(LLM notes) **RAG** (Retrival -Augmented generation)
(Also see 106)

- is a method in LLMs to access external knowlege (in Real time)
  rather than relying on only what they where trained on

  it dose it by combining search (critrical) with language gen

- why: after training cut off date LLMs dont know anything
  and can hallucinate facts Rag solves this by
  • retriving relevant into from external knowlage base → vector simmilarity from DB's
  • generating answer using that retrived into (via LLMs eg GPT4)

A <u>Rag chain</u> is those structured sequence of sheps a RAG
pipeline uses Ex: input→ Embedder → Retriver → filter → LLM → Output

Fine tuning: 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 knolage, hence RAG

[(See page 104) ⇒ **Scaling Laws for LLM**

- describe how preformance improves as model size (prams),
  the data set and compute (FLOPs) change

- More compute → Bigger model + more data → Better Preformance (upto a point)

- loss (perplexity) decreses as log-linearly as model size ↑,
  dahset size ↑ and compute ↑ until diminishing returns or overfitting

- key points/laws
  • more params = better capacity to learn complex patterns
  • more Data = Prevents overfitting, improves generalization
  • more training Shps = more updates = more learning (upto a point)
  • more compute = enables longer training, more updates

**Balance** (Optimal ration Between)
• model size • data sizes • training compute
- if model too big for Dataset = overfit
- too much data for small model = underfit
- this Limit of LLM ie getting better rapidly then slowing down →
  is the "Bending curve" and Sol is mentioned above

**Limits**
• eventualy you hit diminishing returns
- need better architecture (transformers, MoE)
- alignment + fine tuning (RLHF, RAG)



74