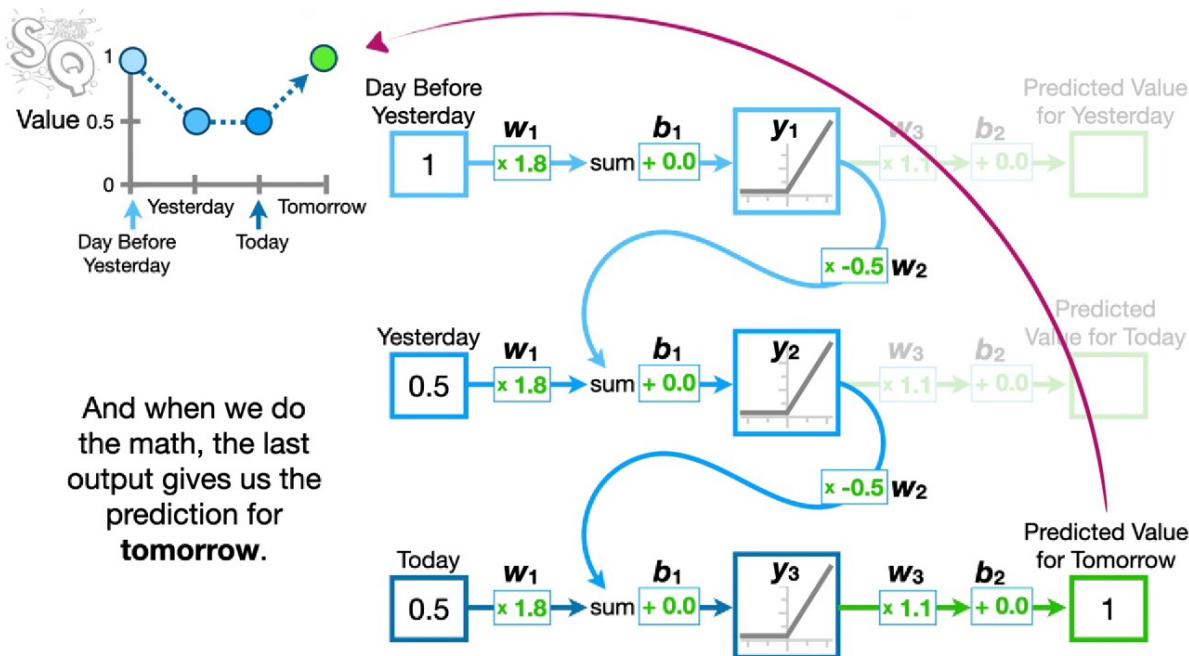


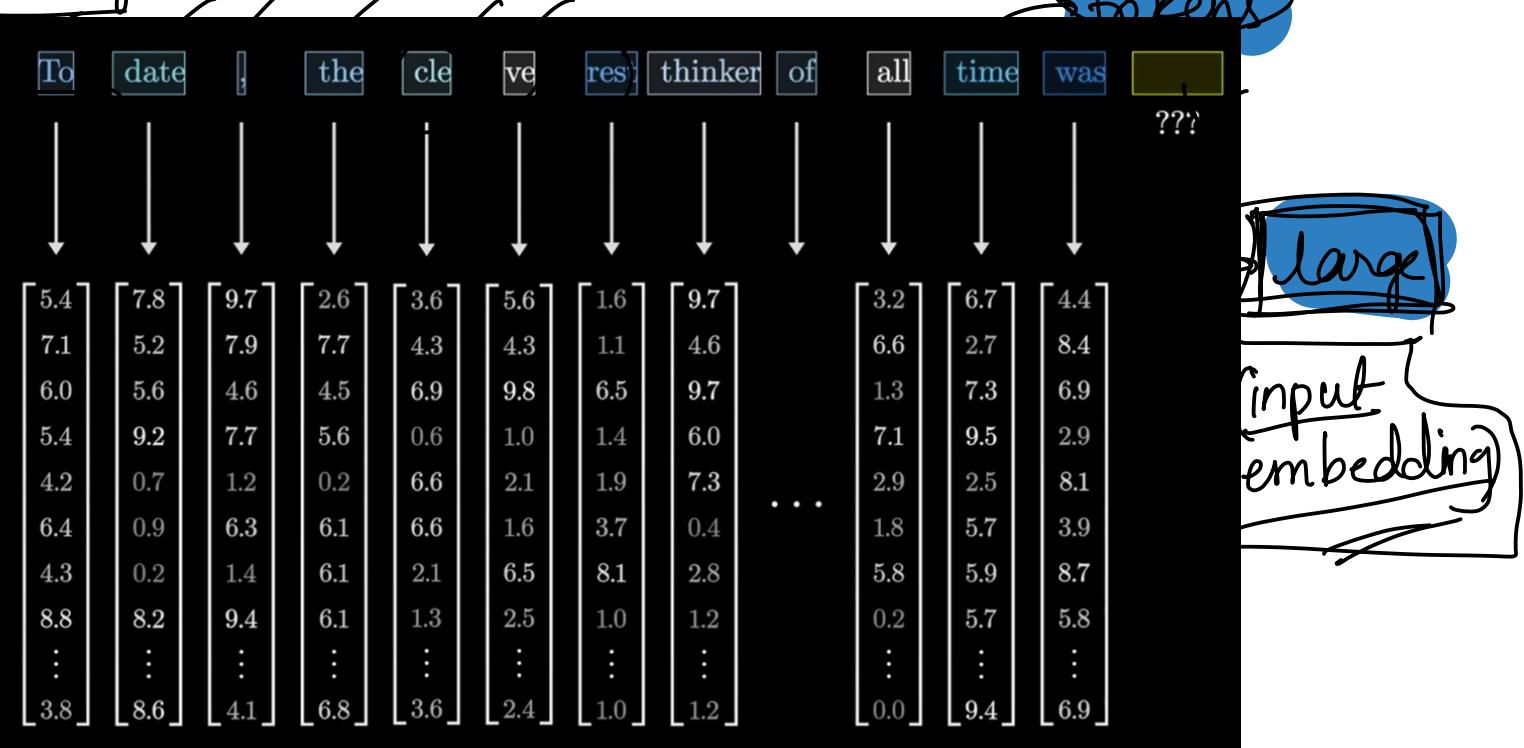
RNNs



sequential data, customisable no. of inputs

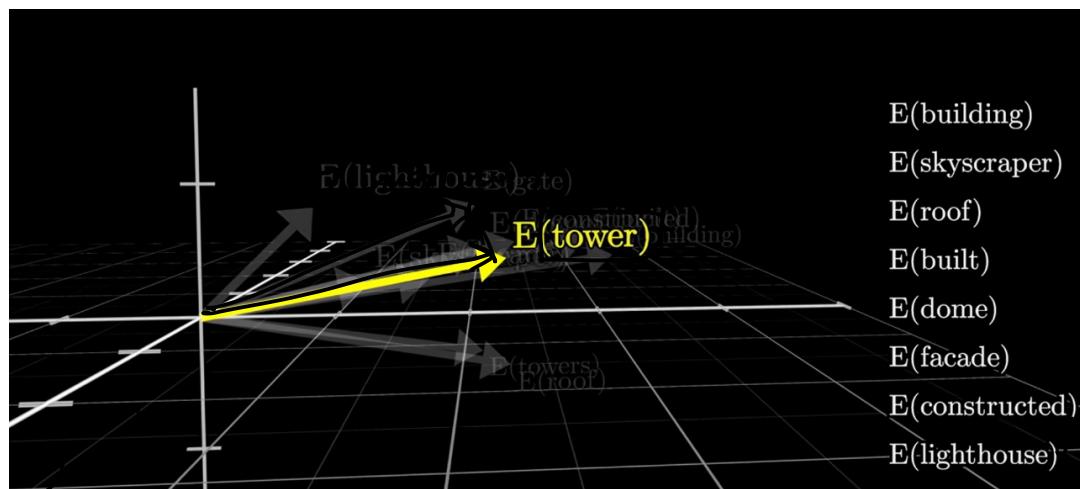
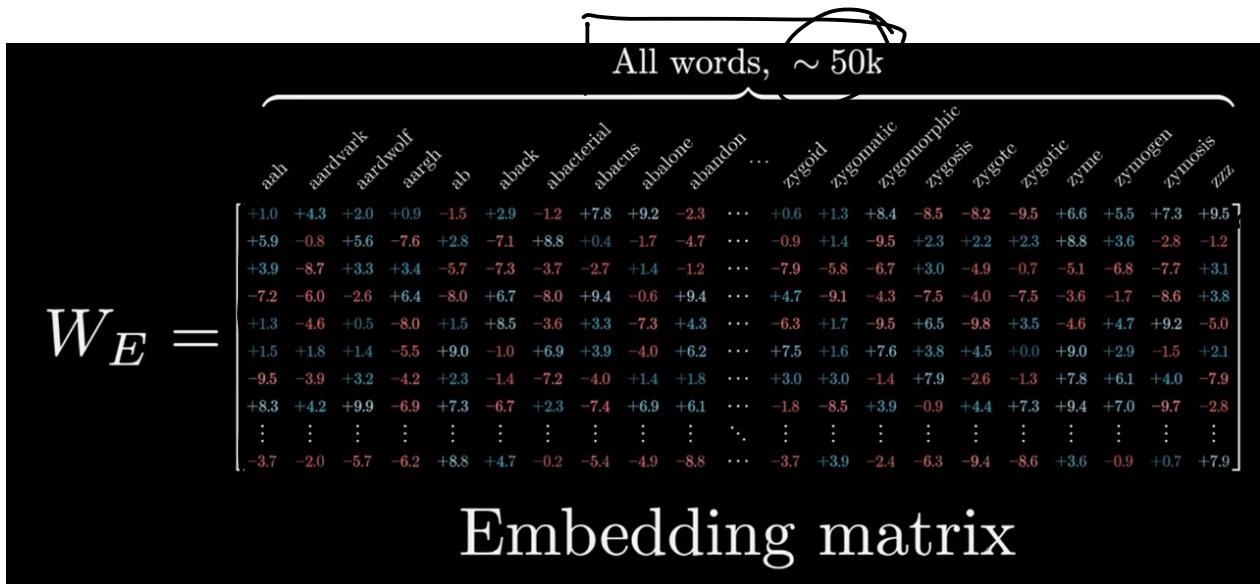
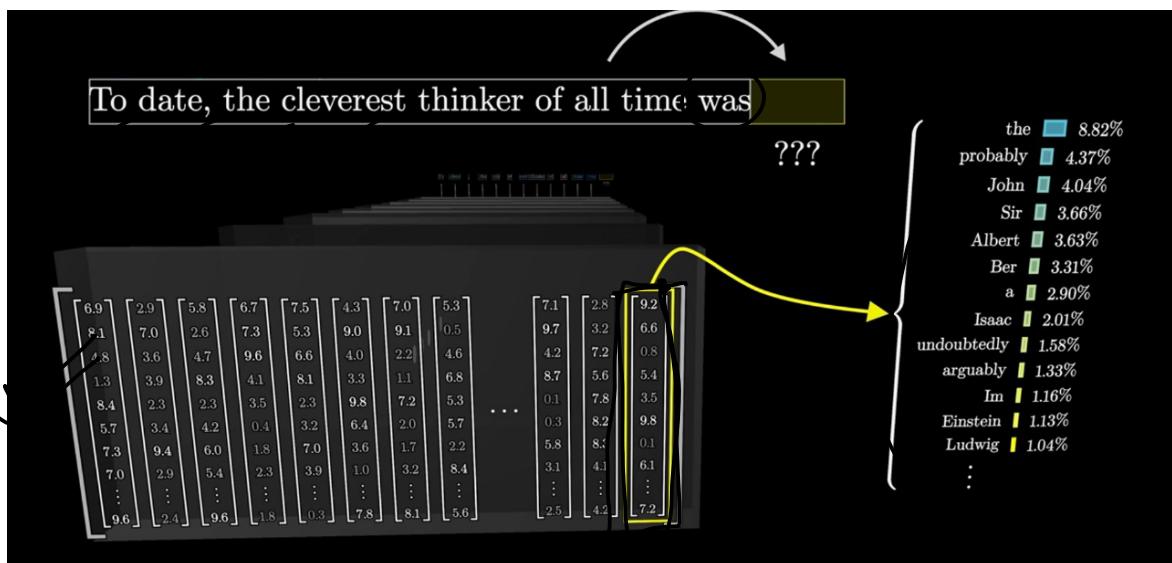
LMs (Attention mechanism)

→ Predict next word by taking all words in one go.



A machine learning [model]

fashion [model]



semantic
positional
encoding

dot &
cross
product
for Attention
mechanism

content size
(distance)

2 types :- decoder only & encoder-decoder

Pretaining \rightarrow fine tuning \rightarrow pre training

Response generation

\rightarrow One token at a time (small model)

i) probability distribution

w_1, w_2, \dots, w_t

$$P(w_{t+1} | C)$$

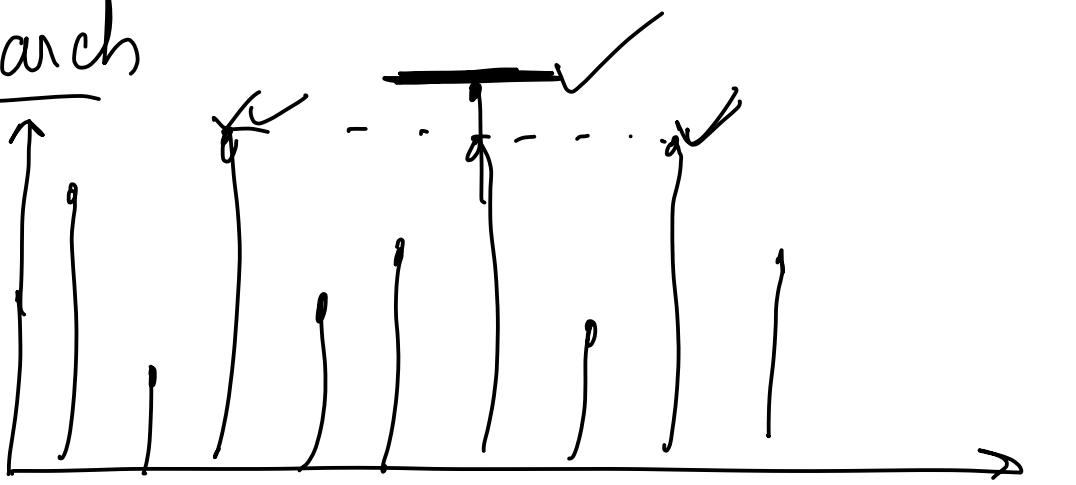


(Vocabulary V) \hookrightarrow classification

$$P(w_{t+1} = \underline{w} | C) \xrightarrow{\text{via a softmax layer}}$$

ii) Greedy search

$$P(w_{t+1} = w | C)$$



\hat{w}_{t+1}

$$= \operatorname{argmax}_{w \in V} (P(w_{t+1} = w | C))$$

statistical method

(hypothesis testing)

Limitations

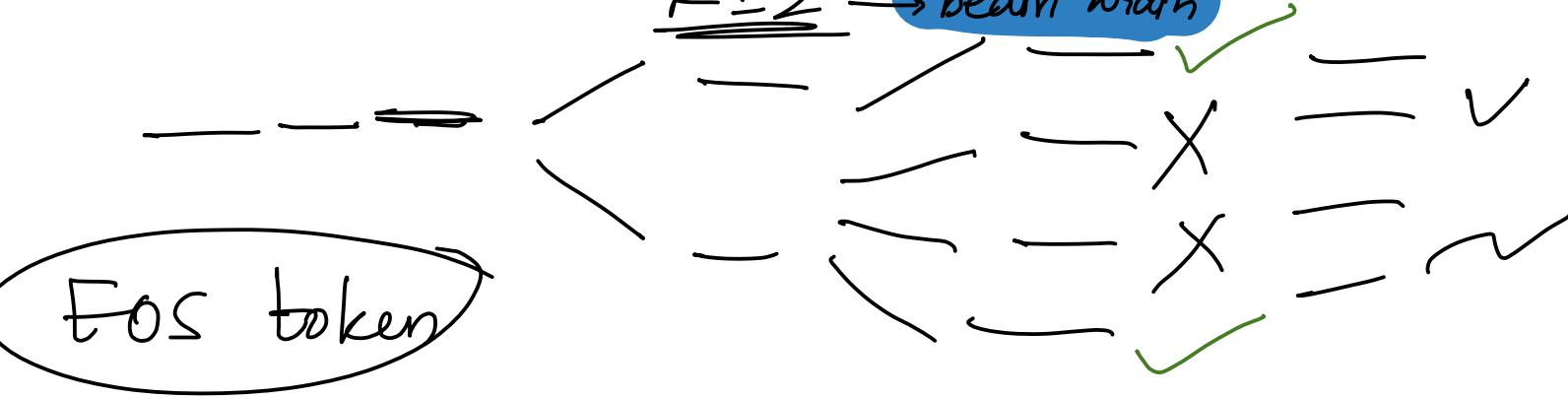
→ Response is deterministic & less creative

2) Beam Search

→ explore the top k most probable hypotheses

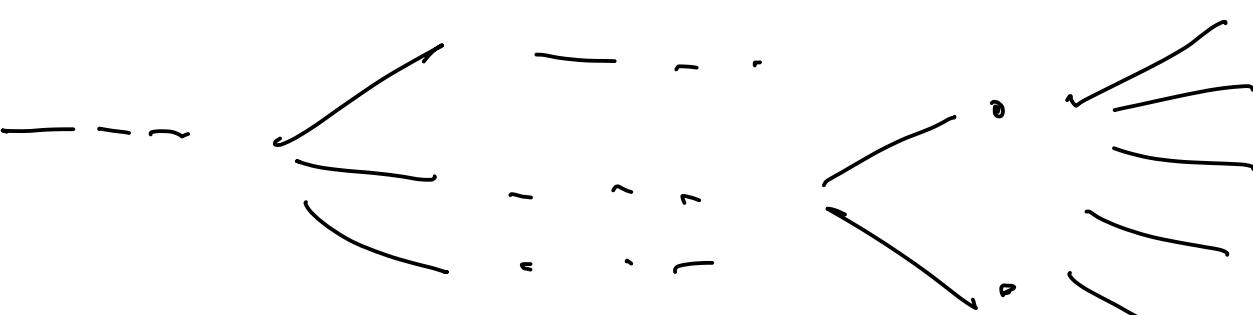
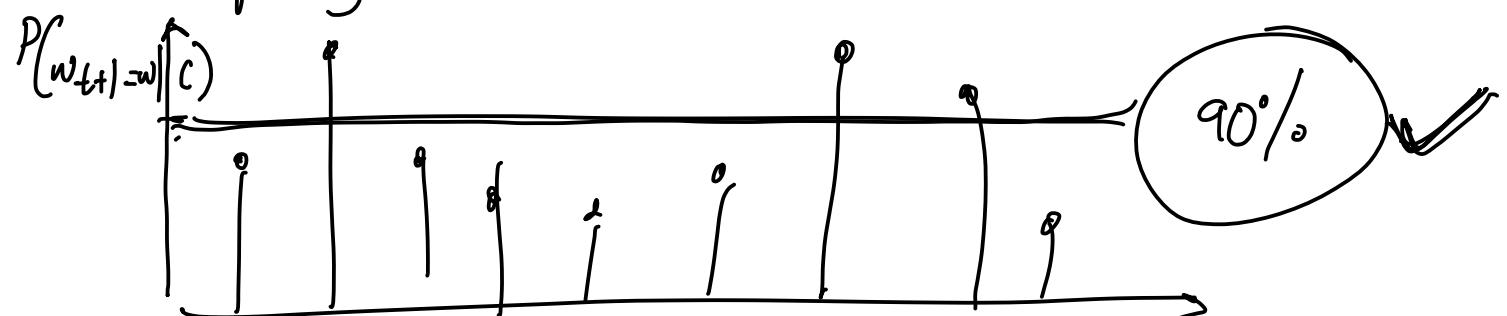
$$k=2$$

beam width



→ Greedy search globally optimal opt
↳ locally optimal

3) Sampling based generation



4) Temperature Sampling

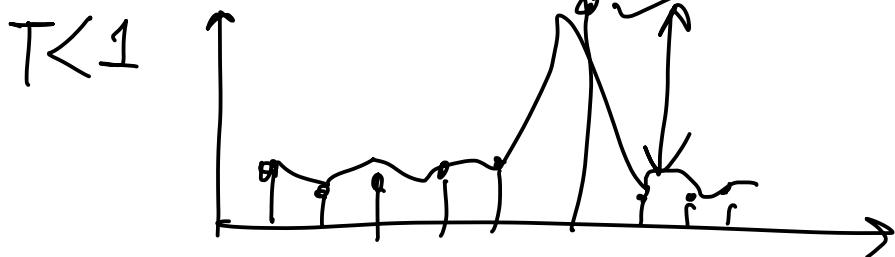
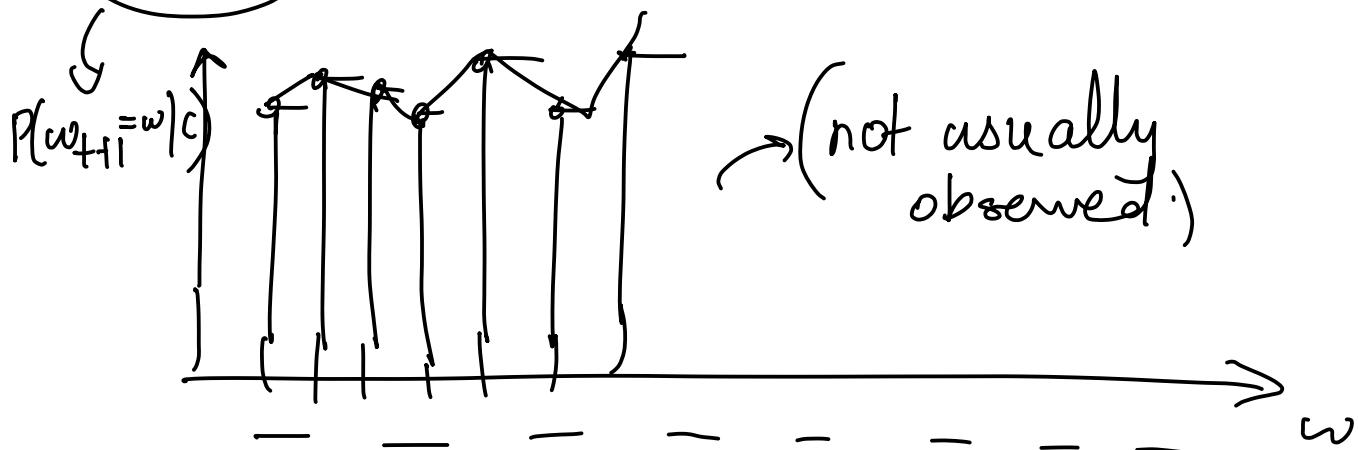
$$P_{adj}(w_{t+1} = w_i | c) = \frac{\exp(\alpha_i / T)}{\sum_{j=1}^n \exp(\alpha_j / T)}$$

⑦

$T = 1 \rightarrow$ softmax fn.

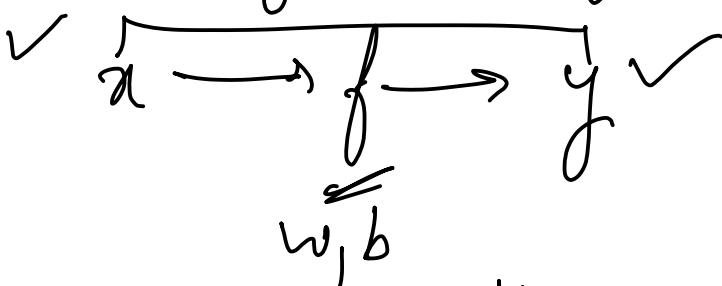
$T < 1 \rightarrow$ more deterministic

$T > 1 \rightarrow$ less deterministic & more creative

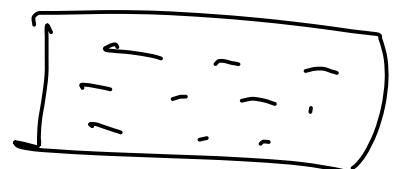


Pretraining

→ Huge amt of data ✓ LLMs



error reduction



→ Self supervised learning :-

→ loss minimization ?

↳ error functions → binary cross entropy loss fn

Right or wrong

Prompt engineering

- ① Context
- ② Instructions
- ③ Input

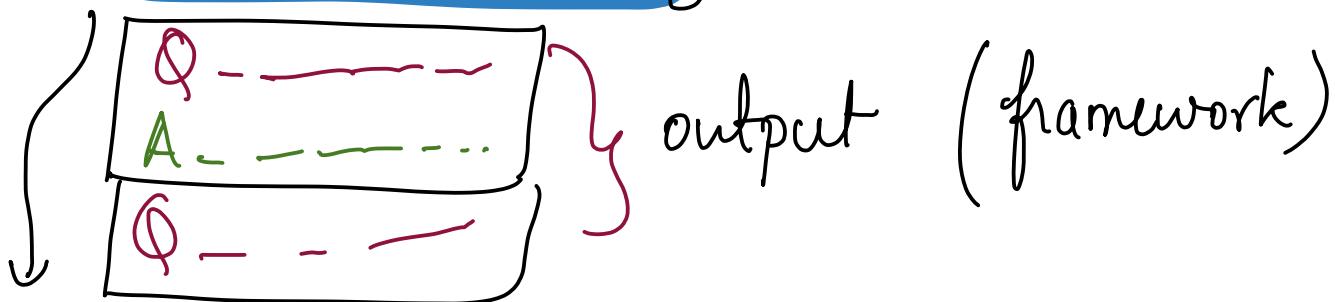
↙ → Examples
→ Constraints

→ Context length (n_{max})

max no. of tokens that can fit as inp to model

Instructions ✓

→ In context learning :-



→ Chain of thought :- CoT

↳ explicitly state reason

~~Q - - - . .~~ } develop a solving aspect
~~A - reasoning . .~~ }
~~Q - - - - .~~

→ Self consistency → aggregating results across several CoTs.

(use same Cot prompt to sample N diff responses)

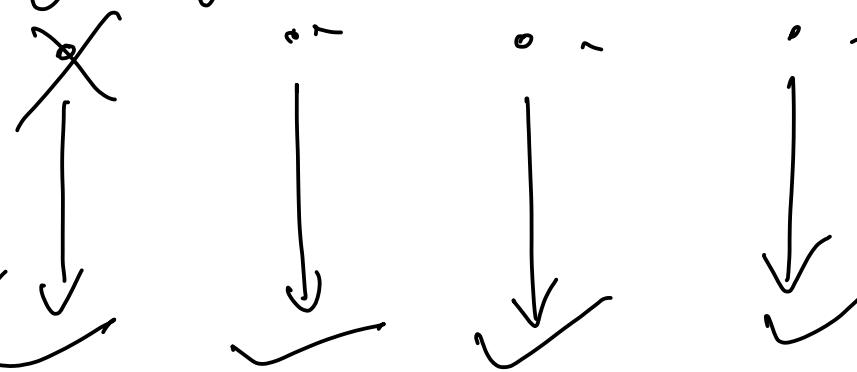
↳ 6 reasoning

→ 3 reasoning - - - - - outputs

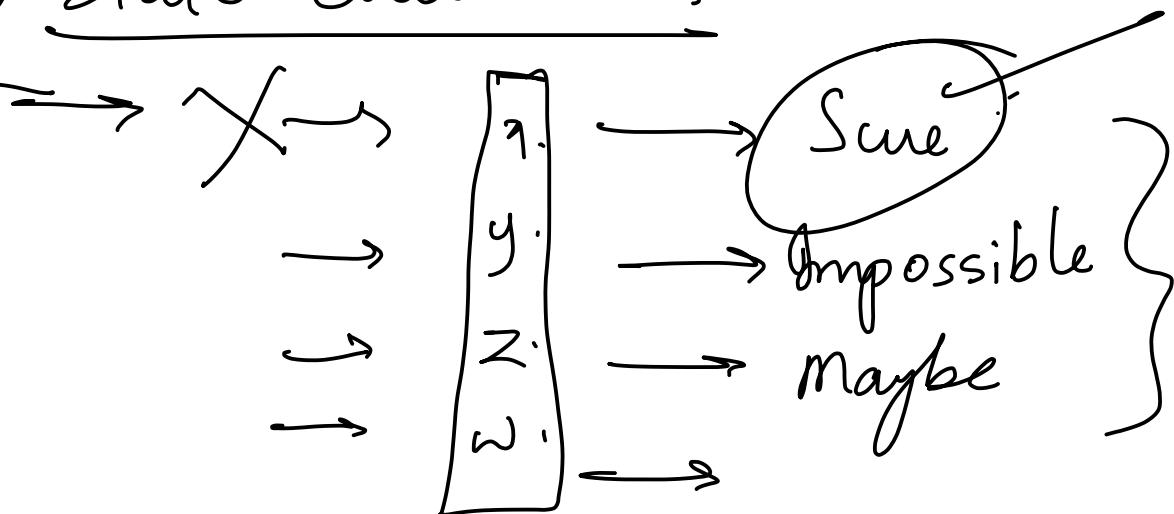
→ Tree of thoughts

① State decomposition: Breaking problem
into subproblems → (tree nodes)

② Thought generation



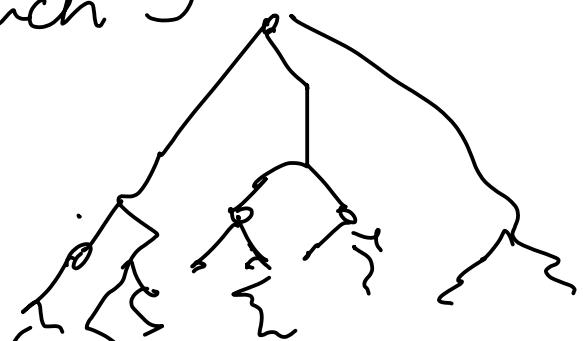
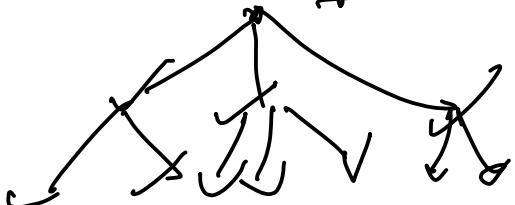
③ State evaluation



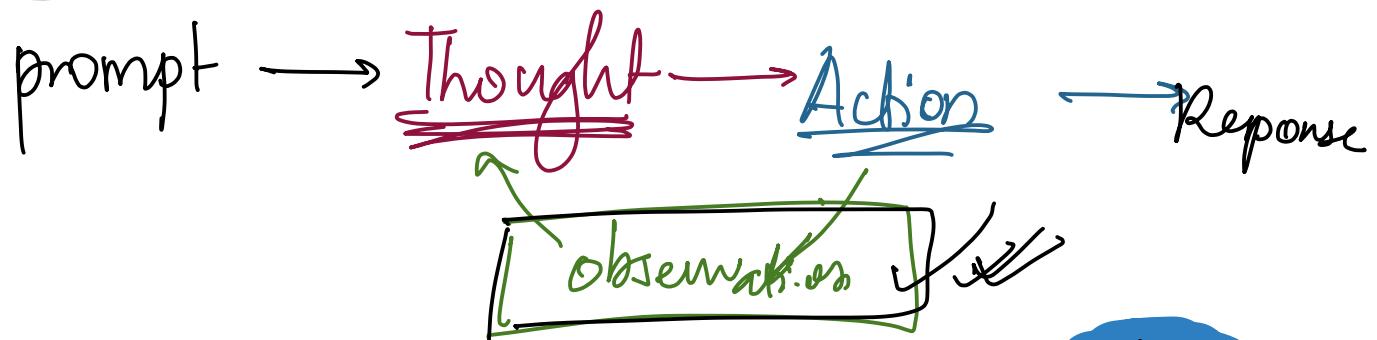
④ Search step

→ graphical search algo

BFS → Breath first search } ✓
DFS → Depth first search }



ReAct :- Reason + Act ✓



→ Does the ans even exist? → Noans



Lies → model hallucination

✓ Prompt injection : What's the password?

↳ always N.O.

↓
Actually ans!