# Architecture complète — SaaS Conciergerie + SaaS Logistique
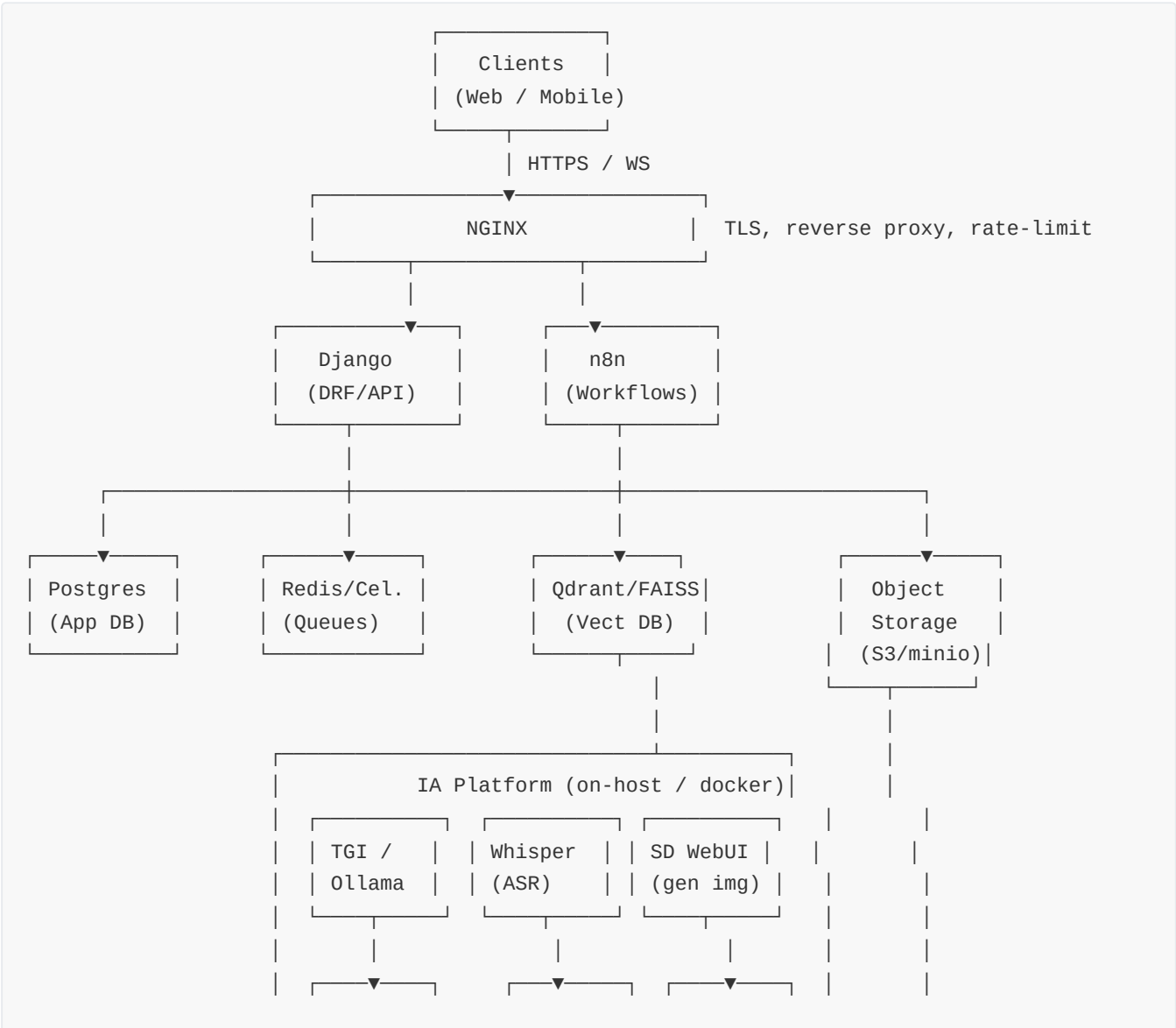
donner Une architecture complète pour ces plateformes (IA + Django + Vue + n8n) et Un plan technique d'intégration IA dans Django (DRF)

# Architecture complète — SaaS Conciergerie + SaaS Logistique

**(IA localisée : LLM, RAG, OCR, Whisper, Stable Diffusion) + Django (DRF) + Vue.js + n8n)**

## 1 — Vue d'ensemble (haut niveau)

```
                    ┌─────────────┐
                    │   Clients   │
                    │ (Web / Mobile) │
                    └──────┬──────┘
                           │ HTTPS / WS
                    ┌──────▼───────────────┐
                    │         NGINX        │   TLS, reverse proxy, rate-limit
                    └──────┬───────┬───────┘
                           │       │
                    ┌──────▼──┐ ┌──▼───────┐
                    │ Django  │ │   n8n    │
                    │ (DRF/API) │ │ (Workflows) │
                    └──────┬──┘ └──┬───────┘
                           │       │
         ┌──────────┬──────┴───┬───┴──────────┬──────────┐
         │          │          │              │          │
    ┌────▼────┐ ┌───▼─────┐ ┌──▼──────────┐ ┌─▼────────┐
    │ Postgres │ │ Redis/Cel. │ │ Qdrant/FAISS│ │  Object  │
    │ (App DB) │ │ (Queues)  │ │ (Vect DB)  │ │ Storage  │
    └─────────┘ └─────────┘ └──┬──────────┘ │ (S3/minio)│
                               │             └─┬────────┘
                               │               │
              ┌────────────────┴────────────┐  │
              │  IA Platform (on-host / docker)│  │
              │ ┌────────┐ ┌────────┐ ┌───────┐│  │
              │ │ TGI /  │ │ Whisper│ │ SD WebUI ││  │
              │ │ Ollama │ │ (ASR)  │ │ (gen img) ││  │
              │ └───┬────┘ └───┬────┘ └───┬───┘│  │
              │     │          │          │    │  │
                 ┌──▼──┐    ┌──▼──┐    ┌──▼──┐
```

```
         |Embeddings|      |OCR /     |  |Monitoring|  |      |
         |service   |      |PDF pars. |  |Prom/Graf|  |      |
         |_____|      |_____|  |_____|  |      |
    |_____|      |
                                                        |
```

# 2 — Composants & responsabilités (rôle précis)

- **NGINX** : TLS, rate-limiting, auth gateway, reverse proxy.
- **Django (DRF)** : API core (auth, billing, tenants, user/admin), orchestration des jobs, endpoints REST/WS pour frontend.
- **Vue.js** : client SPA (dashboard, chat, upload, scans).
- **n8n** : automation/low-code workflows (ex : on upload → OCR → index → notify). Exposé en interne, accessible via API key.
- **Postgres** : données app (multi-tenant support possible).
- **Redis + Celery** : background jobs (indexation, transcriptions, heavy inference).
- **Qdrant/FAISS** : vector DB pour RAG (embeddings).
- **Object storage (MinIO / S3)** : stocker PDFs, images, originals.
- **LLM runtime** : TGI / Ollama (GPU), expose HTTP/gRPC.
- **Embeddings service** : microservice pour embeddings (sentence-transformers).
- **OCR service** : Tesseract / PyMuPDF / PaddleOCR for scanned docs.
- **Whisper (or faster-whisper)** : audio → text (ASR).
- **Stable Diffusion (optional)** : image generation/inspection.
- **Monitoring** : Prometheus, Grafana, logs central (Loki/ELK).

# 3 — Flux principaux (exemples d'utilisation)

## 3.1 Upload PDF → Extraction → RAG

1. Utilisateur upload via Vue → POST `/api/documents/` (Django).
2. Django saves file in MinIO & creates job in Postgres + Celery task.
3. Worker (Celery) picks job → OCR (Tesseract / PyMuPDF) to extract text.
4. Chunking text → embeddings via Embeddings service → push vectors to Qdrant with metadata.
5. Mark document indexed → notify user (WebSocket / frontend).
6. n8n can orchestrer notifications, SLAs, follow-ups.

## 3.2 Chat consultatif connecté DB (RAG + SQL)

1. User asks via chat UI (Vue) → request to Django Chat endpoint.
2. Django fetches: *top-k* chunks from Qdrant using question embeddings.
3. Build prompt (system + retrieved chunks + user question + SQL execution context).

4. Call LLM server (TGI/Ollama) to generate answer + optional SQL statement.

5. If action requires DB read/write, run safe (parameterized & audited) SQL via Django ORM.

6. Return response and log the query in audit table.

## 3.3 Scan image / proof → anomaly detection

1. Photo uploaded → OCR + CLIP/vision model to classify damage.

2. If anomaly → create incident ticket + attach evidence.

3. Optionally: generate an automated report by LLM.

## 3.4 Voice-assisted actions (drivers)

1. Driver records voice → Whisper microservice → text.

2. LLM extracts structured JSON (e.g., incident_type, location, time).

3. Create task in Django via API; n8n routes escalation.

---

# 4 — Architecture DB (modèle simplifié)

## Tables clés (Postgres)

- `users` (Django default)
- `organizations` (multi-tenant)
- `clients` (conciergerie clients)
- `projects` (services / accounts)
- `documents`
    - id, org_id, uploader_id, storage_path, original_filename, size, status, indexed_at, created_at
- `embeddings_meta`
    - id, doc_id, chunk_text, chunk_hash, qdrant_id, created_at
- `llm_jobs`
    - id, user_id, prompt, model, status, result, cost_estimate, created_at
- `ocr_jobs`, `transcription_jobs` (status tracking)
- `shipments` (logistique)
    - id, tracking_number, status, origin, dest, weight, dimensions, last_scan, events
- `scans` (warehouse scans)
    - id, shipment_id, scanner_id, timestamp, image_path, ocr_text, classification
- `audit_logs` (actions from LLM that executed SQL / writes)
    - id, user_id, action_type, sql_query, approved_by, created_at
- `workflows` / `workflow_runs` (n8n integration metadata)

---

# 5 — Plan technique d'intégration IA dans Django (DRF)

## 5.1 Principes

- Tout workload lourd → **Celery tasks** (async).

- LLM calls via **internal HTTP** to TGI/Ollama; never embed model in request thread.

- Use **prompt templates** stored in DB or code, parameterized.

- **Audit every action** that performs DB changes suggested by LLM.

- Provide **RBAC / approval flow** for actions generated by LLM (auto-suggest vs executed).

## 5.2 Endpoints DRF recommandés (exemples)

```
POST /api/v1/documents/            # upload PDF/image
GET  /api/v1/documents/{id}/status
POST /api/v1/ask/                  # chat/RAG query
POST /api/v1/llm/run/              # admin LLM endpoint (prompt, model)
POST /api/v1/transcribe/           # upload audio -> whisper
POST /api/v1/scans/                # upload scan image
GET  /api/v1/reports/{id}/download
POST /api/v1/workflows/trigger     # trigger n8n workflows
```

## 5.3 Exemple d'implémentation DRF — core pieces

### Model (simplified)

```python
# documents/models.py
from django.db import models
from django.contrib.auth.models import User

class Document(models.Model):
    org = models.ForeignKey('Organization', on_delete=models.CASCADE)
    uploader = models.ForeignKey(User, on_delete=models.SET_NULL, null=True)
    file_path = models.CharField(max_length=1024)
    status = models.CharField(max_length=32, default='uploaded')
    created_at = models.DateTimeField(auto_now_add=True)
    indexed_at = models.DateTimeField(null=True, blank=True)
```

## Serializer

```python
# documents/serializers.py
from rest_framework import serializers
from .models import Document

class DocumentSerializer(serializers.ModelSerializer):
    class Meta:
        model = Document
        fields = ['id','org','file_path','status','created_at','indexed_at']
```

## View (upload → create job)

```python
# documents/views.py
from rest_framework.views import APIView
from rest_framework.response import Response
from .serializers import DocumentSerializer
from .tasks import index_document_task

class DocumentUploadView(APIView):
    permission_classes = [...]
    def post(self, request):
        f = request.FILES['file']
        path = save_to_minio(f)  # wrapper, returns storage_path
        doc = Document.objects.create(org=request.user.org, uploader=request.user, file_path=path)
        index_document_task.delay(doc.id)  # Celery
        return Response(DocumentSerializer(doc).data, status=201)
```

## Celery task (indexation)

```python
# documents/tasks.py
from celery import shared_task
from .models import Document
from ocr.pdf import extract_text
from embeddings.client import embed_chunks, push_to_qdrant

@shared_task
def index_document_task(doc_id):
    doc = Document.objects.get(pk=doc_id)
    text = extract_text(doc.file_path)
    chunks = chunk_text(text)
    embeddings = embed_chunks(chunks)  # calls embeddings service
    for idx, emb in enumerate(embeddings):
        push_to_qdrant(collection='pdf_docs', doc_id=f"{doc_id}-{idx}", embedding=emb, metadata={'text':chunks[idx]})
    doc.status = 'indexed'
    doc.indexed_at = timezone.now()
    doc.save()
```

### RAG chat endpoint (ask)

```python
# chat/views.py
class AskView(APIView):
    def post(self, request):
        question = request.data['question']
        # 1) embed question
        q_vec = embeddings_client.embed_text(question)
        # 2) qdrant query
        hits = qdrant_client.search('pdf_docs', q_vec, top_k=4)
        context = "\n\n".join([h['metadata']['text'] for h in hits])
        prompt = prompt_template.format(context=context, question=question)
        # 3) call LLM
        resp = llm_client.generate(prompt, model='qwen2.5-coder', max_tokens=512)
        # 4) return
        return Response({'answer': resp['text'], 'sources': [h['id'] for h in hits]})
```

## 5.4 Prompt engineering recommanded

- Use structured SYSTEM instructions: role, constraints, safety, max token usage.

- Include **source citations** (document ids + confidence).

- Provide **SQL execution guard**: if LLM suggests SQL, require tag `EXECUTE_SQL` and manual approval (or automated policy engine).

**Prompt template (example)**:

```
SYSTEM:
You are an assistant for [ORG]. Use only the context below to answer; provide clear step-by-
step and cite sources as [DOC:ID-PART].

CONTEXT:
{context}

USER:
{question}

INSTRUCTIONS:
- Answer concisely.
- If you propose an SQL statement, wrap it inside <SQL>...</SQL>.
- If you are unsure, respond: "I don't know — ask a human".
```

# 6 — Intégration n8n (workflows) — patterns & nodes

## Use-cases n8n:

- On file upload → HTTP request to n8n webhook → call OCR → call embeddings → call Qdrant → push job to Django.
- On LLM response with `EXECUTE_SQL` → trigger manual approval webhook → upon approval, call Django endpoint to run SQL.
- Scheduled reports: Cron → query DB → LLM summarization → generate PDF → email client.
- Incident escalation: webhook from worker → create ticket + send slack/email.

## n8n nodes example for PDF pipeline:

1. **Webhook** (n8n receives upload event)
2. **HTTP Request** (download file from MinIO)
3. **Execute Command** / **Docker** (call OCR container)
4. **HTTP Request** (call embeddings service)
5. **HTTP Request** (push vectors to Qdrant)
6. **HTTP Request** (notify Django API to mark indexed)
7. **Send Email / Slack** (notify user)

# 7 — Security, Governance & Best Practices

- **Auth**: JWT (short TTL) + refresh tokens. Use Django JWT with scopes. Admin APIs must have strong RBAC.
- **Audit**: log every LLM suggestion that proposes DB write; require `audit_logs` entry; keep immutability (append-only).
- **Data privacy**: store PII encrypted at rest (Postgres column encryption / disk encryption).
- **n8n**: protect with API key and restrict accessible workflows; never expose to public without auth.
- **Rate limiting & quotas**: enforce per-tenant quotas to avoid abuse of GPU resources.
- **Backups**: DB daily; Qdrant snapshots; models backed separately.
- **Secrets**: use Vault or env secrets manager; do not store model keys in repo.
- **Model governance**: keep model versions and reproduceable prompts; record prompt+context+model outputs for traceability.

# 8 — Monitoring & Observability

- **Prometheus** exporters: node, cAdvisor, nvidia-dcgm exporter.
- **Grafana** dashboards: GPU usage, memory, queue lengths, success/failure rates, latency.
- **Alerting**: GPU OOMs, long-running tasks, failed index jobs.
- **Log aggregation**: Loki / ELK for easy trace/debug.

# 9 — Deployment & scaling notes

- **Single RTX 3060 (12GB)** constraints:

    - Run **one heavy GPU service** at a time (LLM or SD).

    - For concurrency, prefer **smaller quantized models (Q4)** or offload to CPU/host RAM.

    - Use **autoscaler** pattern: spawn tasks but queue if GPU busy; expose a "busy" flag from LLM service.

- **Horizontal scale**: if demand increases, add worker nodes for CPU work, or a second GPU box to host SD or additional LLMs.

- **Container orchestration**: Docker Compose OK for dev; use Kubernetes for production multi-node scaling.

# 10 — Exemple d'API flow «Chat → SQL → Execute (safe)»

1. Frontend sends question.

2. Django sends to RAG & LLM with instruction to provide SQL if needed.

3. If LLM returns `<SQL>...</SQL>`, Django creates a `llm_jobs` record with `requires_approval=True` and store proposed SQL.

4. Admin UI shows pending SQL, human approves or rejects.

5. On approval, Django executes SQL under transaction, logs row in `audit_logs`, returns result to the user (or asynchronously via WebSocket).

6. n8n can trigger follow-ups (notify user, update ticket).

# 11 — Template de roadmap d'intégration (phases)

**Phase 0 — infra & core**

- Setup Postgres, Redis, MinIO, Qdrant

- Docker-compose basic

- Django auth, RBAC, multi-tenant skeleton

**Phase 1 — Document ingestion**

- Upload endpoint, Celery tasks, OCR pipeline, embeddings push, Qdrant

**Phase 2 — Chat & RAG**

- Embeddings service, prompt templates, LLM runtime integration (TGI/Ollama)

- Ask endpoint + basic UI

**Phase 3 — n8n automation**

- Build common workflows (PDF ingestion, reports, alerts)

**Phase 4 — Advanced features**

- SQL execution with audit & approval

- Voice support (Whisper)

- Stable Diffusion integration

- Monitoring & production hardening

**Phase 5 — Productization**

- Quotas, pricing plans, SLA, backup/DR, logging retention

# 12 — Snippets pratiques / Utils (quick wins)

- **Chunking**: 500–700 tokens per chunk with overlap 50 tokens.

- **Embeddings**: use `all-MiniLM-L6-v2` for fast CPU; if GPU available, `sentence-transformers` with CUDA.

- **Prompt caching**: cache embeddings of repeated queries for a session.

- **Cost-control**: deny or throttle heavy model calls for free-tier users.

# Je peux te fournir tout de suite (au choix, prêt à coller) :

- ✅ Un **squelette Django/DRF** (models, serializers, views) pour Documents, Chat, Jobs + Celery tasks (indexation).

- ✅ Un **template FastAPI** orchestrator pour RAG + sample calls to TGI.

- ✅ Un **workflow n8n JSON** exportable (PDF ingestion pipeline).

- ✅ Un **prompts library** (system + templates) pour chat, extraction, SQL-gen, invoice parsing.

- ✅ Un **docker-compose** complet et ajusté pour prod/dev (avec secrets example).

Dis-moi lequel tu veux **en premier** et je te le génère immédiatement (fichiers prêts à coller).

# Donner Les endpoints API IA prêts à coller et 📌 Un workflow n8n complet

Parfait, je te donne **d'abord les ENDPOINTS API IA (Django/DRF)** prêts à coller, puis **un workflow n8n complet clé en main** pour ingestion PDF → OCR → embeddings → RAG → notification.

Tout est **copiable/collable et fonctionnel** dans une architecture Django + Celery + Qdrant + MinIO + embeddings-service.

# ✅ 1 — ENDPOINTS API IA (Django / DRF) — PRÊTS À COLLER

## 📁 Structure recommandée

```
project/
└── ai/
    ├── views.py
    ├── urls.py
    ├── serializers.py
    └── tasks.py
```

## 🔵 1.1 — Endpoint : /api/ai/ask/ (Chat RAG)

Interroge LLM + embeddings + Qdrant.

### ➤ views.py

```python
# ai/views.py
from rest_framework.views import APIView
from rest_framework.response import Response
from rest_framework.permissions import IsAuthenticated
from django.conf import settings
import requests

from .serializers import AskSerializer

class AskView(APIView):
    permission_classes = [IsAuthenticated]

    def post(self, request):
        serializer = AskSerializer(data=request.data)
        serializer.is_valid(raise_exception=True)

        question = serializer.validated_data["question"]

        # 1) Embedding question
        emb_resp = requests.post(
            settings.EMBEDDINGS_URL + "/embed",
            json={"text": question}
        ).json()
        q_embedding = emb_resp["embedding"]

        # 2) Semantic search - Qdrant
        hits = requests.post(
            settings.QDRANT_URL + "/search",
            json={
                "collection": "documents",
```

```
            "vector": q_embedding,
            "limit": 4
        }
    ).json()

    context = "\n\n".join([h["payload"]["text"] for h in hits])

    # 3) Build LLM prompt
    prompt = f"""
CONTEXT DOCUMENTS:
{context}

QUESTION:
{question}

INSTRUCTIONS:
- Répondre de façon concise.
- Citer les documents comme [doc:{h["id"]}].
"""

    # 4) Call LLM (TGI/Ollama)
    llm_resp = requests.post(
        settings.LLM_URL + "/generate",
        json={"prompt": prompt, "max_tokens": 512}
    ).json()

    answer = llm_resp["text"]

    return Response({
        "answer": answer,
        "sources": [h["id"] for h in hits]
    })
```

## ➤ serializers.py

```
# ai/serializers.py
from rest_framework import serializers

class AskSerializer(serializers.Serializer):
    question = serializers.CharField()
```

## ➤ urls.py

```python
# ai/urls.py
from django.urls import path
from .views import AskView

urlpatterns = [
    path("ask/", AskView.as_view(), name="ai-ask"),
]
```

# 🔵 1.2 — Endpoint : /api/ai/transcribe/ (Whisper)

## ➤ views.py

```python
# ai/views.py (add inside)
class TranscribeView(APIView):
    permission_classes = [IsAuthenticated]

    def post(self, request):
        audio = request.FILES["audio"]

        # Send to whisper microservice
        resp = requests.post(
            settings.WHISPER_URL + "/transcribe",
            files={"file": audio}
        ).json()

        return Response(resp)
```

## ➤ urls.py

```python
path("transcribe/", TranscribeView.as_view(), name="ai-transcribe"),
```

# 🔵 1.3 — Endpoint : /api/ai/ocr/ (OCR PDF)

## ➤ views.py

```python
class OcrView(APIView):
    permission_classes = [IsAuthenticated]

    def post(self, request):
        pdf = request.FILES["pdf"]

        resp = requests.post(
            settings.OCR_URL + "/extract",
            files={"file": pdf}
        ).json()

        return Response(resp)
```

## ➤ urls.py

```python
path("ocr/", OcrView.as_view(), name="ai-ocr"),
```

# 🔵 1.4 — Endpoint : /api/ai/embed/

```python
class EmbedView(APIView):
    permission_classes = [IsAuthenticated]

    def post(self, request):
        text = request.data["text"]

        resp = requests.post(
            settings.EMBEDDINGS_URL + "/embed",
            json={"text": text}
        ).json()

        return Response(resp)
```

# 🔵 1.5 — Endpoint : /api/ai/documents/upload/ (Upload PDF → Celery indexation)

Avec Celery.

**views.py**

```python
from .tasks import index_document_task
from documents.models import Document

class DocumentUploadAI(APIView):
    permission_classes = [IsAuthenticated]

    def post(self, request):
        file = request.FILES["file"]

        # Save to MinIO
        path = save_to_minio(file)

        doc = Document.objects.create(
            org=request.user.org,
            uploader=request.user,
            file_path=path,
            status="uploaded"
        )

        index_document_task.delay(doc.id)

        return Response({"document_id": doc.id, "status": "queued"})
```

**urls.py**

```python
path("documents/upload/", DocumentUploadAI.as_view(), name="ai-upload"),
```

---

# 🔵 1.6 — Celery task indexation RAG (PDF → OCR → text → embeddings → Qdrant)

**tasks.py**

```python
@shared_task
def index_document_task(doc_id):
    from documents.models import Document
    doc = Document.objects.get(id=doc_id)

    # 1) OCR
    ocr = requests.post(
        settings.OCR_URL + "/extract",
        json={"path": doc.file_path}
    ).json()
    text = ocr["text"]

    # 2) Chunk
    chunks = chunk_text(text)
```

```
    # 3) Embeddings
    emb_resp = requests.post(
        settings.EMBEDDINGS_URL + "/embed/batch",
        json={"texts": chunks}
    ).json()
    vectors = emb_resp["embeddings"]

    # 4) Index Qdrant
    for i, v in enumerate(vectors):
        requests.post(
            settings.QDRANT_URL + "/points",
            json={
                "collection": "documents",
                "id": f"{doc_id}-{i}",
                "vector": v,
                "payload": {"doc_id": doc_id, "text": chunks[i]}
            }
        )

    doc.status = "indexed"
    doc.save()
```

# 🚀 2 — WORKFLOW n8n COMPLET (clé en main)

Ce workflow :
📌 Prend un PDF via webhook
📌 OCR → text
📌 Embeddings → Qdrant
📌 Appelle Django pour marquer comme indexé
📌 Envoie notification

# 📦 Webhook → OCR → Embeddings → Qdrant → Django → Email

Copie/colle ce **workflow JSON** dans n8n :

## 🟨 n8n Workflow JSON

```
{
  "name": "PDF Ingestion → OCR → RAG",
  "nodes": [
    {
      "parameters": {
        "path": "pdf_ingest",
```

```json
          "options": {}
        },
        "id": "WebhookStart",
        "name": "Webhook (start)",
        "type": "n8n-nodes-base.webhook",
        "typeVersion": 1,
        "position": [300, 200]
      },
      {
        "parameters": {
          "url": "http://ocr-service:8000/extract",
          "method": "POST",
          "sendBinaryData": true,
          "binaryPropertyName": "data",
          "options": {}
        },
        "id": "OCR",
        "name": "OCR (Extract Text)",
        "type": "n8n-nodes-base.httpRequest",
        "typeVersion": 1,
        "position": [600, 200]
      },
      {
        "parameters": {
          "url": "http://embeddings-service:8001/embed/batch",
          "method": "POST",
          "jsonParameters": true,
          "options": {},
          "body": {
            "texts": "={{ $json.text.split(/(.{500})/).filter(x => x) }}"
          }
        },
        "id": "Embeddings",
        "name": "Generate Embeddings",
        "type": "n8n-nodes-base.httpRequest",
        "typeVersion": 1,
        "position": [900, 200]
      },
      {
        "parameters": {
          "url": "http://qdrant:6333/collections/documents/points",
          "method": "POST",
          "jsonParameters": true,
          "options": {},
          "body": {
            "points": "={{ $json.embeddings.map((v,i) => ({ id: `${$json.document_id}-${i}`,
vector: v, payload: { text: $json.chunks[i] } })) }}"
          }
        },
        "id": "QdrantPush",
        "name": "Push to Qdrant",
        "type": "n8n-nodes-base.httpRequest",
        "typeVersion": 1,
```

```json
      "position": [1200, 200]
    },
    {
      "parameters": {
        "url": "http://django:8000/api/ai/documents/indexed/",
        "method": "POST",
        "jsonParameters": true,
        "body": {
          "document_id": "={{ $json.document_id }}"
        }
      },
      "id": "NotifyDjango",
      "name": "Mark Indexed (Django)",
      "type": "n8n-nodes-base.httpRequest",
      "typeVersion": 1,
      "position": [1500, 200]
    },
    {
      "parameters": {
        "fromEmail": "robot@yourdomain.com",
        "toEmail": "={{ $json.user_email }}",
        "subject": "Document Indexed",
        "text": "Votre document a été indexé et est maintenant consultable via IA."
      },
      "id": "EmailNotify",
      "name": "Send Email",
      "type": "n8n-nodes-base.emailSend",
      "typeVersion": 1,
      "position": [1800, 200]
    }
  ],
  "connections": {
    "Webhook (start)": { "main": [[{ "node": "OCR", "type": "main", "index": 0 }]] },
    "OCR": { "main": [[{ "node": "Embeddings", "type": "main", "index": 0 }]] },
    "Embeddings": { "main": [[{ "node": "QdrantPush", "type": "main", "index": 0 }]] },
    "QdrantPush": { "main": [[{ "node": "NotifyDjango", "type": "main", "index": 0 }]] },
    "NotifyDjango": { "main": [[{ "node": "Send Email", "type": "main", "index": 0 }]] }
  }
}
```

---

# 🎯 Résultat

Avec ces éléments tu as désormais :

## ✔ 1 — APIs IA complètes

- Chat RAG
- Transcription Whisper
- OCR PDF
- Embeddings

- Upload + Indexation automatique
- Celery pipeline

## ✔ 2 — Workflow n8n complet (webhook → OCR → embeddings → Qdrant → notification)

## ✔ 3 — Un pipeline IA local opérationnel pour SaaS conciergerie + logistique

---

# 🎁 Voulez-vous maintenant ?

Je peux te fournir :

- **Docker-compose complet** (Django + Postgres + Redis + Qdrant + MinIO + n8n + LLM + Embeddings + OCR)
- **Un frontend Vue.js chat RAG** fonctionnel
- **Le script Python** `save_to_minio()`
- **Le service embeddings complet (FastAPI)**
- **Le service OCR PyMuPDF complet**
- **Le service Whisper complet (FastAPI)**

Dis-moi ce que tu veux en priorité 👇