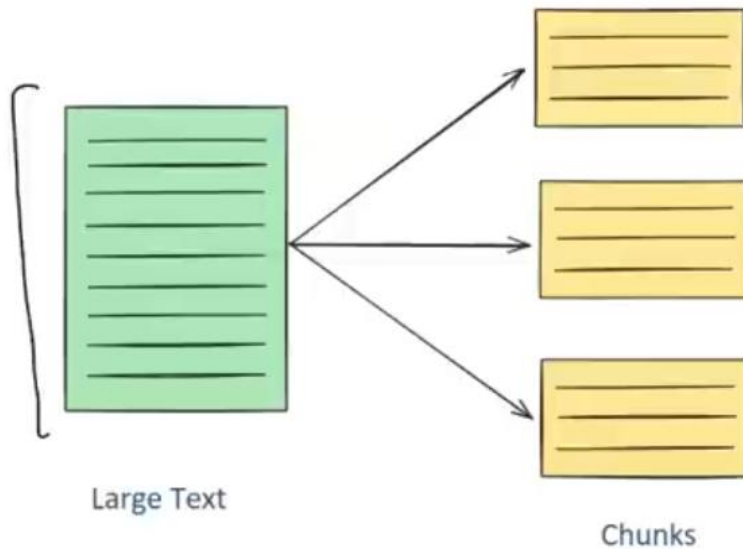


Text Splitting

01 April 2025 18:10

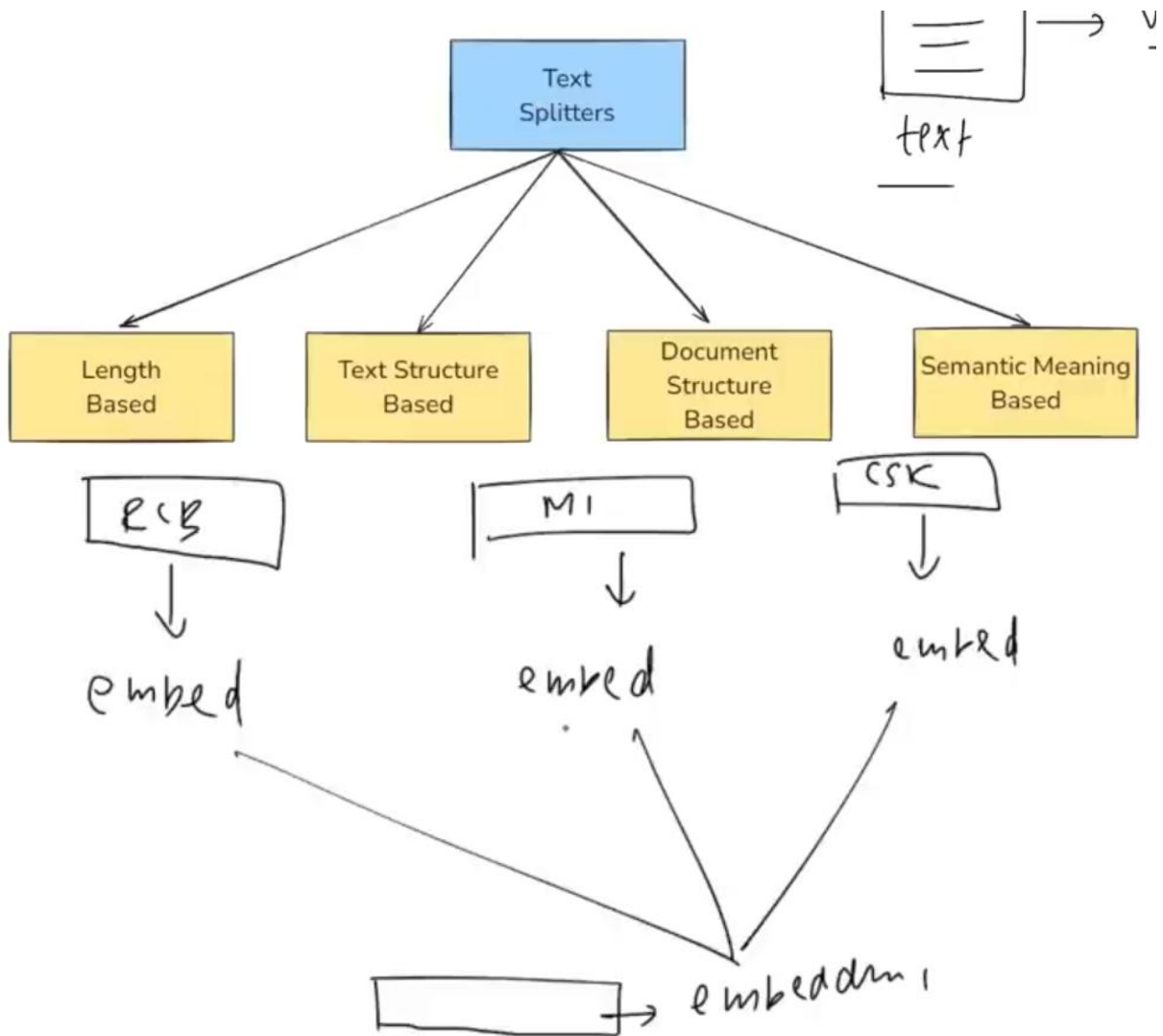
Text Splitting is the process of breaking large chunks of text (like articles, PDFs, HTML pages, or books) into smaller, manageable pieces (chunks) that an LLM can handle effectively.



- **Overcoming model limitations:** Many embedding models and language models have maximum input size constraints. Splitting allows us to process documents that would otherwise exceed these limits.
- **Downstream tasks - Text Splitting improves nearly every LLM powered task**

| 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.



1. Length Based Text Splitting

01 April 2025 18:10

Space exploration has led to incredible scientific discoveries. From landing on the Moon to exploring Mars, humanity continues to push the boundaries of what's possible beyond our planet.

These missions have not only expanded our knowledge of the universe but have also contributed to advancements in technology here on Earth. Satellite communications, GPS, and even certain medical imaging techniques trace their roots back to innovations driven by space programs.

Space exploration has led to incredible scientific discoveries. From landing on the Moon to explorin

g Mars, humanity continues to push the boundaries of what's possible beyond our planet. These missi

ons have not only expanded our knowledge of the universe but have also contributed to advancements in

n technology here on Earth. Satellite communications, GPS, and even certain medical imaging techniqu

es trace their roots back to innovations driven by space programs.

2. Text-Structured Based

01 April 2025 18:10

Structure

My name is Nitish
I am 35 years old

I live in Gurgaon
How are you

\n\n, \n, ' ', '+'

para

line

word

character

[chunk_size = 10]

My name is Nitish (17)
I am 35 years old (17)

I live in Gurgaon (17)
How are you (11)

My name is Nitish (17)/
I am 35 years old (17) X
(34)

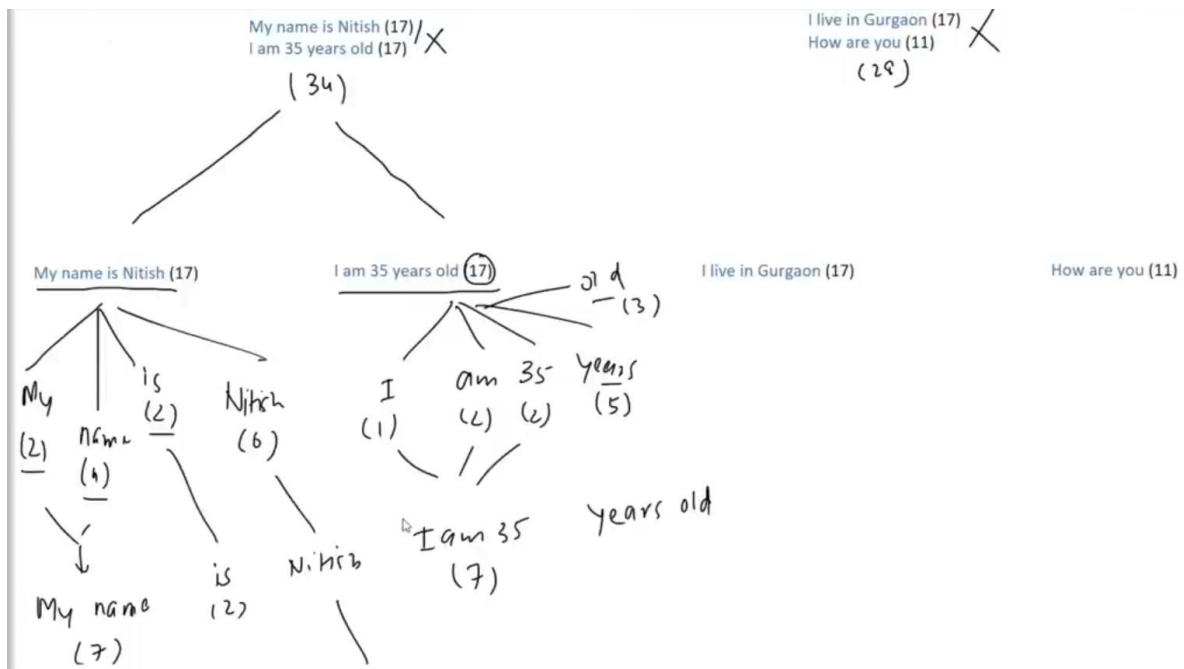
I live in Gurgaon (17) X
How are you (11) X
(28)

My name is Nitish (17)

I am 35 years old (17)

I live in Gurgaon (17)

How are you (11)



chunk size = 25 .

My name is Nitish (17)
I am 35 years old (17)

I live in Gurgaon (17)
How are you (11)

My name is Nitish (17)
I am 35 years old (17)

(34)

I live in Gurgaon (17)
How are you (11)

(28)

My name is Nitish (17)
I am 35 years old (17)

(34)

I live in Gurgaon (17)
How are you (11)

(28)

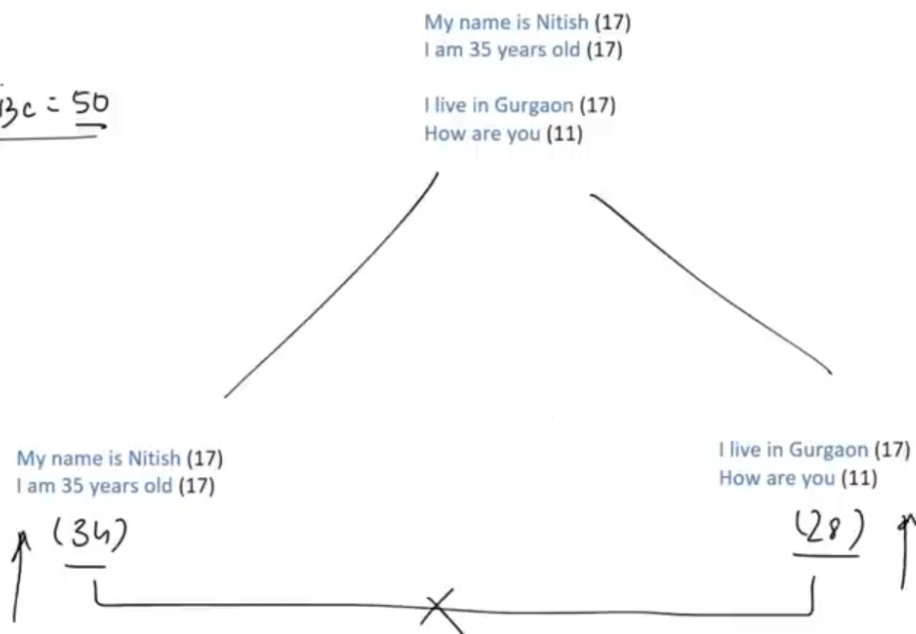
My name is Nitish (17) I am 35 years old (17)



I live in Gurgaon (17) How are you (11)



Chunk size = 50



3. Document-Structured Based

■ Project Name: Smart Student Tracker

A simple Python-based project to manage and track student data,

🔍 Features

- Add new students with relevant info
- View student details
- Check if a student is passing
- Easily extendable class-based design

✨ Tech Stack

- Python 3.10+
- No external dependencies

```
# First, try to split along Markdown headings (starting with level 2)
"\n#{1,6} ",
# Note the alternative syntax for headings (below) is not handled here
# Heading level 2
# -----
```

call →

↓

```
* class Student:
    def __init__(self, name, age, grade):
        self.name = name
        self.age = age
        self.grade = grade # Grade is a float (Like 8.5 or 9.2)

    def get_details(self):
        return f"Name: {self.name}, Age: {self.age}, Grade: {self.grade}"

    def is_passing(self):
        return self.grade >= 6.0

# Example usage
student1 = Student("Aarav", 20, 8.2)
print(student1.get_details())

if student1.is_passing():
    print("The student is passing.")
else:
    print("The student is not passing.")
```

e

```
# First, try to split along class definitions
"\nclass ",
"\ndef ",
"\n\ndef ",
```

```

# First, try to split along Markdown headings (starting with level 2)
"\n#{1,6} ",
# Note the alternative syntax for headings (below) is not handled here
# Heading level 2
# -----
# End of code block
"```\n",
# Horizontal lines
"\n\\*\\*\\*\\*+\\n",
"\n---+\\n",
"\n__+\\n",
# Note that this splitter doesn't handle horizontal lines defined
# by *three or more* of ***, ---, or __, but this is not handled
"\n\n",
"\n",
" ",
" ",

```

↘

```

# First, try to split along class definitions
"\nclass ",
"\ndef ",
"\ntdef ",
# Now split by the normal type of lines
"\n\n", —
"\n", —
" ", —
" ", —

```

4. Semantic Meaning Based

01 April 2025 18:11

Farmers were working hard in the fields, preparing the soil and planting seeds for the next season. The sun was bright, and the air smelled of earth and fresh grass. The Indian Premier League (IPL) is the biggest cricket league in the world. People all over the world watch the matches and cheer for their favourite teams.

Terrorism is a big danger to peace and safety. It causes harm to people and creates fear in cities and villages. When such attacks happen, they leave behind pain and sadness. To fight terrorism, we need strong laws, alert security forces, and support from people who care about peace and safety.

4. Semantic Meaning Based

01 April 2025 18:11

2 chunks

Farmers were working hard in the fields, preparing the soil and planting seeds for the next season. The sun was bright, and the air smelled of earth and fresh grass.
The Indian Premier League (IPL) is the biggest cricket league in the world. People all over the world watch the matches and cheer for their favourite teams.)

Terrorism is a big danger to peace and safety. It causes harm to people and creates fear in cities and villages. When such attacks happen, they leave behind pain and sadness. To fight terrorism, we need strong laws, alert security forces, and support from people who care about peace and safety.

