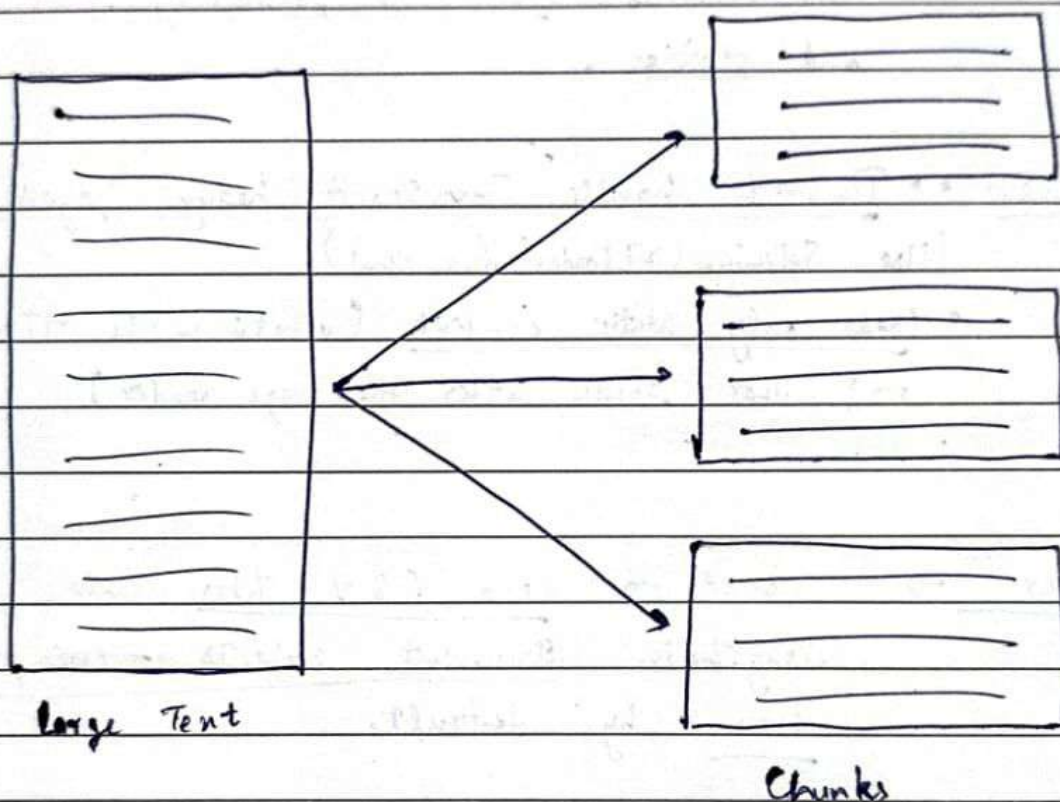


B. TEXT SPLITTER →

Text splitting is the process of breaking large chunks of text (like articles, PDFs, HTML pages, or books) into smaller, manageable pieces (chunks) that an LLM can handle effectively.



• Overcoming Model Limitations → Many ~~limitations~~ embedding models and language models have maximum input size constraints. Splitting allows us to process documents that otherwise exceed these limits.

• Downstream tasks → Text splitting improves nearly every LLM powered tasks.

Task	Why splitting helps
Embedding	Short chunks yield more accurate vectors
Semantic Search	Search results point to focused info, not noise
Summarization	Prevents hallucination and topic drift

• Optimizing Computational Resources → Working with smaller chunks of text can be more memory-efficient and allow for better parallelization of processing tasks.

Ex → docs =

docs[0]	docs[1]	docs[2]
len=48	len=150	len=100

if chunk size = 100

→ C1 → docs[0] (48 chars)
 C2 → docs[1] (100 chars)
 C3 → docs[1] (50 chars)
 C4 → docs[2] (100 chars)

if chunk size = 40

→ C1 → docs[0] (40 chars)
 C2 → docs[0] (8 chars)
 C3 → docs[1] (40 chars)
 C4 → docs[1] (40 chars)
 C5 → docs[1] (40 chars)
 C6 → docs[1] (30 chars)
 C7 → ...
 :
 :

```
from langchain.text_splitter import CharacterTextSplitter
```

```
text =
```

```

splitter = CharacterTextSplitter(
    chunk_size = 100,
    chunk_overlap = 10,
    separator =
)

```

```
result = splitter.split_text(text)
```

Date:

Mo Tu We Th Fr Sa Su

② Text-Structure Based Text Splitting →

It considers the fact that the text is organized as words → sentences → paragraphs.

Hence we use this technique

↓ called

Recursive Character Text Splitting

(Used a lot of times)

It divides paragraphs → lines → ~~words~~ words → characters.

Ex →

My name is Shivam (17)

I am 21 years old (17)

I live in Rohtak (16)

How are you (11)

chunk-size = 10

paragraph

My name is Shivam
I am 21 years old

I live in Rohtak
How are you

My name is Shivam

I am 21 years old

I live in Rohtak

How are you

my name is Shivam

my name

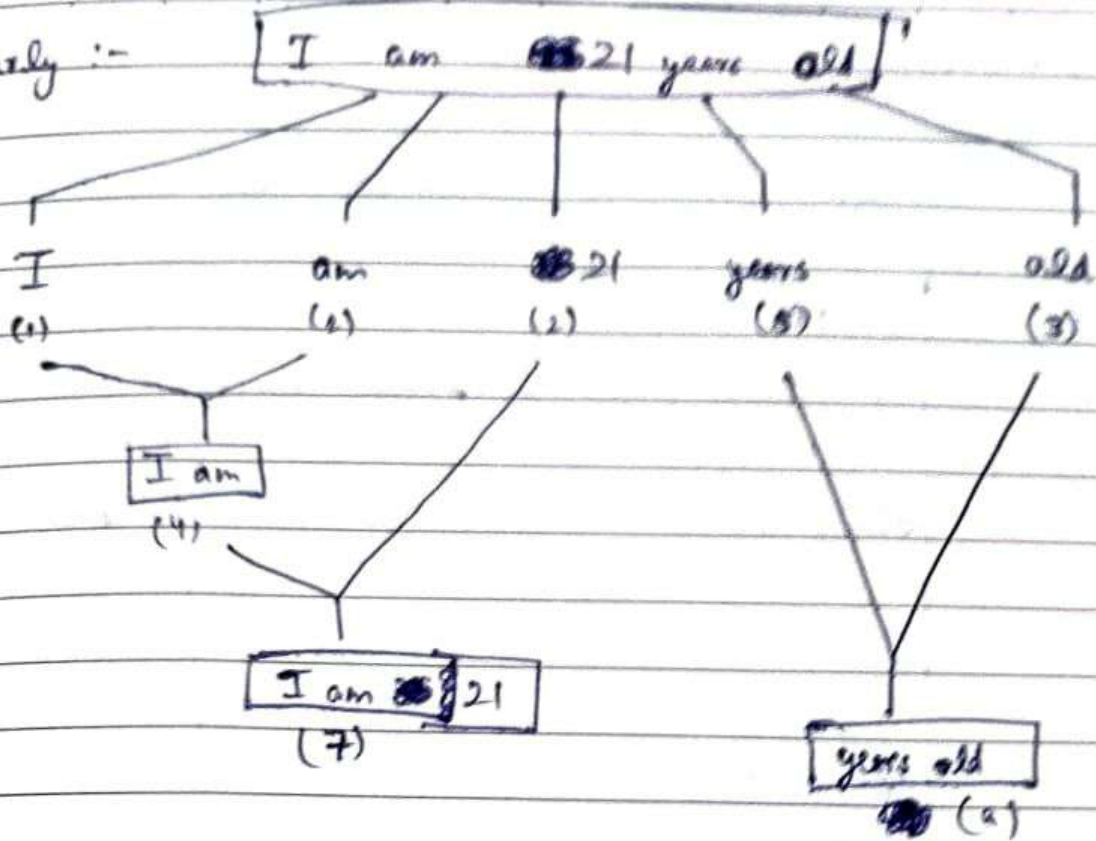
My name is
(10)

Shivam
(6)

Mo Tu We Th Fr Sa Su



Similarly :-



[My name is], [Shivam], [I am 21], [years old], [I live in], [Rohitak],
[How are], [you]

Date :

Mo Tu We Th Fr Sa Su

③ Document-Structure Based Text Splitter →

It is best for the documents that are structured.

Ex →

```
# Title
```

```
A simple project to manage time.
```

```
## Features
```

```
_____
```

```
_____
```

```
_____
```

```
## Tech-Stack
```

```
_____
```

```
_____
```

```
_____
```

markdown

or

code

```
def is_pass(marks):
```

```
    _____
```

```
    _____
```

```
    _____
```

```
def main():
```

```
    _____
```

```
    _____
```

```
    ✓
```



Here separators are:- "`\n`class", "`\n`def", "`\n`!def", etc.
& "`\n\n`", "`\n`", " ", ""

Here separators are: $\backslash n \# \{1, 6\}$, $\backslash n \| ^* \| ^* \| ^* + \backslash n$,
 $\backslash n - - - + \backslash n$, $\backslash n - - - + \backslash n$, &
 $\backslash n \backslash n$, $\backslash n$, $" "$, $" "$, etc.

It is just Recursive Character Text Splitter under the hood, but just with different separators.

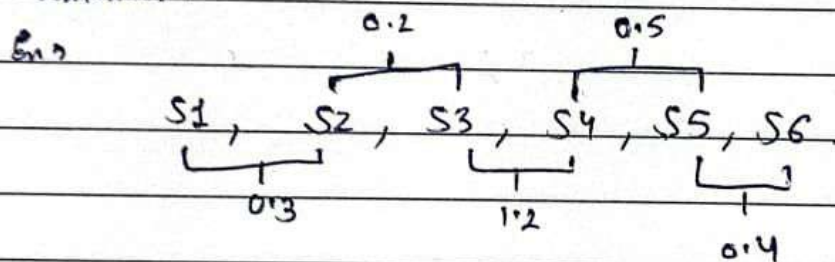
④ Semantic-Meaning Based Text Splitter →

when recursive character text splitter doesn't work, we use this splitter.

Ex → in same paragraph, two completely different sentences.

Hence, this splitter splits text based on its meanings.

It uses embedding models to find the semantic of sentences.



then checks their std. deviations



$[S1, S2, S3]$, $[S4, S5, S6]$

It is currently at experimental stage.