

LLM Programs

LSTM: Long short term memory (LSTM) Network/layers are a type of RNN (Recurrent Neural networks) that is capable of learning long dependencies in sequential data, BiLSTM: is Bidirectional LSTM it processes the input sequence in both forward and backward directions, capturing context from both past and future tokens useful in NLP where understanding context of word depends on both preceding and succeeding words (alt is using transformer as it want to capture global context and large dataset)

- Short term memory (STM) is the facts passed in as part of the prompt like user name, rage etc.
- Long term memory (LTM) in A LLM (more specifically Agents)

Stores important information for the future it saves facts

like past events, preferences and learned skills for better

long term performance, unlike STM it lives beyond session

- Other Architecture patterns: We have seen RAG, React, COT
but other Agent patterns are: Planner Executor (splits into planner and executor to divide work), DAG Agents or Direct asyclic graph where each node is a agent does a task and passes it to the next agent with no cycles (one way), TOT (Tree of thought) is a way to organize agents reasoning the root is the main problem and it branches out creating ideas/steps to lead to a sol.

- Vector Database: when implementing something like a RAG or Agents Vector DB stores and retrieves embeddings efficiently, these embeddings are vector representation of data like documents, text, images etc during query this system converts it to a embedding and searches the vector DB for the most similar/relevant embedding (Ex, related documents/snippets) it does this using methods like ANN to find similar data points in high dim space. this retrieved info from DB is then fed to something like LLM.