

RAG (Retrieval-Augmented Generation)

- is a method in LLMs to access external knowledge (Real time) rather than relying on only what they were trained on
- it does it by combining search (retrieval) with language gen
- Why: after training cut off date LLMs don't know anything and can hallucinate facts. RAG solves this by
 - retrieving relevant info from external knowledge base
 - generating answer using that retrieved info (via LLM) (e.g. GPT4)

A Rag chain is those structured sequence of steps a RAG

pipeline uses. Ex: input \rightarrow Embedder \rightarrow Retriever \rightarrow Filter \rightarrow LLM \rightarrow Output

involves using a pre-trained model on specific dataset to adapt it to a particular task making it better for that task but limited to training knowledge, hence RAG

Scaling Laws for LLM

- describe how performance improves as model size (params), data set and compute (FLOPs) change
- More compute \rightarrow Bigger model + more data \rightarrow Better performance (upto a point)
- loss (perplexity) decreases as log-linearly as model size \uparrow , dataset size \uparrow and compute \uparrow until diminishing returns or overfitting
- Key points/laws
 - more params = Better capacity to learn complex patterns
 - more data = Prevents overfitting, improves generalization
 - more compute = enables longer training, more updates
 - more training steps = more updates = more learning (upto a point)

Balance (Optimal ratio between)

- model size • data size • training compute
- if model too big for dataset = overfit
- too much data for small model = underfit

Limits

- eventually you hit diminishing returns
- need better architecture (transformers, MoE)
- alignment & fine tuning (RLHF, RAG)

- This Limit of LLM is getting better rapidly then slowing down \rightarrow in the "Bending curve" and Sol is mentioned above

