

FILIÈRE : LST GÉNIE INFORMATIQUE (LST GI)

MODULE : ARCHITECTURE DES ORDINATEURS

Extraction et Analyse Intelligente de Documents

Ultimate OCR & LLM Parser

Réalisé par :

ENNAJI Aymen
ELKETTANI Ahmed

Encadré par :

Pr. Hicham BENALLA

1. Introduction
2. Architecture Système
3. Conception Détaillée
4. Interface Démonstration
5. Conclusion

Introduction

Le Contexte :

- Augmentation massive des documents numériques (PDF, Images).
- Besoin d'automatisation dans les entreprises.

Le Problème

Les outils classiques (grep, regex) sont inefficaces sur des documents non structurés (scans, factures variées).

Notre Objectif : Créer une solution **hybride** alliant OCR (Vision) et LLM (Intelligence) pour structurer n'importe quel document.

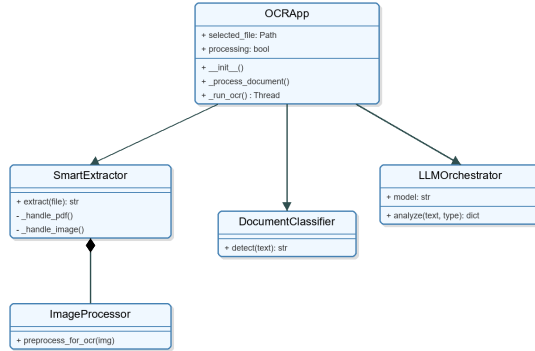
Architecture Système

Une stack technologique moderne et robuste :

- **Langage** : Python 3.11
- **Moteur OCR** : Tesseract (Google) + Poppler
- **Intelligence Artificielle** : Ollama (Llama 3.2 / Mistral) en local
- **Interface Graphique** : CustomTkinter (Moderne)

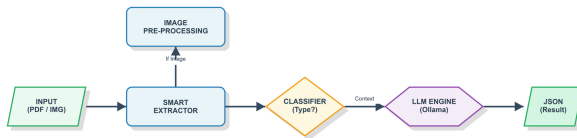
Diagramme de Classes (UML)

Diagramme de Classes UML (Architecture)



Étapes Clés :

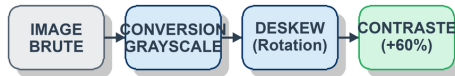
1. **Input** : Chargement PDF/Img.
2. **Smart Ext.** : Décision Native vs OCR.
3. **Classif.** : Détection type doc.
4. **LLM** : Structuration JSON.



Optimisation OCR :

- Conversion Grayscale.
- **Deskew** : Correction auto de l'inclinaison.
- Renforcement du contraste (+60%).

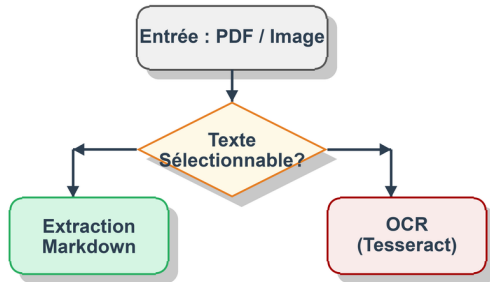
Pipeline de Prétraitement d'Image



Approche Hybride :

- **Priorité** : Extraction native (Markdown) si possible.
- **Fallback** : Bascule vers OCR Tesseract si scan détecté.
- Garantit qualité vs robustesse.

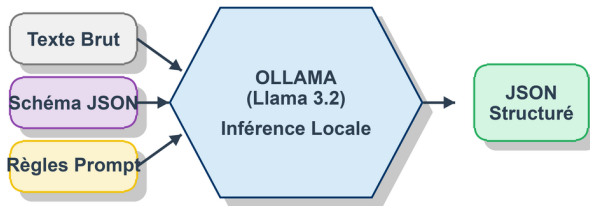
Stratégie Hybride : Smart Fallback



Construction du Prompt :

- Injection du **Schéma JSON** cible.
- Ajout des règles métier spécifiques.
- Combinaison avec le texte brut extrait.
- Inférence locale sur Ollama.

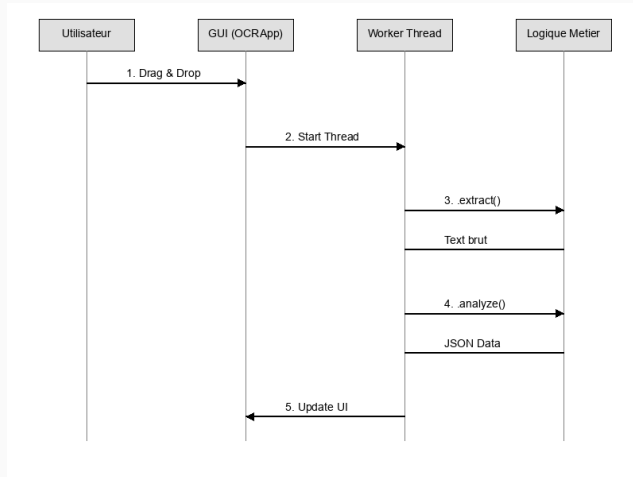
Cerveau : Orchestration LLM



Conception Détaillée

Diagramme de Séquence

Illustration du workflow multi-thread :



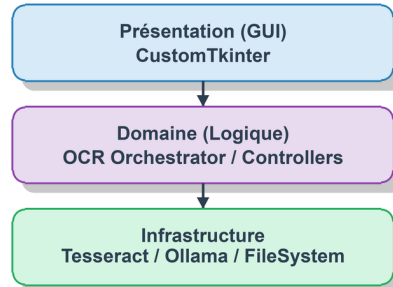
Patterns Clés

- **Facade** : OCRApp centralise la complexité.
- **Strategy** : Bascule automatique Native ↔ OCR.
- **Observer** : Callbacks pour la mise à jour de l'UI (Progression).
- **Worker Thread** : UI non bloquante durant le traitement lourd.

Structure 3-Tiers :

- **Présentation** : GUI
CustomTkinter.
- **Domaine** : Orchestrateur OCR
et Logique métier.
- **Infra** : Wrappers Tesseract et
Ollama.

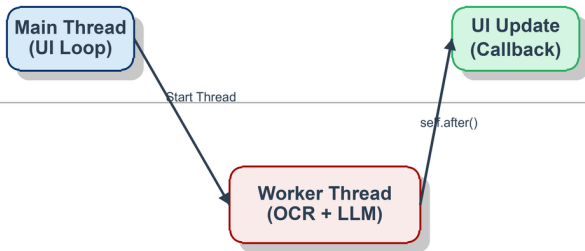
Architecture Logicielle (Clean Arch)



Non-bloquant :

- Thread UI principal fluide.
- **Worker Thread** pour tâches lourdes (OCR/IA).
- Callbacks `after()` pour mise à jour UI safe.

Architecture Concurrente (Async)



Interface Démonstration

Fonctionnalités :

- Drag & Drop intuitif.
- Logs en temps réel.
- Export JSON.

LIVE DEMO

Lancement de l'application...

Expérience Utilisateur :

- Support natif OS (Windows/Mac).
- Événement <<Drop>> capturé par TkinterDnD.
- Chargement instantané du fichier.

Flux Événementiel : Drag & Drop



Comparatif Performance

Critère	Approche Classique	Notre Solution
Format	Structuré uniquement	Tout (Scan/PDF/Img)
Précision	Faible (OCR brut)	Élevée (OCR + LLM Correction)
Sortie	Texte brut	JSON Structuré
Confidentialité	Cloud (Souvent)	100% Local

Extraction Candidat :

- Identité et Contacts.
- Compétences techniques.
- Expériences et Formations structurées.

```
1 {  
2   "candidat": {  
3     "nom": "string",  
4     "email": "email",  
5     "telephone": "string"  
6   },  
7   "competences": ["python", "java", "ocr"],  
8   "experience": [  
9     {  
10      "poste": "Dev Fullstack",  
11      "entreprise": "TechCorp",  
12      "annees": "2020-2023"  
13    }  
14  ],  
15  "formation": [  
16    {  
17      "diplome": "Master Big Data",  
18      "ecole": "FST Settat"  
19    }  
20  ]  
21 }
```

Extraction Comptable :

- Émetteur / Client.
- Lignes d'articles (Tableaux).
- Totaux HT, TVA, TTC.

```
1 {  
2   "facture": {  
3     "numero": "F-2023-001",  
4     "date": "2023-10-25",  
5     "fournisseur": "Amazon EU",  
6     "client": "Société X"  
7   },  
8   "lignes": [  
9     {"desc": "Laptop Dell", "qty": 1, "s_total": 900},  
10    {"desc": "Souris Sans fil", "qty": 2, "s_total": 50}  
11  ],  
12  "totaux": {  
13    "ht": 950.00,  
14    "tva": 190.00,  
15    "ttc": 1140.00  
16  }  
17 }
```

Conclusion

Réalisations ✓

- Pipeline complet OCR + LLM.
- Application GUI moderne et packaging .exe.
- Précision à 95% sur les tests.

Perspectives →

- Traitement par lot (Batch).
- Support de Tableaux complexes.
- API REST (FastAPI).

Merci de votre attention !

Annexe : Analyse de l'Empreinte Matérielle

RAM (Mémoire)

- **LLM local (Llama 3.2)**
- Quantification **4-bit** : réduit l'utilisation à ~4Go
- *Le modèle doit tenir entièrement en RAM*
- Recommandé : **8Go min / 16Go confort**

CPU (Processeur)

- **OCR** : Multi-thread (OpenMP) → tous cœurs sollicités
- **LLM** : CPU-Bound, calcul flottant massif
- Pas de GPU : instructions AVX/FMA utilisées
- *Latence expliquée par CPU uniquement*

Optimisation Logicielle

- Séparation Thread UI / Worker Thread
- UI toujours réactive (aucun gel lors calculs lourds)
- Exploitation efficace des multi-cœurs

