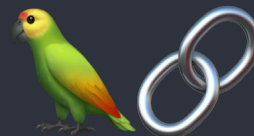


# Análise de textos com Langchain e Gemini

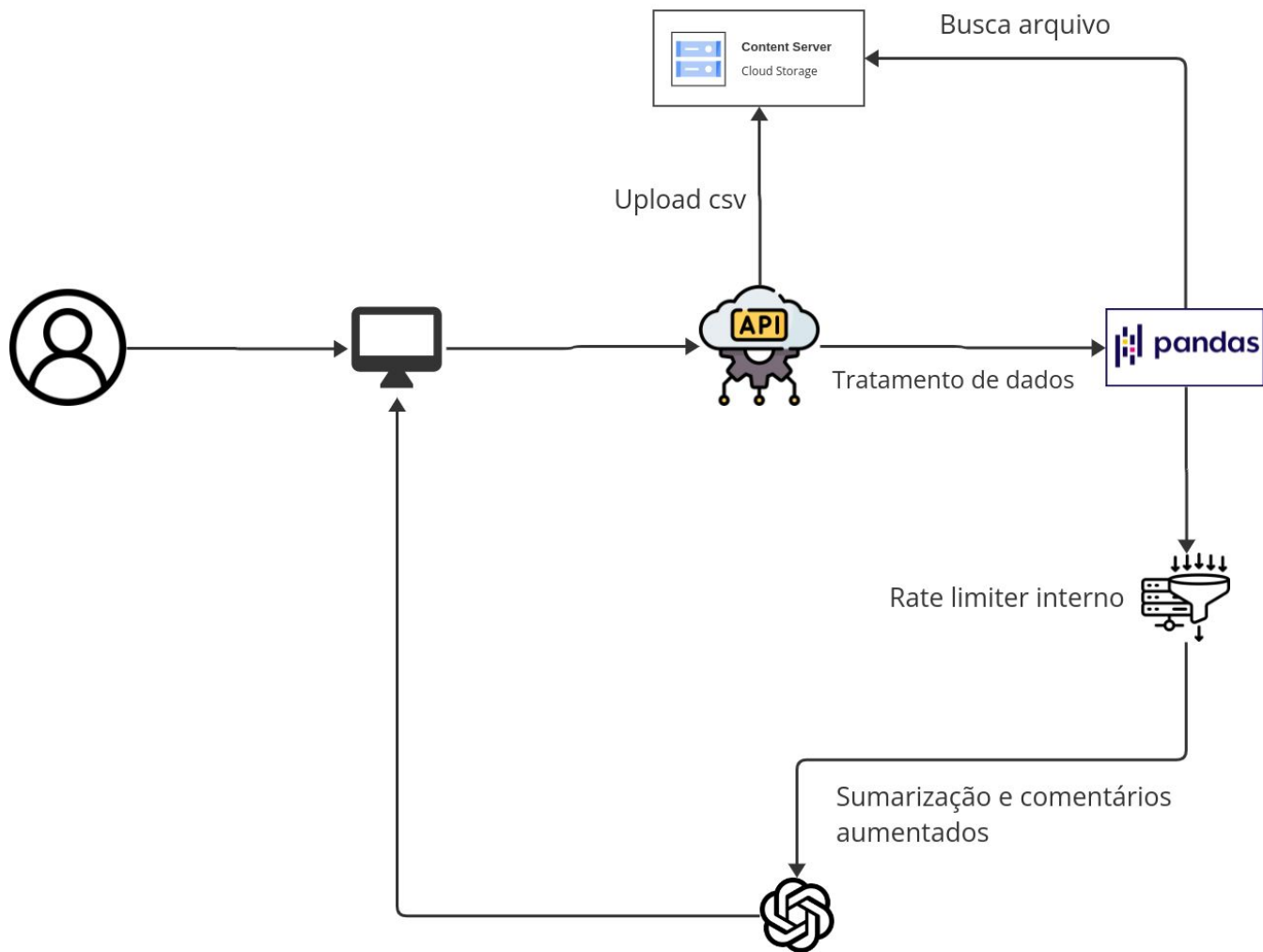


✦  
Gemini

1 0 1 1   0 1 1   0 1   1 0 1 1 0 0 1   1 0   1 1 0 1 1   0 1 1   0 1   1 1 0 1 1 0   1 1 0 1 1 1   1 1 0 1

# Lidando com LLMs no mundo real

1 0 1 1 0 1 1 0 1 1 0 1 1 0 0 1 1 0 1 1 0 1 1 0 1 1 0 1 1 0 1 1 0 1 1 0 1 1 1 1 0 1



# > Problemas enfrentados

- Tempo gasto no tratamento, extração dos dados, aplicação das regras e criação dos prompts;
- Grande números de requests;
- Limite de requests para a LLM;
- Alucinação;
- User Experience.

```

1 # ===== routes.py =====
2 @router.get("/{id}/analyze")
3 async def get_analysis(
4     id: str,
5     df: Annotated[pd.DataFrame, Depends(get_df_from_stored_file)],
6     storage: GCPStorage = Depends(get_file_storage),
7 ) -> DataResponse:
8     data = await process_file(df)
9     return DataResponse(**data)
10
11 # ===== processing.py =====
12 async def process_file(df: pd.DataFrame) -> dict:
13     ...
14     return await mount_data(df)
15
16
17 async def mount_data(df: pd.DataFrame) -> dict:
18     filtered_data = filter_data(df)
19     return await group_data_by_feature(filtered_data)
20
21
22 async def group_data_by_feature(items: list) -> dict:
23     ...
24     response = await asyncio.gather(*[mount_item(item) for item in items])
25     ...
26     return grouped_items
27
28
29 async def mount_item(item: pd.Series) -> Item:
30     ...
31     new_item = Item(...)
32     prompt = build_prompt(new_item)
33     response = await get_answer_from_llm(prompt)
34     ...
35     return new_item
36

```

```

1 # ===== llm.py =====
2 _SEMAPHORE = asyncio.Semaphore(10)
3
4 async def get_answer_from_llm(prompt: list) -> str:
5     chat = OpenAIChatRequest(model='gpt-4o', temperature=0.0, messages=prompt)
6     async with _SEMAPHORE:
7         async with httpx.AsyncClient(timeout=600) as client:
8             response = await client.post(
9                 urljoin(self.url, '/openai/v1/chat/completions'),
10                 headers=self.headers,
11                 json=chat.model_dump(),
12             )
13     ...
14     return response_data['choices'][0]['message']['content']

```

E as requisições em batch?

1 0 1 1   0 1 1   0 1   1 0 1 1 0 0 1   1 0   1 1 0 1 1   0 1 1   0 1   1 1 0 1 1 0   1 1 0 1 1 1   1 1 0 1



```
1 from openai import OpenAI
2
3 client = OpenAI()
4
5 batch_input_file = client.files.create(
6     file=open(FILE_NAME, "rb"),
7     purpose="batch"
8 )
9
10 batch_job = client.batches.create(
11     input_file_id=batch_input_file.id,
12     endpoint="/v1/chat/completions",
13     completion_window="24h"
14 )
```

## > Vantagens e desvantagens

- Diminui o número de requests para a LLM;
- Diminui os custos;
- Entretanto, na OpenAI, o processamento pode levar até 24h;
- Langchain ainda não suporta.



## > Quem sou eu?



YURI SALES

Desenvolvedor Back-End,  
pós-graduado em Engenharia de  
Software

linkedin: /yuriscosta

github: @yurisalesc