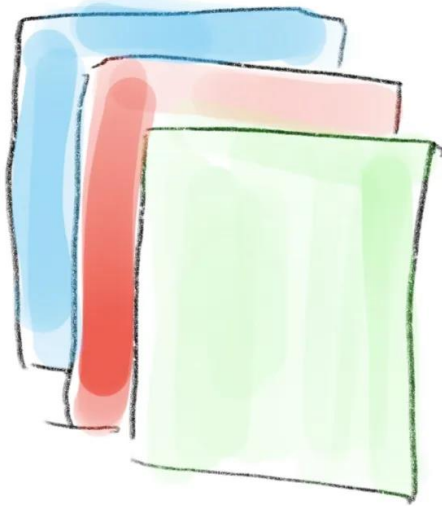
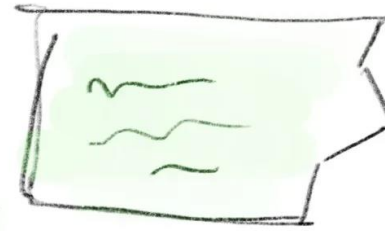


Corpus of Documents



User input

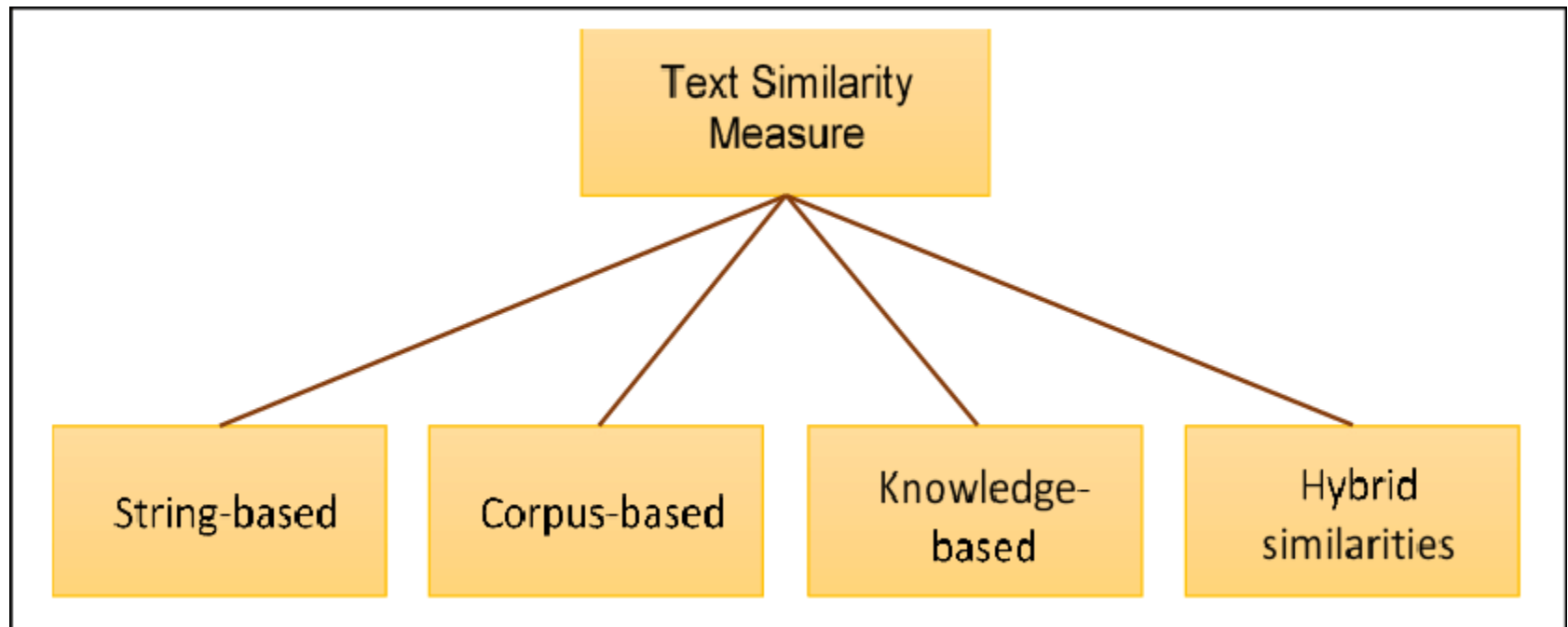


User

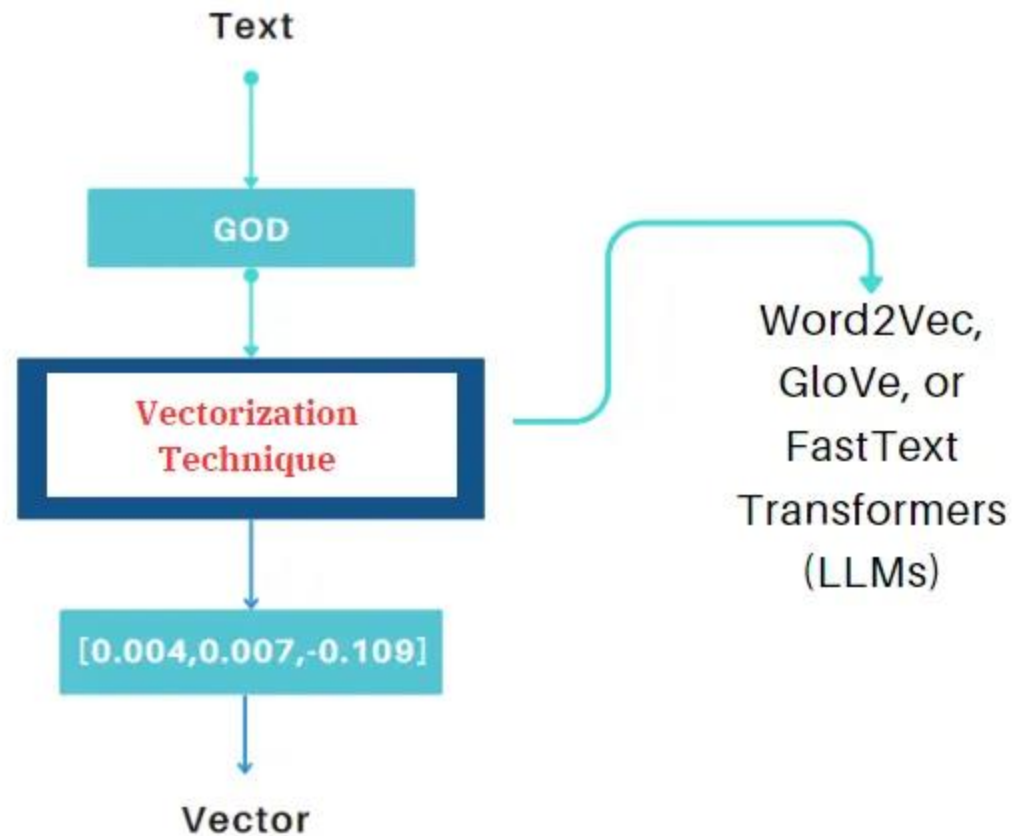


Check Similarity

Return relevant documents

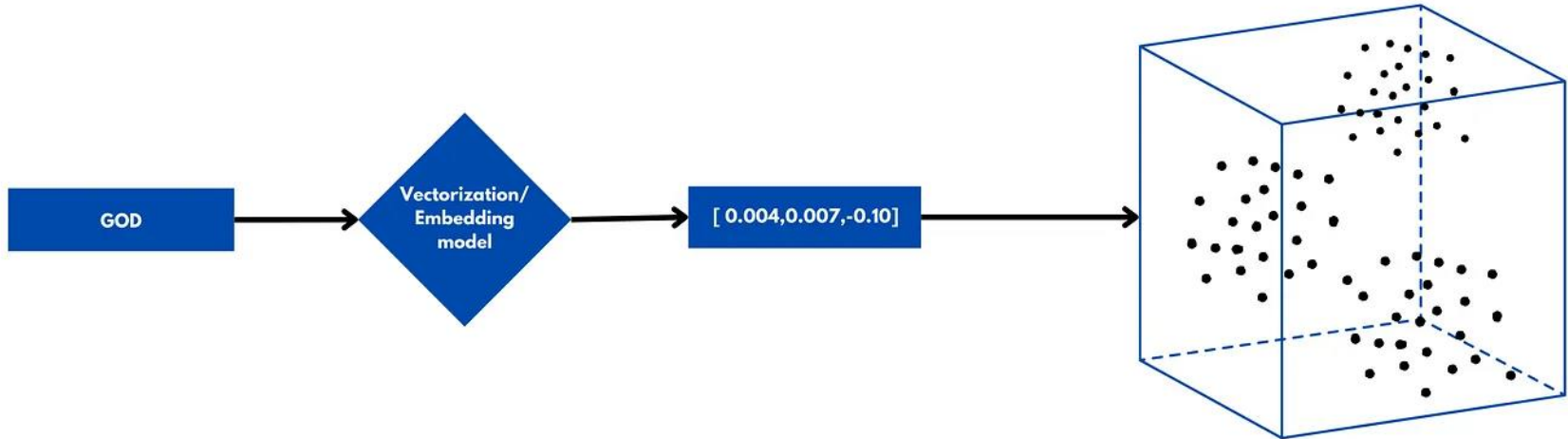


Step 1

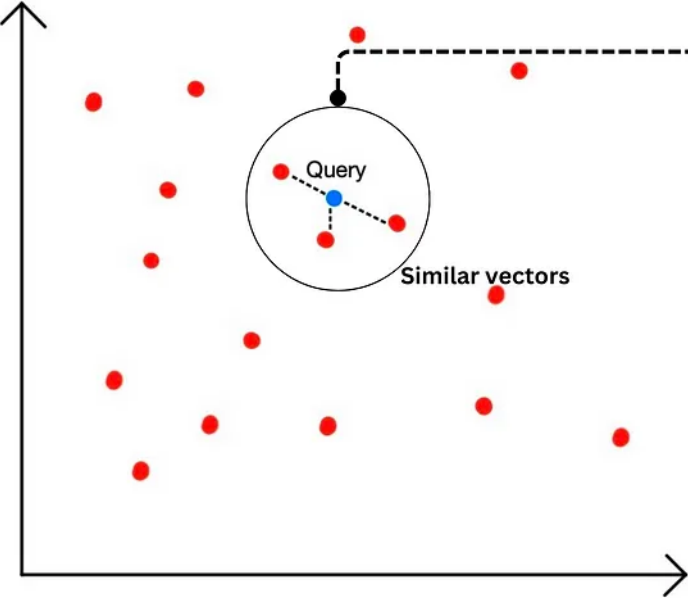


Step 2

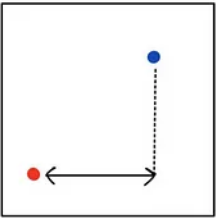
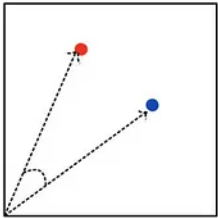
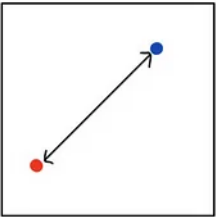
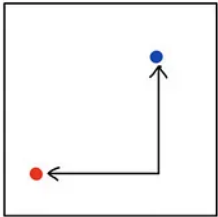
Text → Vectors → Vector DB



Step 3



3 nearest neighbors for a query in a vector space



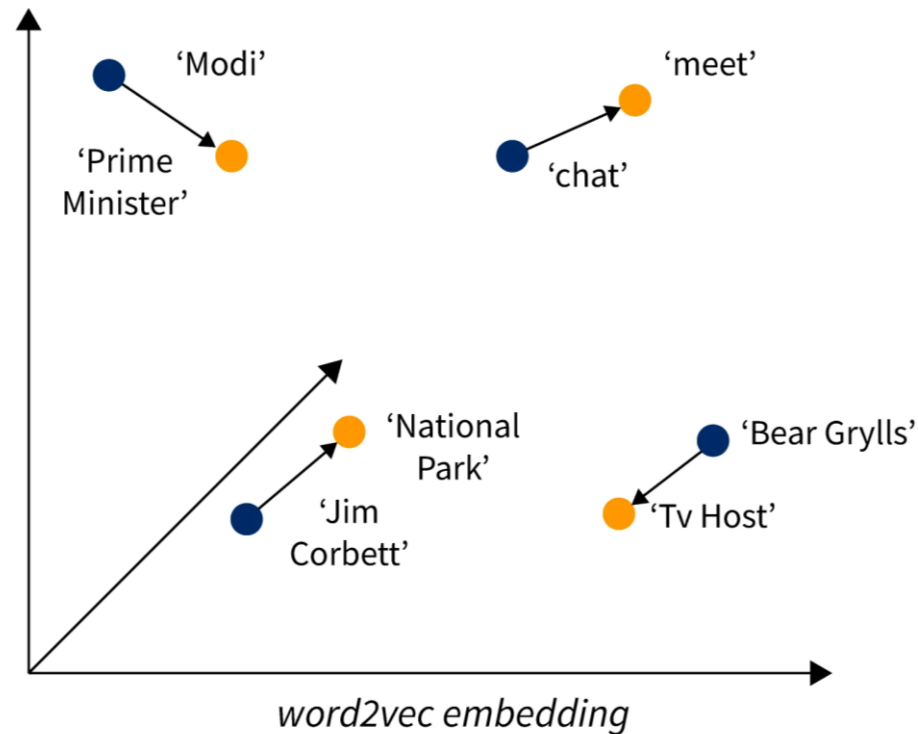
Similarity Measurement Techniques

Document 1

Modi *had a chat with*
Bear Grylls
in **Jim Corbett**

Document 2

The **prime minister**
meets the
TV Host in a National Park



You should have at least 8 GB of RAM available to run the 7B models, 16 GB to run the 13B models, and 32 GB to run the 33B models.

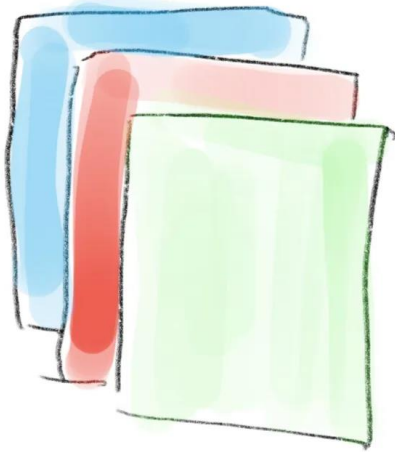
<https://ollama.com/download>

```
(metaverse) c:\Python3\my_project\project_llm>ollama run llama3.2
```

Challenge:

Semantics
of user input

Corpus



Check
Similarity

Return relevant
documents

I don't like to
like

