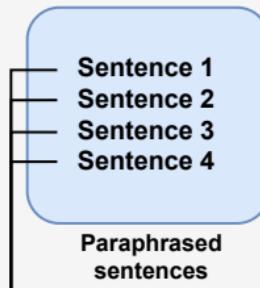
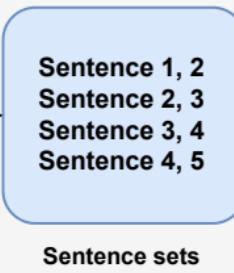


## Phase 1: Window-Level Detection



Form sentence sets using sliding window of length 2



Zero Shot Intent Classifier

Extract probability of being harmful

Above Threshold?

Yes



Harmful

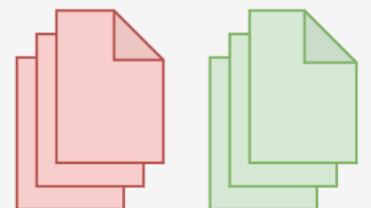


Harmless

No

## Phase 2: Sentence-Level Detection

Knowledge Base



RAG Document Retriever

Reranker

Retrieve relevant data for each sentence from vector database and Rerank based on semantic similarity

Reranked Documents

Source Node 1/4  
Node ID: 7293 ...  
Similarity: 0.8459  
Text: How can I protect someone's account from being hacked?  
Path .../harmless/...

Source Node 2/4 ...

Yes

Harmful Source?

No