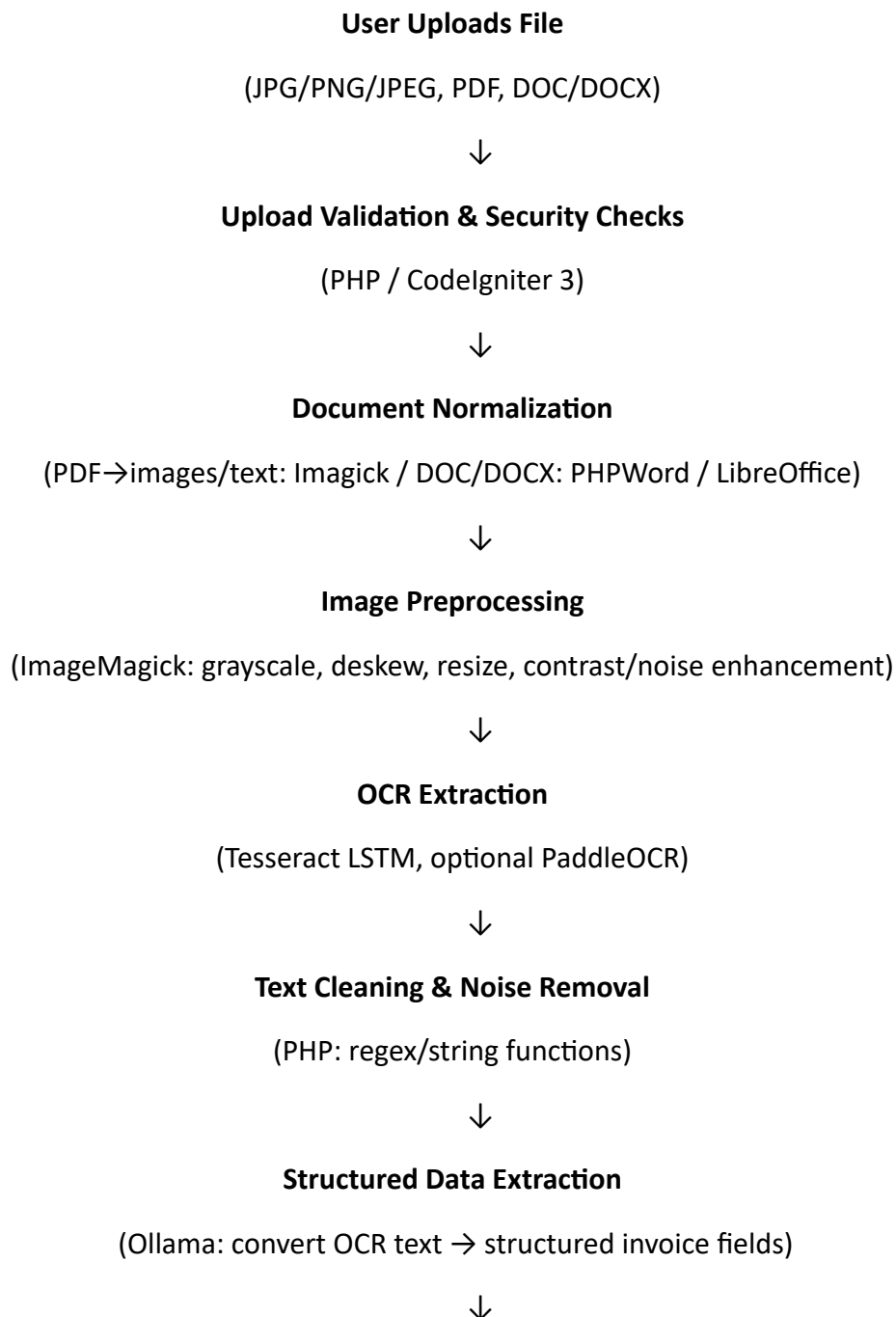


# OCR/AI INVOICE SCANNING FEATURES

AIM:

To automate the extraction of invoice data from user-uploaded files (Image, PDF, DOC) with high accuracy and minimal manual effort.

OVERFLOW



## Confidence Scoring

(PHP: combine OCR + Ollama confidence)



## Prefilled Editable Invoice Form

(PHP Views: show prefilled fields for high-confidence, empty fields for low-confidence; user can manually type/correct)



Save to **Database**

(MySQL )

LIBRARIES /PLATFORMS USED:

Tool	Free / Paid	Accuracy (Typed)	Handwritten	Limits	Offline	Why Use
ImageMagick	Free	⬆️ Boost OCR	⬆️ Boost	No limits	✓	Improves OCR accuracy before recognition
Imagick (PHP)	Free	⬆️ Conversion	⬆️ Conversion	No limits	✓	Best PHP-friendly document handling
Tesseract LSTM	Free	★ ★ ★ ★ ☆	★ ★ ★ ☆ ☆	No limits	✓	Best free OCR for printed invoices
PaddleOCR	Free	★ ★ ★ ★ ☆	★ ★ ★ ★ ☆	No limits	✓	Better handwriting & complex layout support
Ollama (Local LLM)	Free	★ ★ ★ ★ ★ (Semantic)	Fixes OCR errors	No limits	✓	AI understanding & structured

## PAID LIBRARIES (OPTIONAL):

Tool	Paid	Accuracy (Field Extraction)	Structured Invoice Output	Best For
<b>Veryfi OCR API</b>	● Paid	★★★★★ (~98–99% field-level)	✓ Yes (robust fields + line-item extraction)	Fast, invoice-focused automation with high accuracy and sub-3 s processing on average
<b>ABBYY FlexiCapture / FineReader</b>	● Paid	★★★★★ (enterprise-grade)	✓ Yes (excellent table & layout parsing)	Deep accuracy with configurable taxonomy, suited for complex global invoices
<b>Google Cloud Vision / Document AI</b>	● Paid	★★★★☆ (~94–98%)	✓ Yes (key-value extraction)	General OCR + structured extraction with broad language support

## PROCEDURE I DONE SO FAR

### (1) Imagick download

MY current:

(through test.php(localhost/test.php)

**PHP 7.4.1 + Thread Safety enabled + x64**

so installed Imagick: *PHP 7.4 TS x64 (vc15)* =extension

installed- *Ghostscript 10.06.0* for Windows (64 bit)

### (2) Tesseract



[tesseract-ocr-w64-setup-5.3.4.20240503.exe](#)

### (3) installed poppler [25.12.0-0](#)

Python executable:

C:\Program Files\Python311\python.exe

Python packages:

C:\Users\shaha\AppData\Roaming\Python\Python311\site-packages\

Python OCR project:

C:\ocr\_service\app.py

PHP project:

C:\xampp\htdocs\your\_php\_project\

Result:

## 1 Why the OCR text looks messy (this is expected)

Your invoice has **all of these together**:

- Arabic **RTL** text (headers, totals)
- English **LTR** text (items, totals in words)
- **Tables** (Item / Qty / Unit / VAT)
- **Boxes + grid lines**
- Mixed font sizes
- Logos + QR code

Tesseract **does not understand tables**.

It only outputs **reading order text**, not structure.

So this happens:

### What you see in PDF What Tesseract sees

Table rows	Random text blocks
Arabic RTL	Fragmented words
Grid borders	Noise
Item columns	Merged text

That's why:

- Qty / Unit / Price are **not aligned**
- Arabic words look broken
- Some English words are misspelled (Neety, tight)

### Correct decision matrix

Requirement	Tesseract	PaddleOCR
Read normal text	✓	✓
Read table text	⚠ Partial	✓ Better
Detect rows/cols automatically	✗	✗
Rebuild rows via code	⚠ Hard	✓ Easier
Offline	✓	✓
PHP-only	✓	✗ (Python)

#### 1) Upload → save in invoice\_documents (always)

**When user uploads** (image OR pdf), you must create **one row** in invoice\_documents immediately.

##### Save these columns

- original\_file\_path → assets/uploads/invoices/{docId}/original/{file}
- original\_file\_name → uploaded file name
- file\_ext → pdf/jpg/png/...
- mime\_type → detected mime
- file\_size → bytes
- page\_count → initially 1 (update later for pdf)
- status → 'uploaded'
- created\_at, updated\_at

✓ Result: you get **document\_id** (AUTO INCREMENT)

This document\_id will be used in tables 2 and 3.

---

#### 2) If PDF → normalize pages → save in invoice\_document\_pages

After upload, do normalization:

##### Case A: Uploaded file is IMAGE (jpg/png)

- Create **one normalized image**: normalized/page\_1.png (or keep original extension)
- Insert **one row** into invoice\_document\_pages:

### Fields

- document\_id = document\_id from step 1
  - page\_no = 1
  - image\_path = assets/uploads/invoices/{docId}/normalized/page\_1.png
  - ocr\_text = NULL (for now)
  - ocr\_confidence = NULL
  - created\_at = now
- Then update document:
- page\_count = 1
  - status = 'converted'
- 

### Case B: Uploaded file is PDF

- Convert PDF → multiple images:
    - normalized/page\_1.png, page\_2.png, ...
  - Insert **multiple rows** into invoice\_document\_pages (one per page)
- Example:
- page 1 row (page\_no=1, image\_path=..., ocr\_text null)
  - page 2 row (page\_no=2, image\_path=..., ocr\_text null)
  - ...

Then update document:

- page\_count = total pages
- status = 'converted'

✓ Your table already has a unique key (document\_id, page\_no) — perfect.

---

### 3) OCR stage → update invoice\_document\_pages

Now run OCR for each page row:

For each page:

- read image\_path
- run tesseract → get text
- compute confidence (optional)
- update that same row:

#### Update these fields in invoice\_document\_pages

- ocr\_text = extracted text
- ocr\_confidence = confidence number (0–100 or 0–1, but your column is decimal(5,2))

After OCR completes, update document:

- status = 'ocr\_done'
- 

### 4) Prefill → user edits → final save to invoices

When user clicks **Save Invoice**, insert into invoices.

**Must save**

- document\_id (foreign reference)
- all invoice fields: dates, order\_no, reference\_no, subject, totals, etc.
- set is\_confirmed = 0 (or 1 if you want)
- created\_at  
Then update document:
- status = 'saved'

Ethicfin Al Ghamdi Contracting Est. مؤسسة بن علي بن سعد آل محي الغامدي للمظاولات  
 Al-Muzahimiyah, Wasila, 19658 =: VALoA, العام «وسيلة»  
 CR : 798796645324 — ١٠١ NAAN : رقم السجل  
 VAT NO: 1224567456467 TYAN... Y : رقم الضريبة  
 TAX INVOICE / فاتورة ضريبية  
 Customer Name . ee ADIN SS اسم العميل  
 J' 0345 13D Al Manar Dist DAMMAM SAUDI ARABIA a.  
 SS ١ # 7600 العربية السعودية العنوان Due Date | 20-Jan-2026 |  
 تاريخ الاستحقاق  
 Address للعميل  
 ZIP/Postal Code: 32273  
 1. 10 days returns policy. 2. Risk passes to the buyer  
 upon delivery. Ownership of goods remains with the  
 Payment Terms company until full payment is received. 3. All شروط الدفع Our Ref رقم  
 المرجعي  
 applicable taxes, duties, and levies will be charged :  
 extra as per prevailing laws.  
 جموع امل قيمة الضريبة بين سعر الوحدة الوحدة الكمية الاسم والوصف رمز الصنف  
 #  
 \_ — . . Without VAT .  

Item Code	Item & Description	Qty	Unit	Price	VAT	Amount	Total inc. VAT
1}	Plo1168 me lel 100   Box	70.00		5,999.70	899.96	6,899.66	
Total Without VAT الاجمالي بدون ضريبة 7,000.00							
VAT Total 15% ضريبة القيمة المضافة 899.96							
Total With VAT in SAR المجموع شامل الضريبة 6,899.66							

Saudi Riyals Six Thousand Eight Hundred Ninety-Nine and Sixty-Six Halalas  
 Bank Account Details | Sid! تفاصيل الحساب  
 ly ORs 0  
 ank Name اسم الى F ral of Sta  
 BAN We BAN46875 اتات