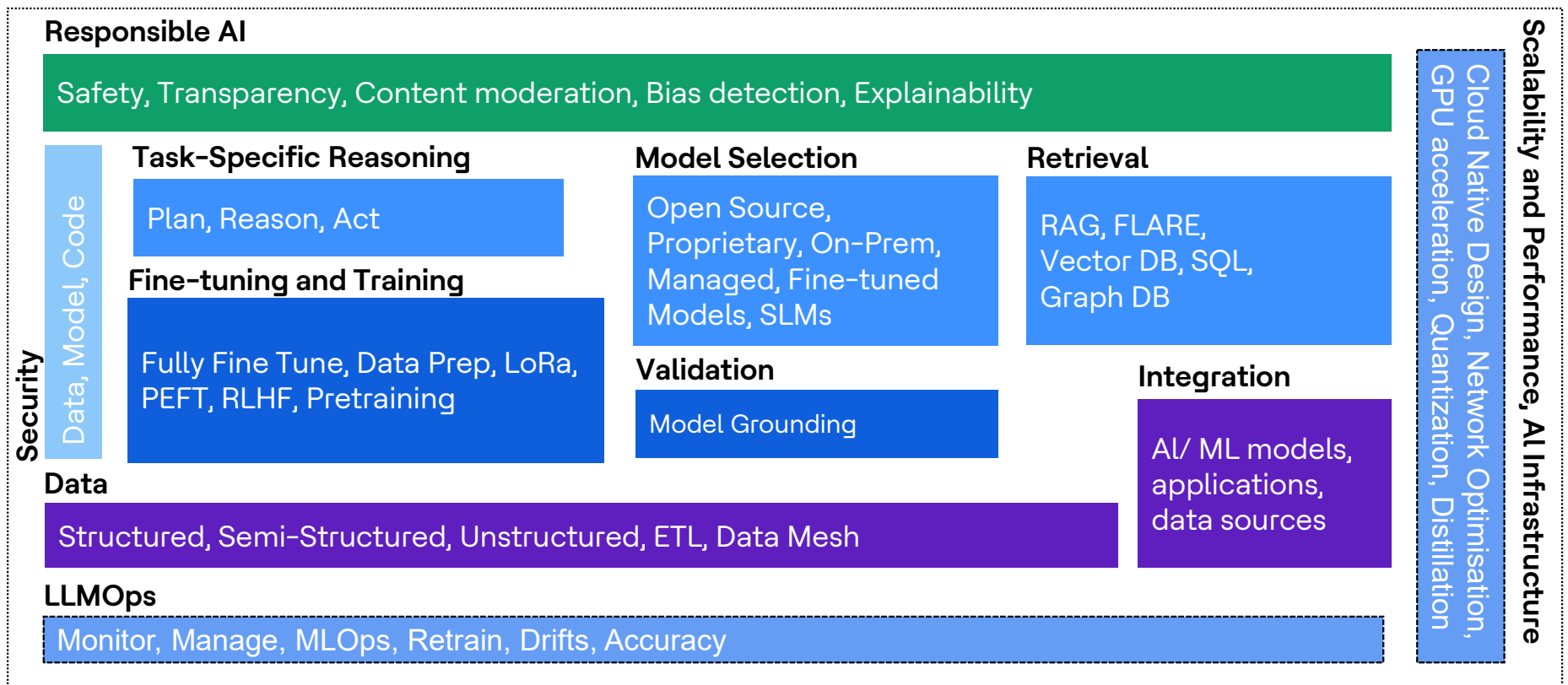
Fine Tuning

- Allows for extensive customization
- Improved Accuracy
- Adaptability
- Expensive
- Technical skills
- Data dependency

Gen AI Application Journey



Common Reference Architecture



A common reference architecture underpins consistency and efficiency across Gen AI solution development, serving as a blueprint for scalable and robust implementation. Whilst following a common logical architecture, it will be specialized for each vendor platform.