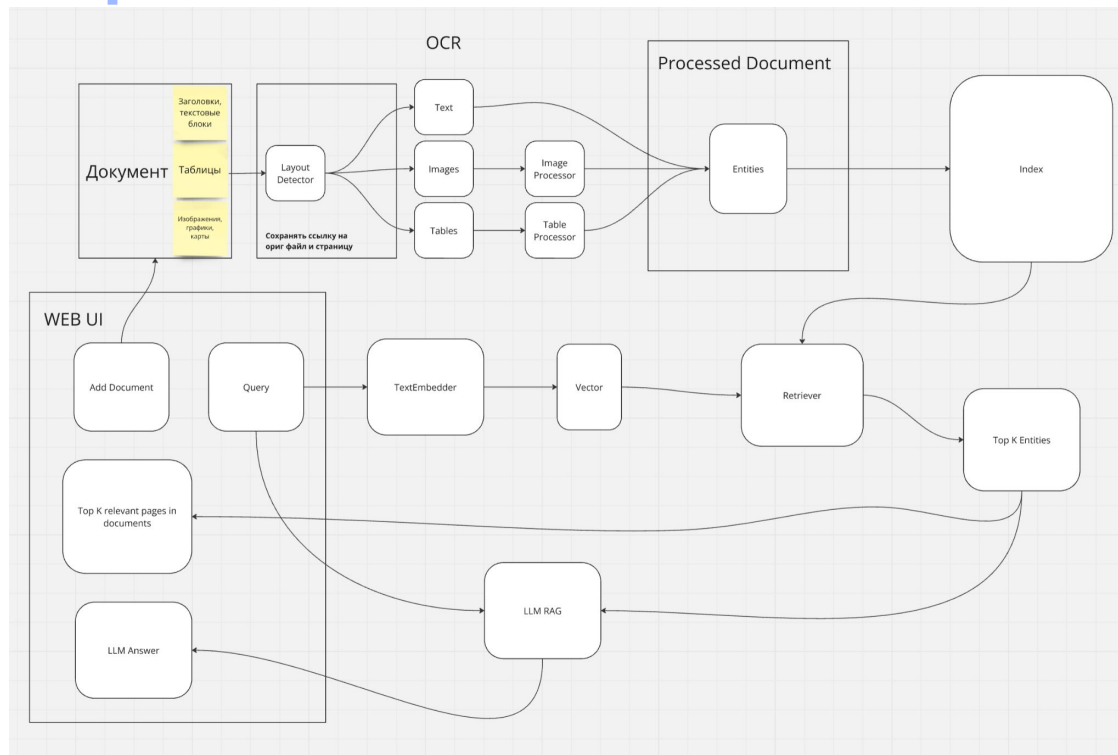


document-search

stay fine-tuned

Кейс №12. Автоматизация поиска и анализа нормативной документации по строительству объектов

4 сентября



OCR

Document reader

5 сентября

Document entities

Document entities + vector

Document entities + vector

Storage

Retriever

```

1 class DocumentReader:
2     ...
3
4     def process_pdf(self, document_path: str) -> ProcessedDocument:
5         ...
6
7     def process_docx(self, document_path: str) -> ProcessedDocument:
8         ...
9
10    ...
    
```

```

1 class ProcessedDocument:
2     name: str
3     id_: str
4     num_pages: int
5     original_format: Literal['pdf', 'docx']
6     entities: list[DocEntity]
7
8
9 class EntityPosition:
10    parent_document_id: str
11    page: int
12    position: float # подумать как лучше
13
14
15 class DocEntity(Protocol):
16    id_: int
17    position: EntityPosition
18
19 class TextDocEntity(DocEntity):
20    id_: int
21    position: EntityPosition
22    text: str
23
24 class TableDocEntity(DocEntity):
25    id_: int
26    position: EntityPosition
27    table: str
28
29 class ImageDocEntity(DocEntity):
30    id_: int
31    position: EntityPosition
32    image: ?
    
```

```

1 class TextEntityEmbedder(Protocol):
2
3     def vectorize(self, item: TextDocEntity) -> VectorizedDocEntity: ...
4
    
```

```

1 class VectorizedDocEntity(DocEntity):
2     vector: np.ndarray
3
4
5 class VectorizedDocument:
6     name: str
7     id_: str
8     entities: list[VectorizedDocEntity]
    
```

```

1 class DocumentStorage(Protocol):
2
3     def add_entity(self, entity: DocEntity) -> None: ...
4
5     def add_document(self, document: VectorizedDocument) -> None: ...
6
7     def get_relevant_indexes(self, document: VectorizedDocument, k: int) -> list[DocEntity]: ...
    
```

```

1 class Retriever(Protocol):
2
3     def __init__(self, storage: DocumentStorage): ...
4
5     def get_relevant_entities(self, query: str, k: int) -> list[DocEntity]: ...
6
    
```

6 сентября

Артём - Document entities (OCR), Web UI

Игорь - CI, readme и интерфейсы, Storage, Backend

Стас - Retriever, Embedder, RAG

Настя - презентация к чекпоинту, проработка идей и гипотез

Репозиторий:



Сложности

- OCR и извлечение сущностей
- Сложная структура документов
- Индексирование разных сущностей (таблицы, графики)

Планы

- Стабильный baseline
- Проработка разных вариантов работы с таблицами и изображениями
- Добавление RAG

Команда



Игорь Павлов
Капитан команды
@bom_bo0m



Артем Иванов
Разработчик
@incncvble



Анастасия Остапчук
ПМ
@anastasiia_ost_v



Станислав Стафиевский
Разработчик
@stafstas