

# Low-Level Design (LLD) – ALL Classifier

---

## 1 Source Tree

```
.
├── airflow/dags/
│   └── all_classification_pipeline.py
├── backend/
│   ├── main.py
│   ├── model_utils.py
│   └── Dockerfile
├── frontend_streamlit/
│   ├── app.py
│   └── Dockerfile
├── monitoring/
│   ├── prometheus.yml
│   └── grafana/provisioning/
└── docker-compose.*.yml
```

## 2 FastAPI Specification

| Path            | Method | Request   | Response  | HTTP Code |
|-----------------|--------|---|---|-----------|
| /predict        | POST   | multipart/form-data {file: image}   | {<br>probabilities:<br>float[4],<br>predicted_class:<br>int } | 200       |
| /feedback       | POST   | multipart/form-data {<br>correct_label:str,<br>predicted_label:<br>str,<br>file:image } | {<br>stored:bool,<br>feedback_count:int<br>}                  | 200       |
| /feedback/count | GET    | –   | {<br>feedback_count:int<br>}                                  | 200       |
| /metrics        | GET    | –   | Prometheus exposition format                                  | 200       |

```
/ping      GET      -      {status:" 200  
ok" }
```

## Exception Handling

```
@app.exception_handler(Exception)  
  
async def fallback_ex(exc:Exception, req:Request):  
  
    logger.exception("unhandled error")  
  
    return JSONResponse({"detail":"server error"}, status_code=500)
```

## 3 SQLite Schema

```
CREATE TABLE feedback(  
  
    id INTEGER PRIMARY KEY AUTOINCREMENT,  
  
    ts REAL,          -- unix timestamp  
  
    correct TEXT,      -- ground-truth class  
  
    predicted TEXT,    -- model prediction at time  
  
    img BLOB          -- orig. JPEG/PNG bytes  
  
);
```

DB file path: \${FEEDBACK\_DIR:- /app/data}/feedback.db (volume-mounted in prod stack for extraction).

## 4 Airflow DAG – all\_classification\_pipeline

| Task-id           | Operator       | Callable / Cmd               | Key output                          |
|-------------------|----------------|------------------------------|-------------------------------------|
| preprocess_images | PythonOperator | resize_images( .<br>..)      | resized images →<br>/data/processed |
| prepare_data      | PythonOperator | make_manifest( .<br>..)      | CSV manifests                       |
| train_classifier  | PythonOperator | train_classifier(...)        | run_id (XCom)                       |
| test_classifier   | PythonOperator | test_classifier(..., run_id) | metrics                             |
| register_model    | PythonOperator | register_model(..., run_id)  | MLflow model reg.                   |

### Retrain-on-Feedback Sensor

```
from airflow.sensors.python import PythonSensor
```

```
def enough_feedback(**_):
```

```
    conn = sqlite3.connect('/app/data/feedback.db')
```

```
cnt = conn.execute('SELECT COUNT(*) FROM feedback').fetchone()[0]
```

```
return cnt >= 100
```

```
wait_feedback = PythonSensor(
```

```
    task_id    = 'wait_feedback',
```

```
    python_callable = enough_feedback,
```

```
    poke_interval  = 3600, # hourly
```

```
    timeout        = 86400, # 1 day window
```

```
)
```

```
wait_feedback >> preprocess_images
```

*(In dev DAG only; omitted from client-facing stack)*

## 5 Monitoring Metrics

| Metric  | Type      | Labels                  | Purpose         |
|---|-----------|-------------------------|-----------------|
| <code>http_requests_total</code>                  | Counter   | method, handler, status | requests volume |
| <code>http_request_duration_seconds_bucket</code> | Histogram | handler                 | latency P50/P95 |
| <code>process_cpu_seconds_total</code>            | Counter   | –                       | CPU usage       |
| <code>node_memory_MemAvailable_bytes</code>       | Gauge     | –                       | free mem        |

Grafana-ready dashboard JSON

(`monitoring/grafana/dashboards/all-api.json`) provisioned on container start.

## 6 Docker Images

| Image       | Base                   | Entrypoint   |
|-------------|------------------------|--|
| back<br>end | python:3.9-slim        | uvicorn main:app --host 0.0.0.0<br>--port 8000                         |
| streamlit   | python:3.9-slim        | streamlit run app.py<br>--server.port 8501<br>--server.address 0.0.0.0 |
| mlflow      | python:3.9-slim        | mlflow server --backend-store-uri<br>sqlite:///mlruns.db ...           |
| prometheus  | prom/prometheus:latest | /bin/prometheus<br>--config.file=/etc/prometheus/prometheus.yml        |

## 7 Environment Variables Summary

| Variable    | Default             | Used By   | Meaning               |
|-------------|---------------------|-----------|-----------------------|
| BACKEND_URL | http://backend:8000 | Streamlit | API base for requests |

|               |           |                   |                               |
|---------------|-----------|-------------------|-------------------------------|
| FEEDBACK_DIR  | /app/data | Backend           | where to put<br>feedback.db   |
| FB_THR<br>ESH | 100       | Airflow<br>sensor | min samples before<br>retrain |

---

*(end of LLD)*