



SEC 10K-
Filings



Preprocessed
Documents

-> Removing
HTML-Attributes,
Empty HTML Tags,
XBLR Tags



Store Full
Documents

-> Right now on Disk,
but maybe Database?



Split Documents

-> Recursive Character
Text Splitter based on
HTML (Chunk Length =
6,000 chars; Overlap =
10 %)



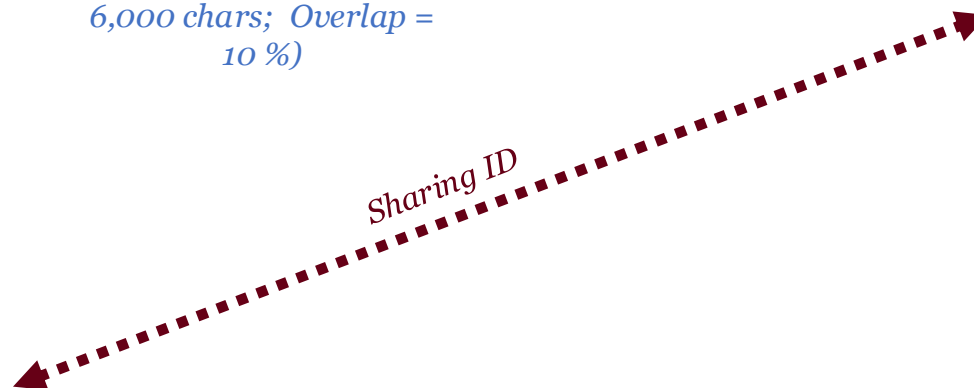
Embeddings

-> Snowflake Arctic-embed-
m-v1.5



Qdrant
Collection

-> Local Docker
Instance



Sharing ID

