

# LLM BASICS & FLOW

So far:

- ✓ 1) Intro to AI / ML - NNs
- ✓ 2) ML model foundations
- ... 3) LLM foundations

## Large Language Model

big scale  
↓  
{ 100B - 5T parameters }  
Pre-trained on internet-scale data

trained on text  
↓  
for most languages  
in the world

ML model,  
specifically, NN

C P T<sub>s</sub>

Generative  
generates text,  
image, video,  
audio

Pre-trained  
already  
trained on large  
datasets of text,  
audio, video,  
images

Transformers  
a type of  
NN architecture  
pioneered first by  
Google in 2017  
↓  
("Attention is all  
you need")

default  
cloud provider  
1) AWS  
2) Azure  
3) BedRock

## OPEN-SOURCE

Meta - <sup>3B | 7B | 70B</sup> Llama 2/3

BERT

Code Llama

Stability - STABLELM  
STABLEImage

MISTRAL

HuggingFace

marketplace of sorts  
for open-source  
AI LLMs in  
particular

## CLOSED (COMMERCIAL)

openAI - GPT3.5, GPT4

Azure - GPT3.5, GPT4

Anthropic - Claude 3?

Google - Gemini Pro

How do LLMs go through generating text:

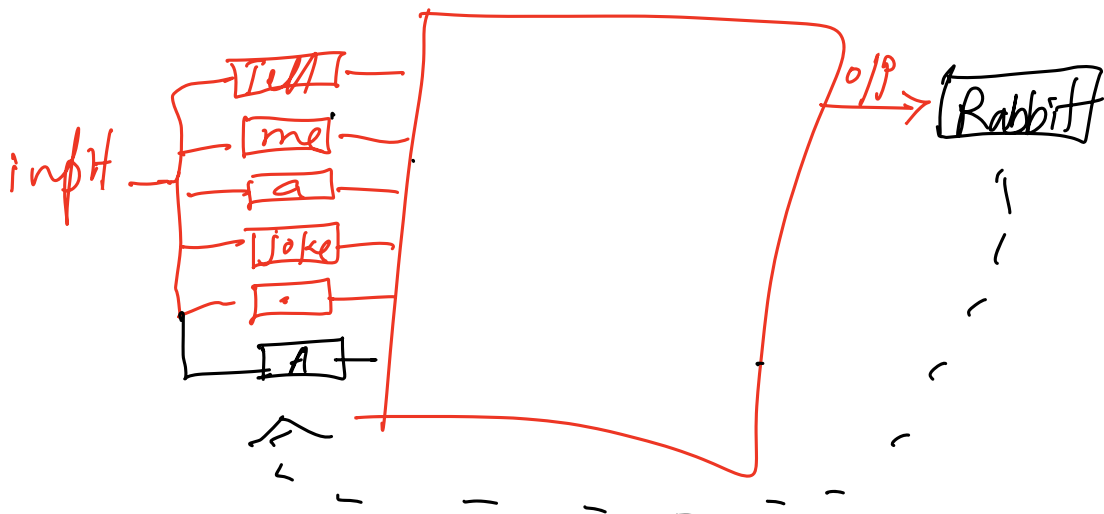
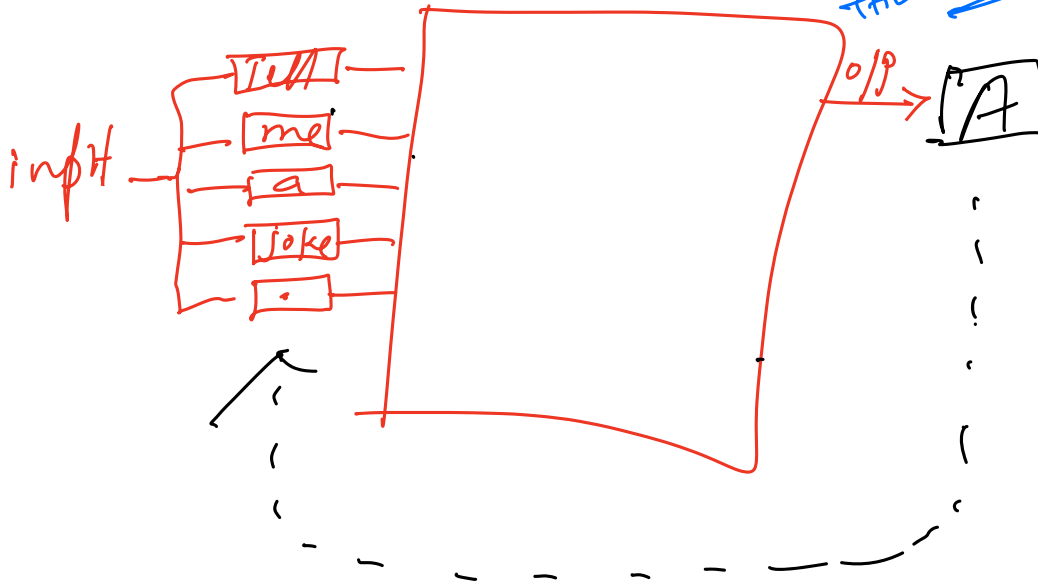
Prompt: Tell me a joke.

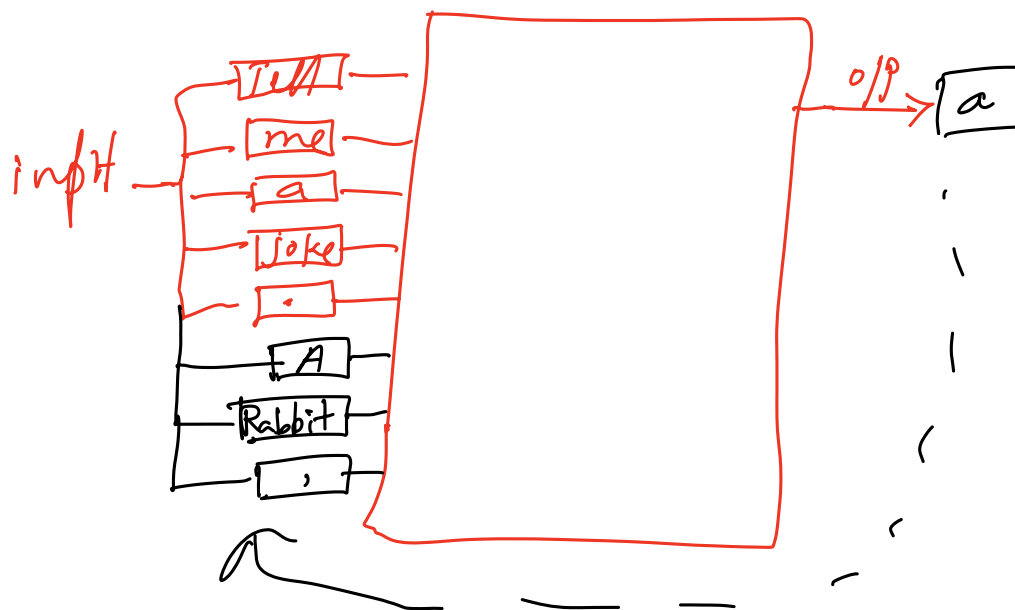
Response: A Rabbit, a lion, and a deer...

<END>

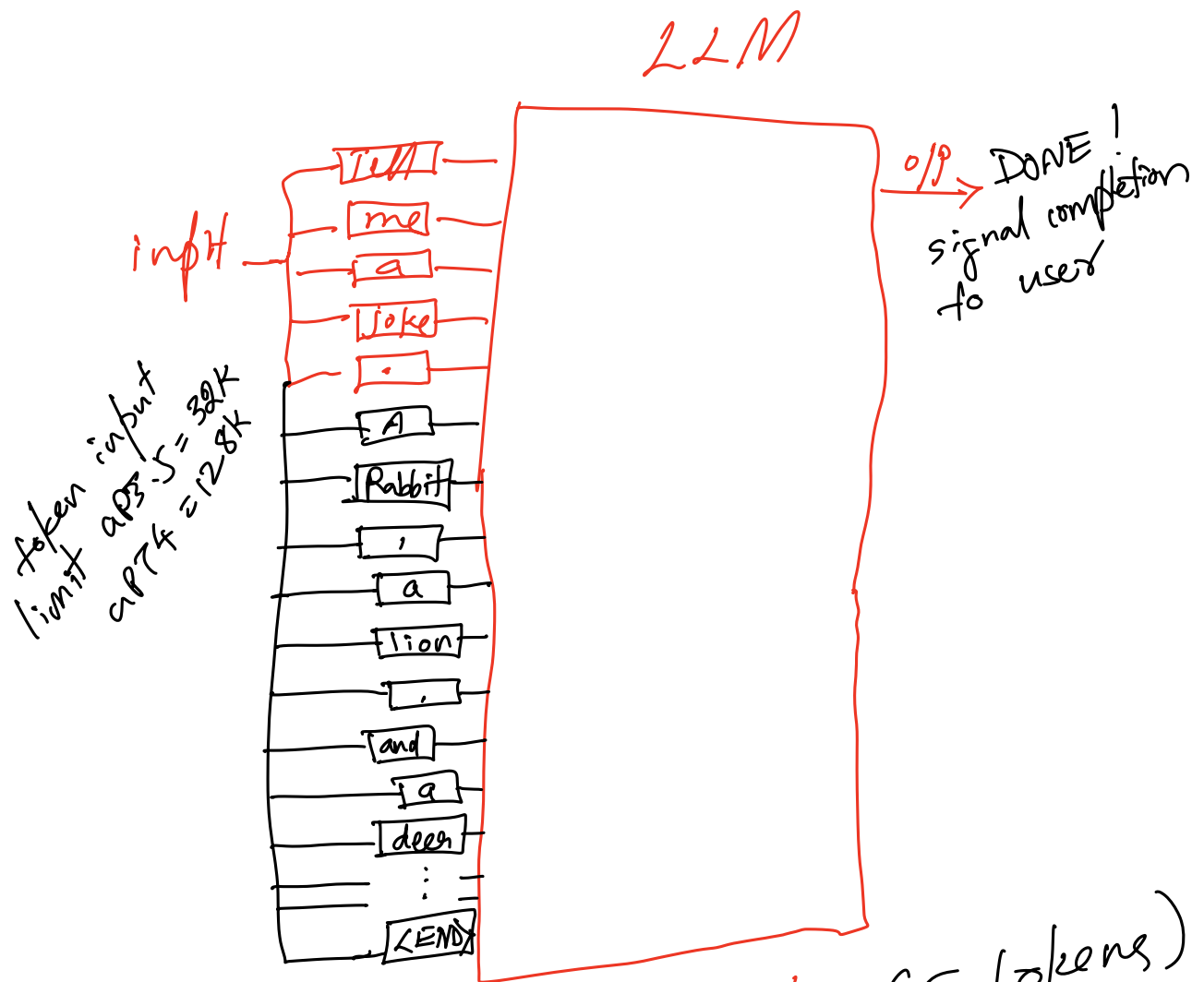
LLM

FOR A GIVEN  
TOKEN SEQ, PREDICT  
THE NEXT WORD





Finally,



Prompt: Tell me a joke. (5 tokens)  
Response: A rabbit, . . . . (26 tokens)

LLMs are good at generating  
"THE NEXT BEST WORD"



The cat is good

The cat is a liar

The cat is fat

George is a good liar

The cat is a good liar