

# OCR Pipeline Architecture

|                       |                 |
|-----------------------|-----------------|
| <b>Document ID</b>    | LT-OCR-ARCH-001 |
| <b>Version</b>        | 1.0             |
| <b>Classification</b> | Internal        |
| <b>Date</b>           | 2026-02-23      |

## 1. Executive Summary

The LT OCR Pipeline is a serverless document processing service responsible for payslip fraud detection, structured data extraction, validation, and final trust decision generation. It combines deterministic file-level forensics, rule-based risk scoring, OCR-powered data extraction, and LLM-based semantic interpretation to produce auditable, explainable outcomes.

The pipeline is designed to **always extract data when technically possible**, while separating **risk detection** from **final trust decisions**. Fraud detection is risk-based, not binary — deterministic rules are combined into a single risk score with outcomes of SUCCESS, NEEDS REVIEW, or FAILED.

This document provides a comprehensive architectural reference covering system design, data flow, processing stages, API contracts, security model, integration patterns, infrastructure, cost analysis, and operational considerations.

|                                      |                         |
|--------------------------------------|-------------------------|
| <b>Service Name:</b> lt-ocr-pipeline | <b>Language:</b> Python |
|--------------------------------------|-------------------------|

|                                     |   |
|-------------------------------------|---|
| <b>Compute:</b> AWS Lambda (Docker) | <b>Auth Provider:</b> AWS IAM / API Gateway |
|-------------------------------------|---|

|   |  |
|---|--|
| <b>API Protocol: REST / JSON (Asynchronous)</b> | <b>Messaging:</b> Amazon SQS                 |
| <b>OCR Engine: AWS Textract</b>                 | <b>LLM Provider:</b> Amazon Bedrock (Claude) |
| <b>Database: PostgreSQL (RDS)</b>               | <b>Storage:</b> Amazon S3                    |

## 2. System Overview

### 2.1 Purpose and Scope

The LT OCR Pipeline processes payslip documents to detect fraud, extract structured payroll data, validate correctness, and produce a final trust decision. It treats documents first as **digital artefacts**, then as **text**, and only later as **business data**.

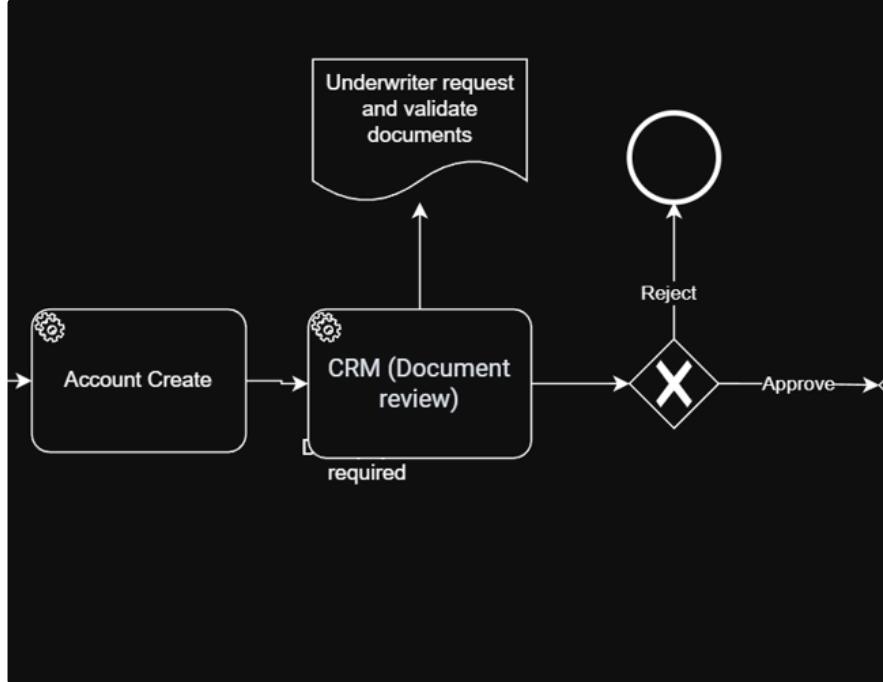
| Capability                           | Description  |
|--------------------------------------|--|
| <b>Deterministic Fraud Detection</b> | File-level digital forensics, layout anomaly detection, encoding analysis, and structural inspection — all rule-based and explainable. |
| <b>OCR Extraction</b>                | AWS Textract-powered text extraction with table and key-value pair recognition for structured payroll data.                            |
| <b>Fuzzy Label Mapping</b>           | Database-driven label dictionary with fuzzy matching to map OCR-detected labels to canonical internal IDs before LLM interpretation.   |
| <b>LLM-Based Semantic Extraction</b> | Amazon Bedrock (Claude) performs value extraction, context reasoning,  |

|  |   |
|--|---|
|  | multilingual handling, and schema completion.   |
| <b>Validation &amp; Trust Evaluation</b> | Deterministic validation rules for extraction usability, authenticity, consistency, and external employer verification. |
| <b>Risk-Based Scoring</b>                | All fraud signals merged into a single clamped score (0–100) with explainable flag breakdown.                           |
| <b>Full Auditability</b>                 | Complete traceability via process_id, OCR evidence references, rule execution logs, and fraud flag details.             |

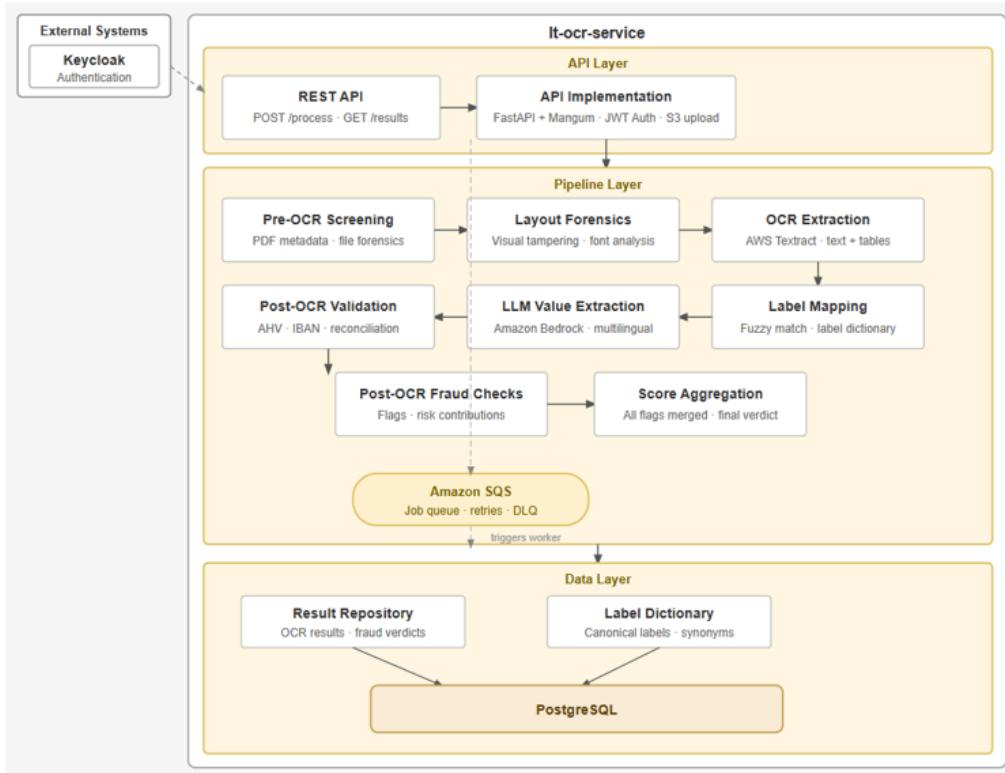
## 2.2 Architecture Diagram

The diagram below illustrates where the OCR pipeline is triggered within the overall process flow. When the process reaches the **CRM Document Review** stage, the system automatically submits the applicant's payslip to the OCR service for processing.

At that point, the OCR pipeline shown in the architecture diagram below takes over.



It accepts the document, runs it through all fraud detection and extraction stages in the background, and returns a structured result including extracted payroll data and a fraud verdict. This result is then surfaced back to the underwriter within the CRM to support the approval or rejection decision.



The pipeline follows these phases on the architecture

| Phase                      | Description   |
|----------------------------|---|
| <b>1 — Ingest</b>          | Customer Portal uploads payslip → S3 Bucket storage                   |
| <b>2 — Fraud Screening</b> | Fraud pattern recognition via deterministic forensic analysis         |
| <b>3 — OCR Extraction</b>  | AWS Textract + Amazon Bedrock (Claude) for structured data extraction |
| <b>4 — Validation</b>      | Business validation rules, identifier checks                          |
| <b>5 — Serve</b>           | Final consolidated data stored in DataHub                             |

## 2.3 Data Flow

The processing lifecycle follows an asynchronous pattern:

| Step | Action   |
|------|--|
| 1    | Client uploads payslip via API<br>Gateway → Lambda A   |
| 2    | Lambda A stores file in S3, inserts PENDING record in PostgreSQL, pushes job to SQS  |
| 3    | Lambda A returns ACCEPTED response immediately   |
| 4    | SQS triggers Lambda B (Worker)   |
| 5    | Lambda B executes full pipeline:<br>Fraud Detection → OCR → Label<br>Mapping → LLM Extraction →<br>Validation → ZEFIX → Score<br>Aggregation |
| 6    | Lambda B writes final result to PostgreSQL   |

## 3. Module Structure

The service follows a layered modular architecture deployed as Docker containers on AWS Lambda. It is organized into three concerns: entry points that handle communication, core services that execute the pipeline, and domain entities that represent the data.

### 3.1 Application Entry Points

The pipeline uses two separate Lambda functions sharing a single Docker image (different CMD handler overrides).

**Lambda A (It-ocr-api):** FastAPI application adapted for Lambda using Mangum. Handles client-facing operations: file upload, S3 storage, PostgreSQL record creation, SQS job dispatch, and immediate ACCEPTED response. Ensures response within API Gateway's 29-second timeout.

**Lambda B (It-ocr-worker):** Background processing Lambda triggered by SQS. Executes the full pipeline with a 15-minute (900-second) timeout to accommodate complete processing duration.

### 3.2 Core Services

| Service                      | Description  |
|------------------------------|--|
| <b>FraudDetectionService</b> | Orchestrates all deterministic fraud detection stages. Analyzes the document as a digital file (PDF metadata, structure, layout, encoding) without interpreting payroll meaning. Produces integrity flags, risk contributions, and hard failures for unreadable files. |
| <b>OCRService</b>            | Manages AWS Textract integration for text extraction. Executes DetectDocumentText for raw OCR and AnalyzeDocument (Tables) for structured key-value pairs and table  |

|                                |  |
|--------------------------------|--|
|                                | <p>data. Stores all OCR outputs unchanged.</p>   |
| <b>LabelMappingService</b>     | <p>Performs deterministic label detection using a database-driven label dictionary. Extracts candidate labels from OCR lines, applies fuzzy matching with configurable thresholds, and handles unmapped labels for continuous dictionary improvement.</p>        |
| <b>LLMExtractionService</b>    | <p>Integrates with Amazon Bedrock (Claude) for semantic value extraction and schema completion. Works with pre-mapped label_ids and OCR output. Handles multilingual content (DE/FR/IT/EN). Does NOT perform OCR, fraud risk assignment, or validation.</p>      |
| <b>ValidationService</b>       | <p>Applies deterministic validation rules for extraction usability (mandatory fields, parseability, schema compliance) and authenticity checks (salary logic, totals reconciliation, identifier validation). Each rule adds deterministic risk contribution.</p> |
| <b>FraudAggregationService</b> | <p>Merges all risk contributions into a single clamped score (0– 100). Derives final decision: SUCCESS, NEEDS REVIEW, or FAILED.</p>   |

### 3.3 Domain Entities

| Entity                  | Description  |
|-------------------------|--|
| <b>ProcessRecord</b>    | Represents a single document processing job. Contains process_id, document reference, status (PENDING, PROCESSING, SUCCESS, NEEDS REVIEW, FAILED), and timestamps. |
| <b>FraudResult</b>      | Contains the aggregated fraud assessment: risk_score (0–100), fraud_flags list, and validation_summary.  |
| <b>ExtractionResult</b> | Structured payroll data extracted via OCR and LLM. Includes label_id, extracted_value, and evidence references for traceability.                                   |
| <b>LabelDictionary</b>  | Configurable label mapping table with canonical names, language context, and known synonyms for fuzzy matching.  |

## 4. Data Model

**Database Type:** PostgreSQL (RDS)    **Storage:** Amazon S3

## 4.1 Key Query Patterns

| Query                           | Description   |
|---------------------------------|---|
| <b>Fetch Pending Jobs</b>       | Retrieves PENDING process records for monitoring and retry handling.  |
| <b>Label Dictionary Lookup</b>  | Retrieves canonical labels filtered by language and document type for fuzzy matching.                                     |
| <b>Process Result Retrieval</b> | Fetches complete processing result including fraud score, extraction data, and validation summary for a given process_id. |

## 4.2 Table: process\_record

Central table tracking the lifecycle of every document processing job.

| Column            | Type      | Description                                 |
|-------------------|-----------|---|
| <b>id</b>         | UUID (PK) | Unique identifier of the process record     |
| <b>process_id</b> | UUID      | External process identifier for correlation |
| <b>created_at</b> | TIMESTAMP | Job creation timestamp                      |
| <b>updated_at</b> | TIMESTAMP | Last status update timestamp                |

|                      |         |  |
|----------------------|---------|--|
| <b>status</b>        | VARCHAR | PENDING,<br>PROCESSING,<br>SUCCESS,<br>NEEDS REVIEW,<br>FAILED |
| <b>document_key</b>  | VARCHAR | S3 object key of the original document                         |
| <b>document_type</b> | VARCHAR | Document type hint (e.g., payslip)                             |
| <b>risk_score</b>    | INT     | Aggregated fraud risk score (0–100)                            |
| <b>decision</b>      | VARCHAR | Final outcome:<br>SUCCESS,<br>NEEDS REVIEW,<br>FAILED          |

|                           |       |   |
|---------------------------|-------|---|
| <b>fraud_flags</b>        | JSONB | List of triggered fraud flags with descriptions     |
| <b>extraction_data</b>    | JSONB | Structured payroll data extracted from the document |
| <b>validation_summary</b> | JSONB | Breakdown of validation checks and results          |
| <b>error_details</b>      | JSONB | Error information if processing failed              |

#### 4.3 Table: label\_dictionary

Stores the canonical label mapping used for fuzzy matching during label detection.

| Column                | Type         | Description                             |
|-----------------------|--------------|---|
| <b>id</b>             | UUID (PK)    | Primary key                             |
| <b>label_id</b>       | VARCHAR (UQ) | Internal canonical identifier           |
| <b>canonical_name</b> | VARCHAR      | Standard payroll field name             |
| <b>language</b>       | VARCHAR      | Language context (DE, FR, IT, EN)       |
| <b>synonyms</b>       | JSONB        | Known label variants for fuzzy matching |
| <b>document_type</b>  | VARCHAR      | Applicable document type scope          |
| <b>is_active</b>      | BOOLEAN      | Toggle for enabling/disabling the label |

#### 4.4 Table: unmapped\_label\_log

Tracks OCR labels that could not be matched to the dictionary, enabling continuous improvement.

| Column                | Type      | Description                      |
|-----------------------|-----------|----------------------------------|
| <b>id</b>             | UUID (PK) | Primary key                      |
| <b>process_id</b>     | UUID (FK) | Links to process_record          |
| <b>raw_label_text</b> | VARCHAR   | Original OCR-detected label text |

|                       |         |  |
|-----------------------|---------|--|
| <b>top_candidates</b> | JSONB   | Top N candidate label_ids with similarity scores |
| <b>ocr_evidence</b>   | JSONB   | OCR snippet and bounding box reference           |
| <b>resolved</b>       | BOOLEAN | Whether the label has been manually reviewed     |

#### 4.5 Table: external\_validation\_log

Logs all calls to external services for audit, debugging, and performance monitoring.

| Column               | Type      | Description                            |
|----------------------|-----------|--|
| <b>id</b>            | UUID (PK) | Primary key                            |
| <b>process_id</b>    | UUID (FK) | Links to process_record                |
| <b>provider_name</b> | VARCHAR   | Name of external system                |
| <b>status</b>        | VARCHAR   | SUCCESS, FAILURE, TIMEOUT              |
| <b>latency_ms</b>    | INT       | Execution time in milliseconds         |
| <b>raw_response</b>  | JSONB     | Raw payload received from the provider |

## 5. API Specification

|                                   |   |
|-----------------------------------|---|
| <b>Protocol:</b> REST / JSON      | <b>Security:</b> AWS IAM / API Gateway      |
| <b>Documentation:</b> OpenAPI 3.0 | <b>Design:</b> Asynchronous (Accept → Poll) |

### 5.1 Document Processing Endpoint

|             |                        |
|-------------|------------------------|
| <b>POST</b> | <b>/v1/ocr/process</b> |
|-------------|------------------------|

Submits a payslip document for asynchronous processing. Returns immediately with an ACCEPTED status.

#### Request Structure

```
{"process_id": "uuid",
"document_key": "string",
"document_type": "payslip", "filename": "string"}
```

#### Response Structure (Immediate)

```
"process_id": "uuid",
```

```
"status": "ACCEPTED"
```

## 5.2 Result Retrieval Endpoint

|     |                             |
|-----|-----------------------------|
| GET | /v1/ocr/result/{process_id} |
|-----|-----------------------------|

Retrieves the processing result for a submitted document.

### Response Structure (Completed)

```
{"process_id": "uuid",
  "status": "SUCCESS", "risk_score": 15,
  "decision": "SUCCESS", "fraud_flags": [
    "PDF_METADATA_SUSPICIOUS_PRODUCER",
    "ENCODING_MIXED_APOSTROPHES"
  ],
  "extraction_data": {
    "gross_salary": 8500.00,
    "net_salary": 6200.00,
    "employer_name": "Example AG"
  },
  "processing_time_ms": 12400}
```

## 6. Security Model

### 6.1 Authentication & Authorization

The LT OCR Pipeline is secured through AWS-native mechanisms. API Gateway handles external client authentication, while Lambda functions operate under IAM roles with least-privilege policies.

|  |  |
|--|--|
| <b>Protocol: AWS IAM / API Gateway</b> | <b>Secrets:</b> AWS Secrets Manager                    |
| <b>Network: VPC with NAT Gateway</b>   | <b>Database Access:</b> VPC Security Groups (TCP 5432) |

### 6.2 IAM Role Permissions

| Role                    | Permission   | Assigned To    |
|-------------------------|--|----------------|
| <b>lt-ocr-api- role</b> | S3 PutObject, SQS SendMessage, RDS Connect, Secrets Manager Read | Lambda A (API) |

|                           |  |                   |
|---------------------------|--|-------------------|
| <b>lt-ocr-worker-role</b> | S3<br>GetObject/HeadObject,<br>Textract<br>AnalyzeDocument/DetectDocument, Bedrock<br>InvokeModel, RDS<br>Connect, Secrets Manager Read, SQS<br>ReceiveMessage | Lambda B (Worker) |
|---------------------------|--|-------------------|

### 6.3 Network Security

Lambda functions are placed inside the VPC with NAT Gateway for outbound access. RDS security group allows inbound TCP 5432 only from Lambda security groups.

Sensitive credentials (PostgreSQL password, API keys) are stored in AWS Secrets Manager and retrieved at runtime — no plaintext values in environment variables.

## 7. Evaluation Workflow

The LT OCR Pipeline executes a multi-phase asynchronous processing workflow. The client receives an immediate acknowledgment, and the full pipeline runs in the background via SQS-triggered Lambda.

| Step | Phase                    | Details   |
|------|--------------------------|---|
| 1    | <b>Request Reception</b> | Client sends document processing request to API Gateway → Lambda A. File is uploaded to S3, |

|   |  |   |
|---|--|---|
|   |  | PENDING record created in PostgreSQL, job message pushed to SQS. Client receives ACCEPTED response immediately.   |
| 2 | <b>File-Level Forensics</b>                | Document analyzed as a digital artefact. S3 metadata (HeadObject) and raw bytes (GetObject) inspected. Checks: file validity, PDF metadata (Author, Creator, Producer, CreationDate, ModDate), structural indicators, layout artefacts, and encoding anomalies. |
| 3 | <b>OCR Execution</b>                       | AWS Textract processes the document. DetectDocumentText for raw OCR. AnalyzeDocument (Tables) for structured key-value pairs and table data. OCR quality gate stops processing on weak results. All outputs stored unchanged.                                   |
| 4 | <b>Label Detection &amp; Fuzzy Mapping</b> | Candidate labels extracted from OCR lines using bounding box proximity and row  |

|   |  |   |
|---|--|---|
|   |  | alignment. Fuzzy matching against label dictionary. Unmatched labels stored with top N candidate suggestions for manual review.   |
| 5 | <b>LLM-Based Extraction</b>              | Amazon Bedrock (Claude) receives OCR output and pre-mapped label_ids. Performs semantic value extraction, context reasoning, multilingual handling (DE/FR/IT/EN), and schema completion.  |
| 6 | <b>Validation &amp; Trust Evaluation</b> | Extraction usability checks (mandatory fields, parseability, schema compliance). Authenticity checks (salary logic, reconciliation, identifier validation).   |
| 7 | <b>Risk Aggregation &amp; Decision</b>   | All risk contributions merged: file-level + layout + arithmetic + identifier. Score clamped to 0–100. Final decision: SUCCESS (low), NEEDS REVIEW (medium/high), FAILED (technical failure only). Result persisted to PostgreSQL. |

## 8. Fraud Detection Detail

### 8.1 Deterministic Checks

| Category                     | What Is Checked                                  | Example Risk Signals   |
|------------------------------|--|--|
| <b>File Validity</b>         | Format, readability, bounds                      | Unsupported format, unreadable PDF, abnormal page count                          |
| <b>PDF Metadata</b>          | Author, Creator, Producer, CreationDate, ModDate | Word/LibreOffice producer, recent modification for old payroll, missing metadata |
| <b>Structural Indicators</b> | PDF internal structure                           | Incremental update chains, unusual object streams, missing font embedding        |
| <b>Layout Artefacts</b>      | Visual rendering consistency                     | Misaligned tables, irregular margins, overlapping text                           |
| <b>Encoding Anomalies</b>    | Text-level irregularities                        | Mixed apostrophes (' vs '), font switches mid-word, mixed decimal separators     |
| <b>Arithmetic Checks</b>     | Payroll consistency                              | Gross vs Net plausibility, totals reconciliation, negative totals                |

|                              |                |   |
|------------------------------|----------------|---|
| <b>Identifier Validation</b> | UID, AHV, IBAN | Format validation,<br>checksum verification |
|------------------------------|----------------|---|

## 8.2 Libraries Used

| Library         | Role                                    |
|-----------------|---|
| <b>PyMuPDF</b>  | PDF parsing, rendering, font inspection |
| <b>pdfminer</b> | Structural analysis                     |
| <b>Pillow</b>   | Image-level layout inspection           |

## 8.3 Final Decision Logic

| Risk Level        | Outcome      |
|-------------------|--------------|
| Low               | SUCCESS      |
| Medium / High     | NEEDS REVIEW |
| Technical failure | FAILED       |

**Note:** Fraud or authenticity concerns **do not invalidate extracted data**. Only technical failures stop processing.

## 9. Design Patterns and Principles

| Pattern                                   | Description  |
|---|--|
| <b>Asynchronous Processing</b>            | Client receives immediate ACCEPTED response. Full pipeline runs asynchronously via SQS-triggered Lambda. Decouples API responsiveness from processing duration.  |
| <b>Digital-First Analysis</b>             | Documents are treated first as digital artefacts (file forensics), then as text (OCR), and only later as business data (payroll extraction). This layered approach maximizes fraud detection coverage. |
| <b>Risk-Based Scoring</b>                 | Fraud detection is probabilistic, not binary. All deterministic rules contribute to a single aggregated score (0–100), enabling nuanced decisions rather than hard pass/fail.                          |
| <b>Separation of Extraction and Trust</b> | Data extraction always proceeds when technically possible. Risk scoring is independent — fraud concerns flag but do not prevent data output.   |

|  |  |
|--|--|
| <b>Fuzzy Matching with Continuous Learning</b> | Label detection uses fuzzy matching with configurable thresholds. Unmapped labels are logged with candidate suggestions, enabling dictionary expansion and continuous improvement.                               |
| <b>LLM as Semantic Layer Only</b>              | The LLM handles value extraction and schema completion but does NOT perform OCR, fraud risk assignment, or validation. Clear responsibility boundaries prevent LLM hallucination from affecting trust decisions. |
| <b>Deterministic Explainability</b>            | All fraud scoring is rule-based and explainable. Every risk contribution is traceable to a specific check, enabling audit and manual review.   |
| <b>Auditability by Design</b>                  | Every processing step is logged: OCR evidence references, label mapping decisions, LLM extraction with evidence linking, validation results, and fraud flags.  |

## 10. Infrastructure

| Service                | Role        | Details   |
|------------------------|-------------|---|
| <b>AWS API Gateway</b> | Entry point | Exposes lt-ocr-api as public HTTPS endpoint. Routes all |

|                                     |                    |  |
|-------------------------------------|--------------------|--|
|                                     |                    | client requests to Lambda A. Hard timeout of 29 seconds.   |
| <b>AWS Lambda A (It-ocr-api)</b>    | Client-facing API  | FastAPI + Mangum. Receives uploads, stores to S3, creates PENDING record, pushes to SQS, returns ACCEPTED.                       |
| <b>AWS Lambda B (It-ocr-worker)</b> | Background worker  | SQS-triggered. Runs full pipeline. 15-minute timeout (900s).   |
| <b>Amazon SQS</b>                   | Job queue          | Decouples Lambda A from Lambda B. Provides automatic retries and DLQ support.  |
| <b>Amazon ECR</b>                   | Container registry | Single Docker image for both Lambdas (different CMD overrides). Required due to dependencies exceeding Lambda's 250MB zip limit. |
| <b>AWS Textract</b>                 | OCR engine         | DetectDocumentText (raw OCR) + AnalyzeDocument Tables (structured KV pairs and tables).  |
| <b>Amazon Bedrock (Claude)</b>      | LLM provider       | Semantic value extraction, schema completion, and visual fraud analysis.   |

|                            |                    |  |
|----------------------------|--------------------|--|
| <b>Amazon S3</b>           | Document storage   | Original files stored unchanged. Pipeline never modifies source documents.     |
| <b>PostgreSQL (RDS)</b>    | Database           | Process records, extraction results, fraud scores, validation summaries.       |
| <b>Amazon CloudWatch</b>   | Observability      | Logs from both Lambdas. Debugging, metrics, alarms. 30-day log retention.      |
| <b>AWS Secrets Manager</b> | Credential storage | PostgreSQL password, API keys. Retrieved at runtime, no plaintext in env vars. |

## 11. Configuration Checklist

| # | Configuration                                       | Service  |
|---|---|----------|
| 1 | Grant Lambda roles access to S3 bucket              | S3 / IAM |
| 2 | Enable Anthropic Claude Sonnet + Haiku model access | Bedrock  |

|   |  |              |
|---|--|--------------|
| 3 | Place lt-ocr-api Lambda inside the VPC | VPC / Lambda |
|---|--|--------------|

|    |   |                       |
|----|---|-----------------------|
| 4  | Place lt-ocr-worker Lambda inside the VPC               | VPC / Lambda          |
| 5  | Open TCP port 5432 on RDS security group for Lambda     | VPC / Security Groups |
| 6  | Confirm NAT Gateway exists in VPC (or create one)       | VPC                   |
| 7  | Share RDS endpoint, DB name, username, password         | RDS                   |
| 8  | Share VPC ID and Subnet IDs                             | VPC                   |
| 9  | Attach IAM policies to lt-ocr-api role                  | IAM                   |
| 10 | Attach IAM policies to lt-ocr-worker role               | IAM                   |
| 11 | Set CloudWatch log retention to 30 days on both Lambdas | CloudWatch            |

Regardless of any pending decisions on the overall infrastructure setup, the following two AWS services are very much needed. These are independent of deployment, VPC, or any other architectural choices still under discussion.

**AWS Textract — AnalyzeDocument (Tables)** Used for structured data extraction from payslip documents. Must be accessible via our IAM role.

**Amazon Bedrock — Claude Model Access** Used for intelligent data interpretation and schema completion. The following models must be explicitly enabled in the Bedrock console under

**Model access:**

| Model | Version |
|-------|---------|
|       |         |

Claude Sonnet

anthropic.claude-3-5-

sonnet-20241022-v2:0

**Amazon S3 — Development Bucket** We will need access to an S3 bucket for development and testing purposes. To avoid any risk to production documents, we would propose access to the dev bucket with read/write access. Please advise on what works best on your side.