# Normalization in LLMs

- Normalization in LLMs is typically applied between the attention and feed forward layers (and vise versa) its usually applied before or after each sub layer (depending on variants):
  - Post norm: normalization after Attention/MLP (used in original transformer paper)
  - pre norm: normalization before Attention/MLP (used in modern LLMs as it stabalizer training in Deep models)

Why? :
  - Prevents exploding/vanishing activations: means the numbers (output of layers) dont become too big or small as they pass through many layers - this keeps the models computation stable

  - Makes optimizations smoother: means the model learns more steadily during training, the loss surface becomes less chaotic, so gradient decent can take more reliable steps towards better performance insted of bouncing around/getting stuck

  - Allows deeper, more stable transformer training: beacuse norm keeps each layers well behaved (balanced and stable not too big or small consistant across layers) you can stack more layers without the models training becoming unstable or diverging, in short normalization makes it possible for very deep NN like LLMs to still learn effectively and converge smoothly, insted of breaking down as they grow in size

(for ray see 94, 93, 106)       Vector Databases Vs knowlage graph and usage in RAG
(for verctor DB see 95)

- Knowlage graphs: are a structured representation of facts and relationships, it stores entities (nodes) and relationships (edges) → like a map EX) (Elon) —[founded]→ (Tesla)
  of how concepts connect, each entity has attributes       (Tesla) —[Makes]→ (EV·s)
  and connections, allowing resoning, merging and context linking       ⤳node ⤳edge

- Vector DB: store embeddings (numeric representations of text, imgs etc) it allows semantic search by simmilarity (cos, Euclidian). Ex: • you embed a document into a 1536-dim vector
  • you query "who made tesla". DB finds most simantically simmiller chunks (even it words diffre)

| Feature | knowlage graph | Vector Databases |
|---------|---------------|------------------|
| Structure | Explicit (nodes, edges, relations) | implicit (numeric emb) |
| Query type | Symbolic/logical (graph traversal) | Semantic simmilarity serch |
| Strength | Precise resoning and explainibility | Flexible meaning based ret |
| Data type | structured Facts | Unstructured media, text |
| Ex | Neo 4s, AWS Neptune | Pinecone, Weavite, FAISS |

- IN RAG: you can use either or Both, in classic RAG (vector DB) you retreive semantically simmiller chunks from embeddings and feed it to LLM, in KG-RAG you retrieve related entites/facts via graph traversal and or Augmen ted emb (embeddings enhanced with additional context or metadata) with graph links
Ex: use KG to find connected info ("founder of Ex company) then vector DB to fetch text about them or use Hybird RAG: use vetor DB-semantic recall, KG for logic, resoning, structure

Ex) KG: Structerd brain of facts, used in search engines
     VectorDB: memory of meanings, used in RAG pipelines
- They complement each other, not replace ments