input $\longrightarrow$ word

input + word $\longrightarrow$ word

input word word $\longrightarrow$ word

Questi

Answer

Question

text
↓
tokens ⎤
↓     |
1-hot | depend on input size
↓     |
encoder ⎦

↓
ANN ← Paner, most of work
↓     size ~ energy
decoder ← [0, 0, 0.7, 0, 0.3, 0.2, ... ]
↓                    ↓      ↓
wave

Output size ~ energy

Energy used by American $\rightsquigarrow$ 140 LLM tokens per second

Question

How can I answer

LLM

human-written
split opn
Algorithm

Opion 1

Opm 2

Which of those best ds Q? res 1, res 2, ...

rest

user or action

## Questions

Agentic vs Generative

next few years