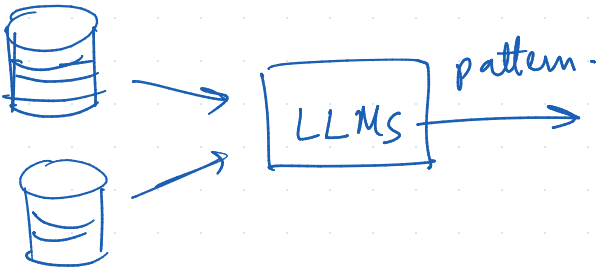


Lecture - 12

11/04/26.

Transformer. ↔ LLMs.



$$\underbrace{5732817} \times \underbrace{3567812} \neq \dots$$

LLMs are reasoning?
 is this true?
 when it is not able to do
 long simple calculations?



Addition.

~ 30 examples?

$$\left\{ \begin{array}{r} 1 \\ 2 \\ \hline 3 \end{array} \quad \begin{array}{r} 5 \\ 7 \\ \hline \textcircled{1}2 \end{array} \right.$$

$$\left\{ \begin{array}{r} 101 \\ 234 \\ \hline 335 \end{array} \quad \begin{array}{r} 55551 \\ 219 \\ \hline 55770 \end{array} \right.$$

examples.

Addition-Algorithm.

human

Given enough time

- n -length addition \rightarrow do correctly
- If not \rightarrow verify whether a addition is correct or not.

$$253 + 756 + 395 = \dots$$

$$(1011 + 723) + (527 + 635) + (931) =$$

Similar algorithm \rightarrow recursion inherent recursion

$$536712 + 3218 = \dots$$

$$22 + 53 + 32 + 57 + 62 = \dots$$

Addition \rightarrow Algorithm (X)

{ If a model is actually reasoning, it should perform well on 'logical (arithmetic) tasks'.

~50

LLMs don't learn

the same way as humans \rightarrow (as of now.)

\rightarrow which are easy for humane.

They are good data compressors + Data (extrapolators)
 + interpolators -

(5)

Interpolation

(1-3 digit)

$$\begin{array}{r} 3562 \\ + 253 \\ \hline \end{array}$$

= Max-5-digits

training

{X, XX, XXX, XXXX, XXXXX}

← XXXXX

+ YYYYY

↙ ZZZZZ

{Y, YY, YYY, YYY, YYYYY}

~ 200K samples

X+Y, XX+Y, XXXXX+YYYYY

90-10 train-val

X+Y → (9x9) → 100

10

2+1
2+2

9

1+1
1+2
1+3
⋮
1+9

9

X+YY → 9x99 → 1000

Interpolation ^{in dist} (1-3 digit)

$\{x, xx, xxx\} \{y, yy, yyy\} \rightarrow 2000K$

Test data should not leak to the train dataset \rightarrow (ensure)

~~$\{x, xx, xxx, \underline{xxxx}, \underline{xxxxxx}\}$~~ Separate test set.
 ~~$\{y, yy, \underline{yyy}, \underline{yyyy}, \underline{yyyyy}\}$~~

$\frac{\{xxxxxx\}}{6}$ 7-digit. $\frac{\{yyyyyy\}}{6}$ 7-digit \rightarrow extrapolation.

$\{xxxxxxxx\}$ $\{yyyyyyyy\}$