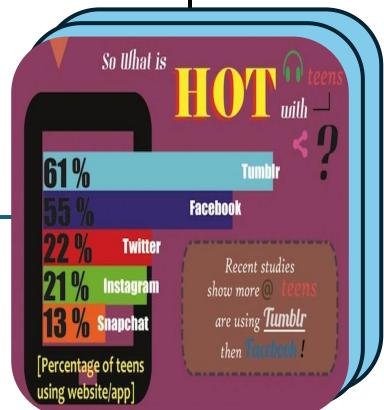


RAG



Parsing & Chunking

So What is  
HOT teens  
with.....13  
% ...  
...Snapcha  
t .....

Document Database

LLM-based  
Retriever

LLM-based  
Retriever

So What is  
HOT teens  
with.....13  
% ...  
...Snapcha  
t .....

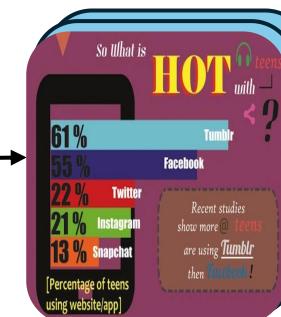
twitter

LLM-based  
Generator

Which is the fifth  
most favorite  
application of  
Teenagers?

VLM-based  
Retriever

VLM-based  
Retriever



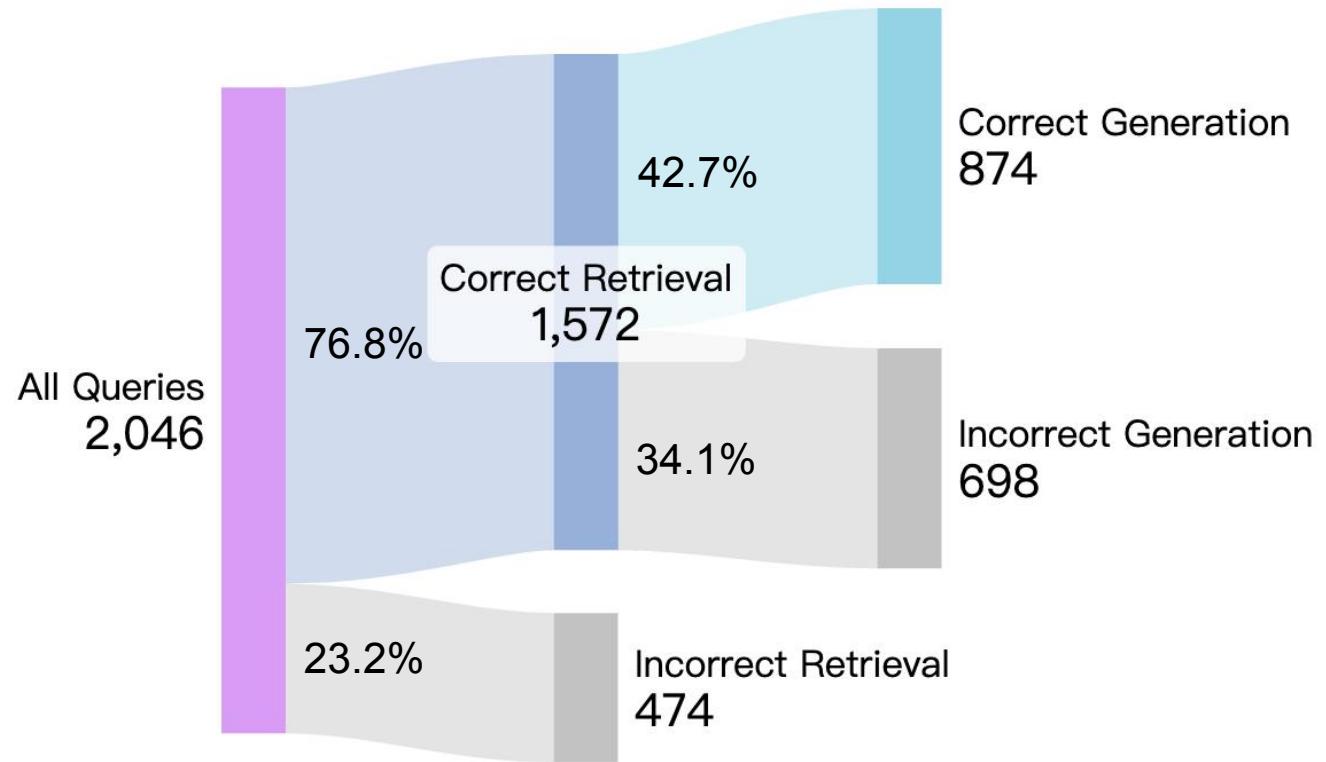
VLM-based  
Generator

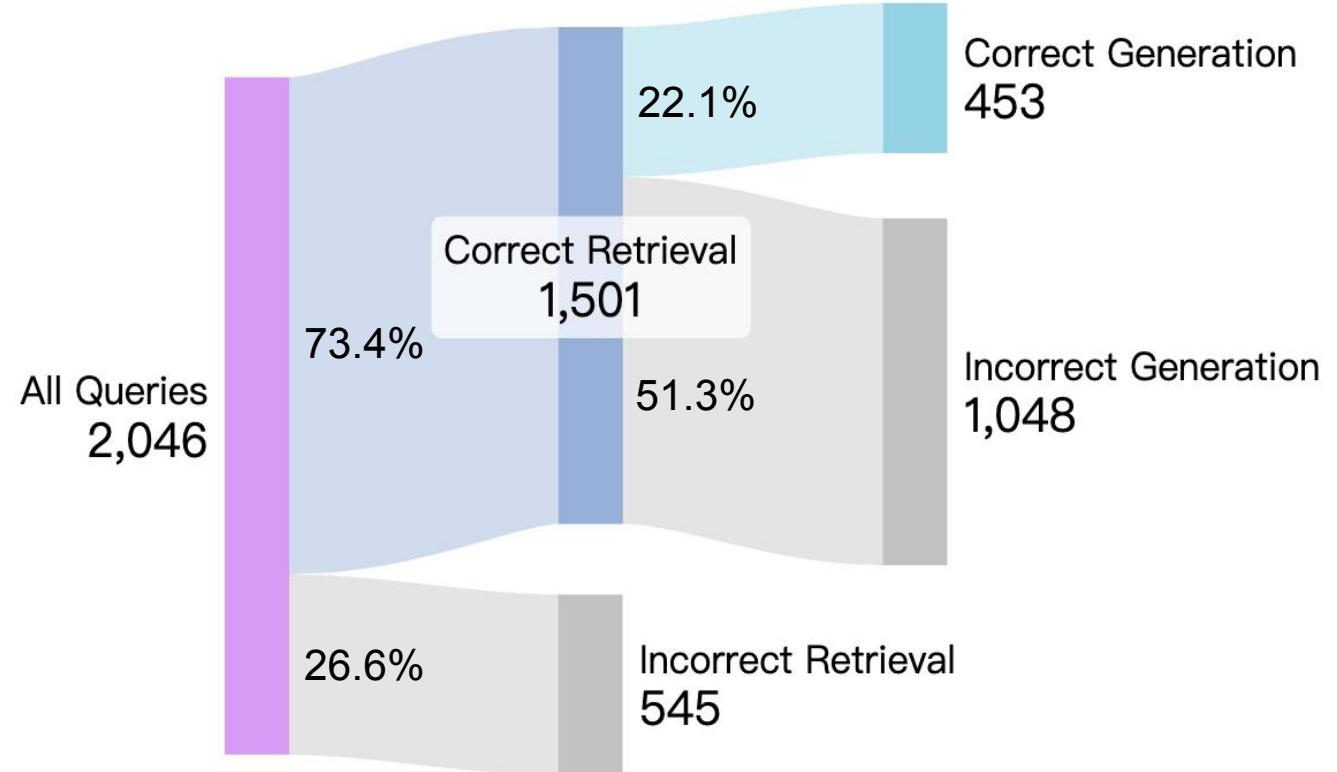
snapshot

RAG-V

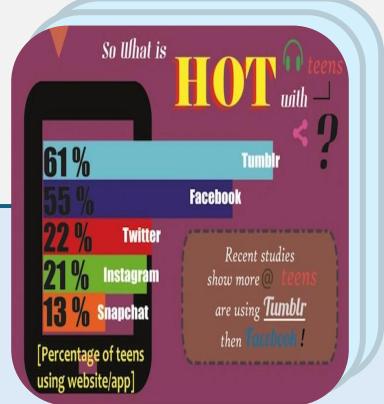
Document Embedding Inference

RAG Pipeline Inference





RAG



Parsing & Chunking

So What is  
HOT teens  
with.....13  
% ...  
...Snapcha  
t .....

LLM-based  
Retriever

Vector Database

LLM-based  
Retriever

So What is  
HOT teens  
with.....13  
% ...  
...Snapcha  
t .....

twitter

LLM-based  
Generator

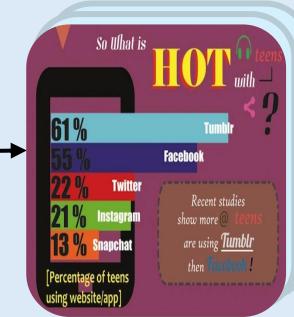
Which is the fifth  
most favorite  
application of  
Teenagers?

VLM-based  
Retriever



Vector Database

VLM-based  
Retriever



snapshot

VLM-based  
Generator

RAG-V

Document Embedding Inference

RAG Pipeline Inference