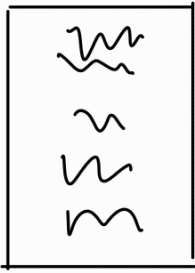
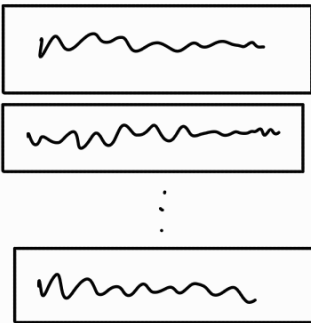


PDF

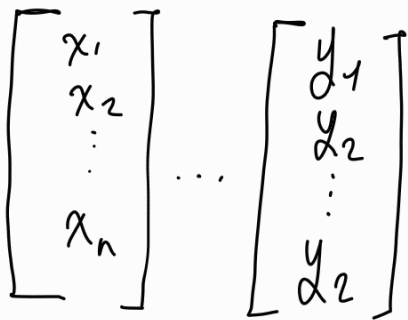


extract the text into a single string

split the string into chunks

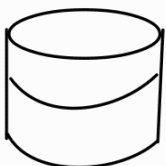


↓ embeddings*



↓ ②

storing in a vector database

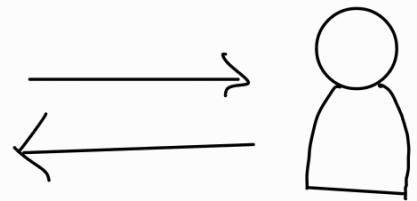


① In our Case we take the ① path, the other one is for create a database of Knowledge with more information

* embedding Can be made using different models, with different approaches

* In a production environment it's recommended to separate the LLM execution and the interactive chat, depending on the number of simultaneous users. That's why use Openai, Anthropic, etc Api endpoint could be a better option.

Create Streamlit app



Pass the info to an Open Source LLM (Meta-Llama-3)

Future improvements:

- * fine tuning the LLM with Questions and answers
- * Generate different LLM-answers and use Cosine Similarity with the embedding Question to find the best answer.