# Problems w/ LSTM & RNN

→ Slow

→ Requires a lot of data

→ Accuracy (Vanishing Gradient Problem)

→ $[\ \_\ [\ [\ \_\ |\ \_|\ \_)\ \_[\ \_|\ \_|\ \_\ ]\ ]$

Tokenize then Embed

⇒ <u>Transformer</u> : Picks up entire sequence.

→ { Query : What is the text      (Systematic
    Key   : Where is the text        Syntax
    Value : Here is the text      & Derivable
                                      Semantics )
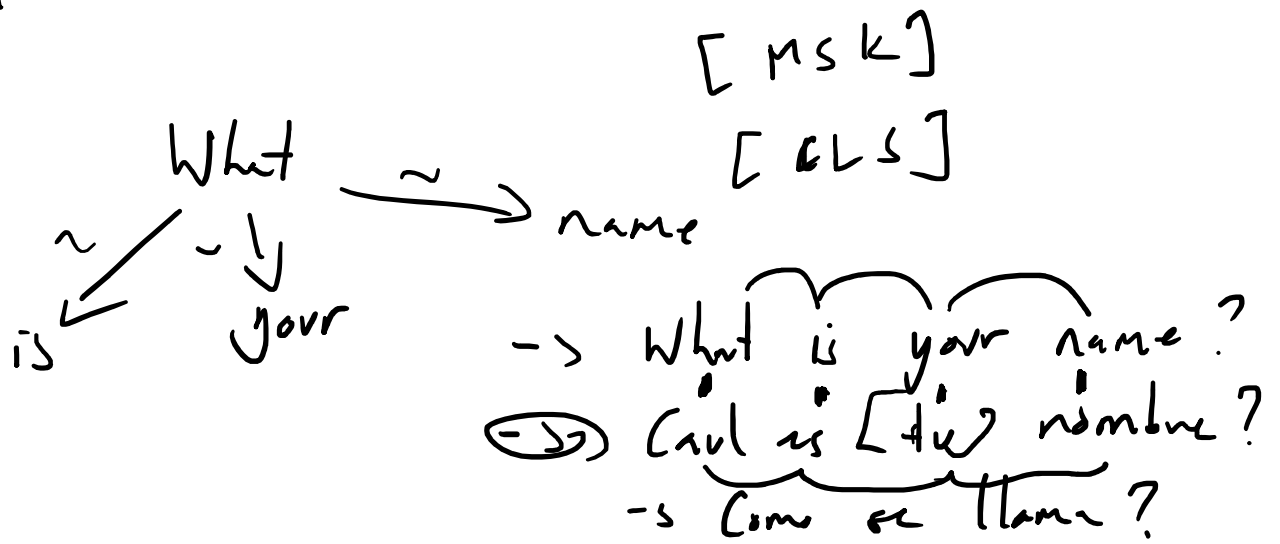
→ Encoder & Decoder

$q, k, V$         Tries to understand $q, k, V$.

→ GAN ~ Generative Adversarial Networks.

→ <u>Masking</u>      <u>What</u> <u>is</u> ⬭ name ?
                              <u>your</u>

→ Multi-Headed Attention your

[MSK]

[CLS]

What ~ → name

~ / ~↓

is  your

→ What is your name?

⊙→ Cual es [die] nombre?

→ Como se llama?

→ Separate Encoders & Decoders

Google

BERT
Bidirectional
Encoder
Representation
Transformer

Open AI   GPT

Generative
Pretrained
Transformer

→ Machine Translation & Text-Generation (GPT)

→ Text Classification & Text Summarization (BERT)
150M    7B

→ GPT-1, GPT-2, GPT-3, GPT-4 ~ Not o.Source.

→ BERT, RoBERTa, XLNet, DistilBERT ~ Hugging Face

→ BERT, RoBERTa, XLNet, DistilBERT ~ Hugging Face
110M    135M    1024 GPU ~ V100

→ GPT-3.5 & GPT-4                    ~ FAIR ~ Yann LeCun

→ Meta - LLaMa - 3.5 B parameters
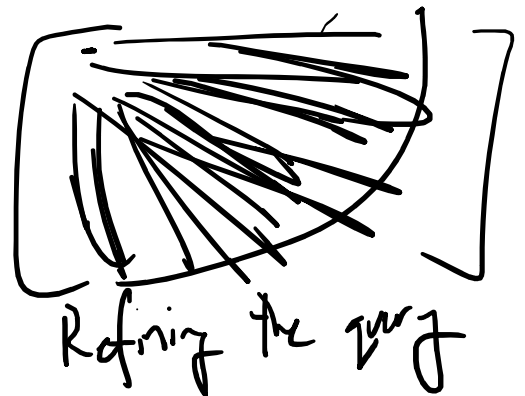V100 & 1100              Wikipedia & News ~ 80 B tokens

→ LLM ~ Large Language Model (High tech transformer)
  ~ Manifests - Use & calibrate this model
        [ (Agri) / (Automobiles) / (Electronics) ]
            [ 7 0 0 0 0 1 0 0 0 ]

  ~ Tell me about Lambo.              
                                      Refining the query
  → Prompt Engineering

    ┌ GPT ~ Open AI - Paid
    ├ Claude ~ Anthropic ~ Paid         → Llama Index
    ├ (Llama) Meta ~ OpenSource         → Lang chain ✓
    └ Gemma ~ Google ~ OS
          Mistra AI ~ Paid              → Colab
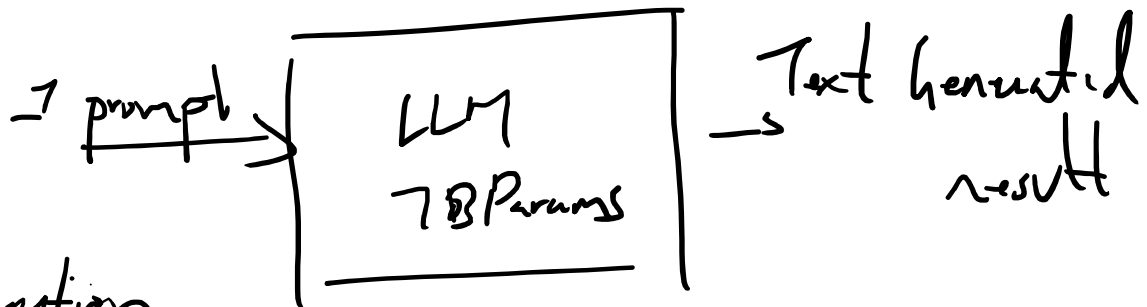                                        ↳ Local ~ Ollama

Vistra AI v Paid

↳ Local ~ Ollama

→ RAG : Retrieval Augmented Generation

→ Information is spread

→ prompt → LLM
7B Params → Text Generated
result

Data Ingestion

→ docx, pdf

🌐 Internet

Search Engines

Excel (SS)

→ Transform
Embed

(ETL)

Load

~

→ Load, Transform & Embed converts
data into vectors.

⤷ Chroma DB
⤷ FAISS
⤷ Object Dex

[ | | . . . | ]

~ Vector Store
or
↳ Database

}-> 1-AISS
Es Object Dox
-> Cassandra
(ASTRA)

Map & Reduce

Vector Database

-> Similarity Search

Task: Create a query engine for paper.