

Symbolic AI. — a kind of good old-fashioned AI (GOFAI)

- ↳ Represents knowledge as explicit symbols and rules
- ↳ System manipulate these symbols using logic and search algorithms
- ↳ Knowledge is (usually) hand-crafted by experts.

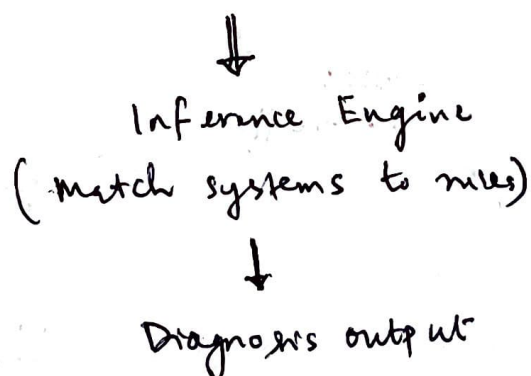
Ex-1

Knowledge Base (KB)

IF fever AND cough AND fatigue
THEN diagnosis = flu

IF fever AND rash
THEN diagnosis = measles

IF chest_pain AND breath_shortness
THEN diagnosis = heart_issue



Q. What about new symptoms ?

Q. What about old systems but not needed by the KB ?

Ex-2

Chess position



Hand-coded
eval fun

Queen = 9 points
Rook = 5 points
...
Pawn = 1 point



min-max
tree

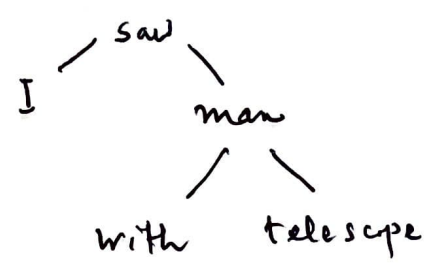
+ centre control
+ King safety
+ Pawn structure

Ex-3

Sentence: "I saw the man with the telescope"

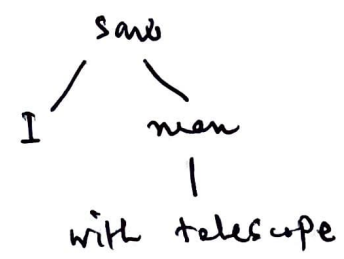
2 interpretations

①



(I used the telescope)

②

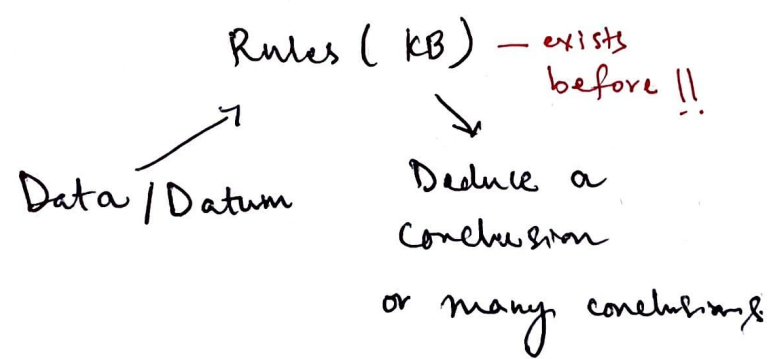


(the man used telescope)

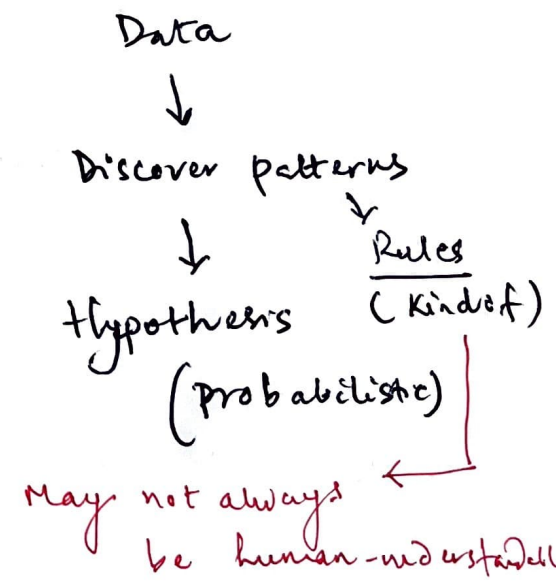
Q. Can symbolic AI "Learn" from data? (YES: INDUCTIVE LOGIC PROGRAMMING)

Q. What is learning?

□ Reasoning (deductive)



"Inductive" learning



Let's, for now, restrict ourselves to text data:

a. Language patterns exist in data.
Instead of writing rules, can we learn patterns from billions of text examples?

Foundation Models: (FMs)

FMs learn representations* and patterns* directly from massive amount of data.

A single model trained on broad data, then adapted to many tasks.

⇒ shifts from rule-based AI to learning based AI.

(Read: Brown et. al 2020, "Language models are few shot learners" (GPT3))

The man went to the bank...

which bank? for what?

GPT-4: IF bank AND deposit THEN FinBank.

NEWFAI: Learn from millions of sentences with "bank" in different context.

Automatically learn associations.

FMs example:

modality

models

→ Text

GPT-4, BERT, LLaMA

Vision

CLIP, DINOv2

Audio

Whisper, AudioLM

Code

Codex, CodeGen

Multiple models

Gemini, GPT-4o.

Biology

ESM, AlphaFold*

Chemistry

ChEMBERTa

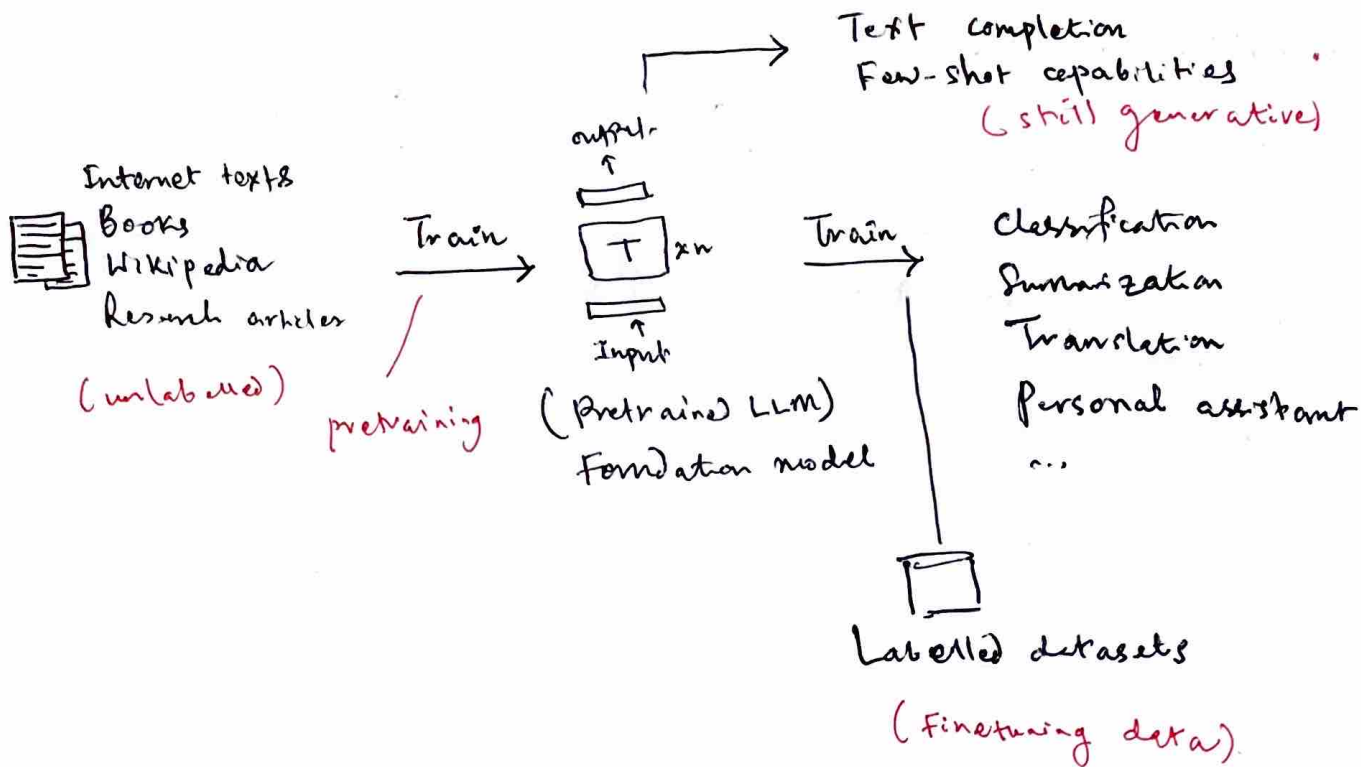
→ Pretrained on language data (texts)

also called "language model" or

"Large language model" or

"Pretrained LLM"

FMs: Serve as foundation for many downstream tasks



Q. What is being modelled by a FM?

FM is a generative model that learns a probability distribution over sequences.

→ Given context, predict next token.

Bayes' rule

$$P(x|y) = \frac{P(y|x)P(x)}{P(y)}$$

or

$$P(\underline{x} | \underline{y}) \propto P(\underline{x}|y)P(y)$$

Remember this

$P(\text{"the cat sat on the mat"})$

$P(\text{"the mat sat the car on"}) \leftarrow \text{very low prob.}$

$x_1 \quad x_2 \quad x_3 \dots$

autoregressive decomposition:

$$P(\underline{x}) = P(x_1, x_2, \dots, x_n)$$

$$= P(x_1) P(x_2 | x_1) \dots P(x_n | x_{n-1}, \dots, x_1)$$

product of text token (word) probabilities
here.
given all prev. words.

Q. How does one estimate $p(x_i | x_{i-1}, \dots, x_1)$

Traditional models: n-gram counts

FMs: Deep networks approximate

Training objective:

Input: The cat sat on the ← prompt

Predict: mat (high prob.)

floor (medium prob.)

water (low prob.)

Q. Can we say the model must be learning:

Grammar: "cat" follows an article "the"
Semantic: Cat is an animal (meaning) (syntax)
World Knowledge: Cats do not like water.

Prompt: The cat sat on the ?

temperature
scaling.

temperature = 0
(deterministic)

mat


temperature = 1
(creative)
random !!

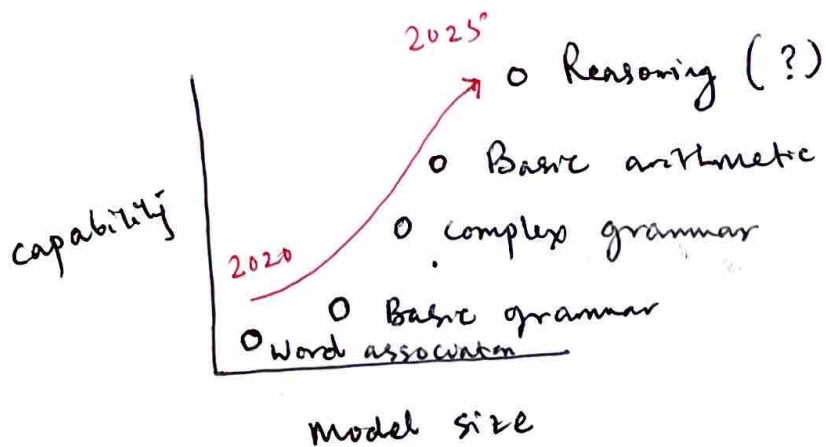
flower bed

$$p(x_i) = \frac{e^{(z_i/T)}}{\sum_j e^{(z_j/T)}}$$

$T > 1$: Flatter distribution over next word probs.

$T < 1$: Sharper distribution over "

$T = 20$ 
mat floor flower....



Different types of models and training objectives.

* GPT style (causal / autoregressive) Decoder only

The cat sat on ?

└──┬──┬──┬──┘

← attention flow

→ Generation task

"the"

* BERT style (masked language model) Encoder only

The cat [mask] on the mat

← bidirectional attention →

"sat"

→ Understanding task

* TS-style (Encoder-Decoder)

→ Seq to Seq task

Translate to French "Hello"

↓

Bonjour

Q. ChatGPT or Claude ?

How are these built ?

GPT — FM

↓

RLHF

+ Instruction Tuning

(Future lectures)

Let's walk through a quick text generation example with GPT-style model:

USER: "~~the~~ The cat sat on the" (complete the sentence)

Stage 1 Tokenization Text \rightarrow Numerical tokens
(mostly BPE: Byte-Pair encoding)
subwords

"The cat sat on the"

\downarrow

["the", "cat", "sat", "on", "the"]

\downarrow \downarrow \downarrow \downarrow \downarrow

[2061, 374, 21, 53, 2061] Token ID

Stage 2 map or convert token IDs to vectors

Token ID \rightarrow [2061] \rightarrow Embedding lookup table \rightarrow [0.23, -0.89, ..., 0.12]

d-dimensional vector
(for GPT $\Rightarrow d = 4096$)

each token \rightarrow d-dim.

Potentially, capture semantic meaning.

additional positional info: (more on this in the DL class)

Pos 1: [0.00, 1.00, 0.00, 1.00, ...] (d-dim.)

Pos 2: [0.84, 0.54, ...]

\vdots

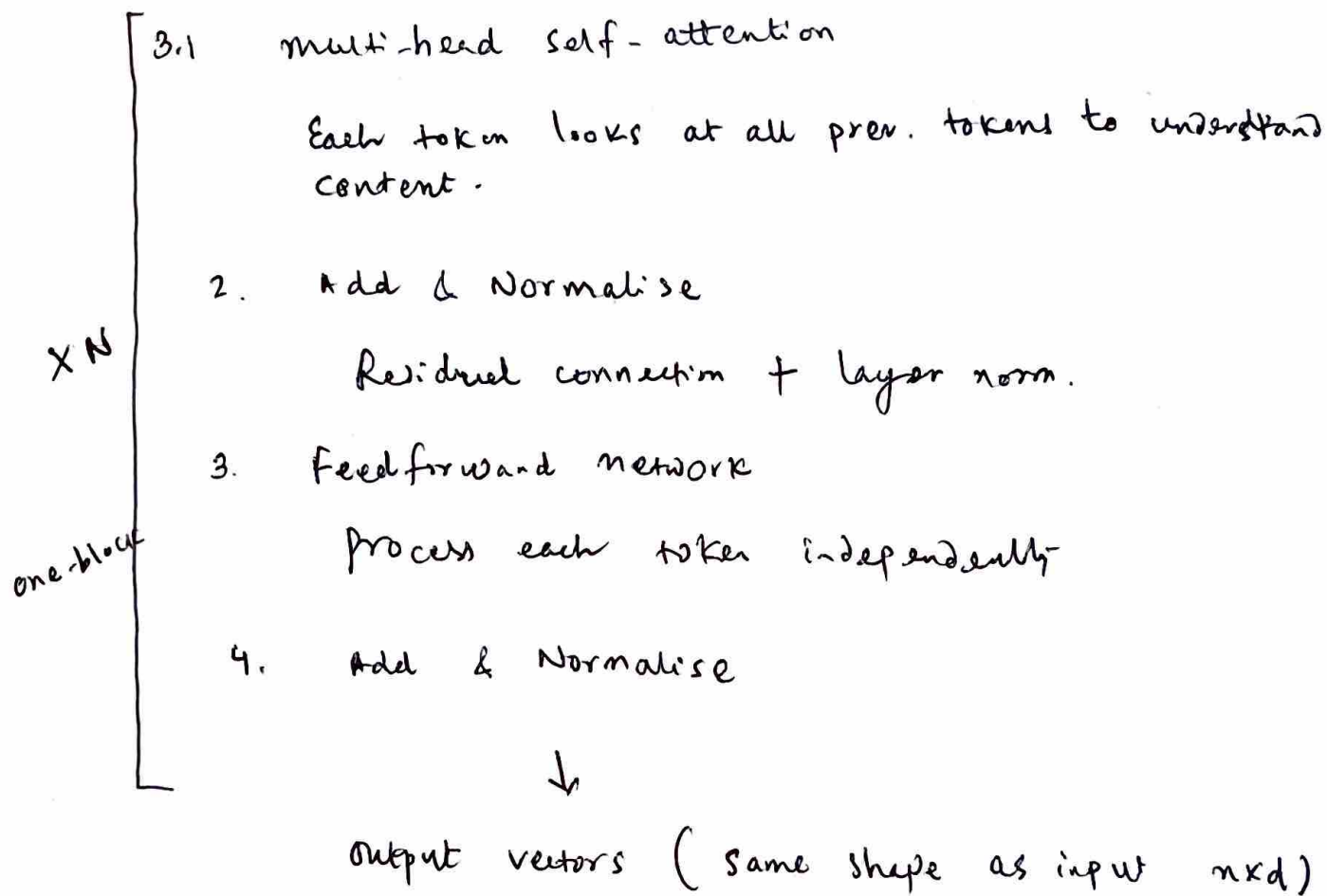
Final emb for each token:

"The": [0.23, -0.89, ...] + [0.00, 1.00, ...]

only in classical GPT inp

Stage 3 Transformer block.

Input: Sequence to vectors (one per token)



$N = 96$ (in GPT-3)

$N = \sim 120$ (GPT4)
Claude

① Self attention: "The cat sat in the"

$$\begin{aligned} \text{att}(\text{sat}, \text{cat}) &= \text{softmax}(\text{att_score}) \approx \frac{0.70}{0.1} \approx 0.95 \\ \text{att}(\text{sat}, \text{the}) &= \text{softmax}(\text{att_score}) \approx 0.01 \end{aligned}$$

② weighted sum of token rep.

strong relevance

$$\text{new_repr}(\text{sat}) = 0.01 \times \text{repr}(\text{The}) + 0.95 \times \text{repr}(\text{cat})$$

Multiple heads: Doing these in parallel multiple times

Note
Original
token dim
"d" is
SPLIT across
heads.

Head 1: Learn syntactic relations

Head 2: Learn semantic relations

Head 3: ...



Concat all heads → Final representation.

Stage 4

After N transformer layers.

Final hidden state for last token → prediction



Final transformer output

d-dim → v-dim.



vocab. size (50,000)

logits = [v-dim entries]



softmax



probs.

water 0
0 0.001
0

mat 0 0.85 → $p(\text{mat} \mid \text{"The cat sat on the"})$
= 0.85

floor 0 0.12