

- 1) Decoder ~~is~~, Transformers are more scalable for LLM's because they are single transformer with self-attention (causal) which directly matches the next token prediction objective used in large-scale language modeling.

This avoids extra encoder stack and cross attention required in encoder-decoder architectures.

- 2) Next-Token prediction is sufficient because many NLP tasks can be expressed as predicting the next word given a context.

When trained on massive text, the model learns the structure of language

To correctly predict the next token, the model must learn meaning and logical dependencies.

A single-next token objective acts as a universal ~~is~~ training signal that enables the emergence of complex language abilities.