## **Indexing Embedding** Query **Chunking Input:** query from the user for which the domain are the The uploaded documents. document is split into chunks. These chunks are E.g) What is the current buy then converted to price of the described stock? vectors via Each vector is then *embedded*. embedding. into a Vector Database. Generation Retrieval User The LLM combines the final ordered relevant documents **Output** and generates an output. The LLM returns generated output in the form of a comprehensible response based on the retrieval of relevant chunks, re-ranking, and the combination of those chunks. Similarity Search is used on Reranking the Vector DB and embeddings E.g) The current buy price of the stock within to *retrieve* chunks (split for XYZ Corp is \$125.75. Please note After relevant that stock prices fluctuate constantly documents) that are most chunks are retrieved. due to market conditions, so it's **relevant** to the inputted query. important to check a reliable financial they are *re-ranked* news source or brokerage platform for based on relevance the most up-to-date information. to query