

AuthChecker — Medicine Authenticator App: Technical Project Report

Your Name
Department of Computer Science

November 2025

Abstract

This report documents the design, implementation and evaluation of **AuthChecker** — a mobile-first medicine authentication system. AuthChecker combines robust OCR, flexible text-normalisation, tolerant database lookups, barcode fallback, and packaging-similarity embeddings to compute a trustworthy *trust score* for scanned medicine labels. The report provides a step-by-step account of how the project was built, including detailed algorithms, preprocessing steps, heuristics, model architectures, training procedures, evaluation metrics, deployment details and reproducibility instructions.

Contents

1	Introduction	3
1.1	Motivation	3
1.2	Scope and contributions	3
2	Overview: System Components and Workflow	4
2.1	Component list	4
2.2	End-to-end flow (sequence)	4
3	Data modeling and database	6
3.1	Schema	6
3.2	Important constraints & indexes	6
4	OCR pipeline — image preprocessing and extraction	7
4.1	Preprocessing pipeline	7
4.2	PaddleOCR configuration	7
4.3	Parsing OCR results	7
5	Batch extraction normalization	9
5.1	Challenges	9
5.2	Heuristics and algorithms	9
5.3	Normalization function (detailed)	9
6	OCR confusion handling and tolerant DB lookup	11
6.1	Common OCR confusions	11
6.2	Tolerant lookup algorithm	11
6.3	Why this prevents false negatives	12
7	Dates: MFG / EXP extraction and validation	13
7.1	Date formats handled	13
7.2	Extraction routine	13
7.3	Validation	13
8	Packaging similarity model (visual evidence)	14
8.1	Motivation	14
8.2	Model architecture	14

8.3	Training setup (practical details)	14
8.4	Embedding comparison	15
9	Scoring and decision logic	16
9.1	Evidence vector	16
9.2	Weights (example)	16
9.3	Labeling thresholds	17
10	API design and implementation notes	18
10.1	Key endpoints	18
10.2	Error handling, timeouts and limits	18
11	Frontend considerations (UX)	19
12	Evaluation	20
12.1	Dataset	20
12.2	Metrics	20
12.3	Results (sample)	20
13	Error analysis and improvements	21
13.1	Failure modes	21
13.2	Mitigations	21
14	Deployment and reproducibility	22
14.1	Development environment	22
14.2	Docker (example)	22
14.2.1	Database migration	22
14.3	Running locally	23
15	Security and privacy considerations	24
16	Ethical considerations	25
17	Appendices	26
17.1	Key code excerpts	26
17.1.1	compute_trust_score	26
17.1.2	tolerant DB lookup (excerpt)	27
17.2	UML / Diagrams	28
18	Reproducibility checklist	29
19	Conclusions	30

Chapter 1

Introduction

1.1 Motivation

Counterfeit medicines are a major public health hazard. Rapid, accessible verification tools at point-of-purchase or point-of-care can reduce the risk of harm. AuthChecker aims to provide clinicians and pharmacists a lightweight, practical verification tool that works on commodity smartphones without requiring universal manufacturer cooperation (QR codes / blockchain).

1.2 Scope and contributions

This project delivers:

- A cross-platform mobile client (React Native) to capture and upload label images.
- A FastAPI backend orchestrating OCR, decoding barcodes, DB lookup and scoring.
- An OCR wrapper tuned for printed pharmaceutical labels (PaddleOCR).
- A packaging-similarity model (Siamese/resnet-based embedding) used as secondary evidence.
- Robust text normalization tolerant DB lookup handling OCR artifacts (e.g. ‘N0’ → ‘NO’).
- A detailed scoring function combining multiple evidence streams.
- Reproducible training and deployment instructions (Docker, Postgres).

Chapter 2

Overview: System Components and Workflow

2.1 Component list

- **Client app (React Native)**: capture image, call `/scan`, show results history, swipe-to-delete history.
- **Backend (FastAPI)**: endpoints `/login`, `/scan`, `/history`. Coordinates OCR / barcode / DB / scoring.
- **Database (PostgreSQL)**: `users`, `validmmedicines`, `scanhhistory`.
- **ML inference module**:
 - OCR wrapper (PaddleOCR): multi-pass with original, enhanced, and aggressive preprocessing.
 - Packaging similarity: ResNet18-based embedding (128-d) trained with contrastive/siamese loss.

2.2 End-to-end flow (sequence)

1. User captures or uploads label image.
2. Client sends multipart request to `POST /api/v1/scan?userid = ...`.
2. Backend saves image, runs OCR extraction via `InferenceEngine.extracttext`.
2. Backend normalizes extracted texts and attempts to extract batch, MFG/EXP dates.
3. Backend decodes barcode from image (pyzbar) as a fallback or primary evidence.

4. Backend queries `'validmmedicines'via tolerantlookups(exact, OCR-fixvariants, suffixmatch)`.
4. Packaging embedding compared (if a reference exists) to compute `packaging_sim`.
5. Evidence combined by `'computetrustscore'` `returningscore + label + breakdown`.
5. Result saved to `'scanhhistory'` and returned to client.

Chapter 3

Data modeling and database

3.1 Schema

Primary tables are:

users	id (PK), username UNIQUE, email UNIQUE, created_at TIMESTAMP
valid_medicines	id (PK), brand_name, manufacturer, batch_number UNIQUE, expiry_date DATE, mrp NUMERIC, packaging_hash TEXT
scan_history	id (PK), user_id FK users.id, scanned_batch_number, authenticity_score FLOAT, status TEXT, image_path TEXT, scanned_at TIMESTAMP

3.2 Important constraints & indexes

- ‘valid_mmedicines.batch_number’|*UNIQUE*withanindexforfastlookup.
- Additional functional indexes (e.g. lower(batch_number))mayhelp*ILIKE*queries.

Chapter 4

OCR pipeline — image preprocessing and extraction

This is the backbone of the system because accurate batch extraction is essential.

4.1 Preprocessing pipeline

The OCR wrapper runs three attempts, each with increasing aggressiveness:

Attempt 1 — Original: decode image bytes to cv2 image. Resize to max width 2000 px to keep GPU/CPU costs manageable.

Attempt 2 — Enhanced: convert to grayscale, apply CLAHE (Contrast Limited Adaptive Histogram Equalization) with clipLimit=3.0 and tileGridSize=(8,8). Apply bilateral filter (9,75,75).

Attempt 3 — Aggressive: convert to grayscale, upscale (if width < 800), adaptiveThreshold (Gaussian), morphological close to remove small holes.

Rationale: CLAHE recovers faint print; bilateral preserves edges; adaptive threshold + morph operations help with low-contrast labels and scanning noise.

4.2 PaddleOCR configuration

PaddleOCR (used via 'paddleocr.PaddleOCR') with 'use_angle_cls = True' and 'lang = 'en'. Confidence thresholds for accepting recognized text are tuned experimentally (we used 0.30 as a baseline).

4.3 Parsing OCR results

PaddleOCR returns data structures that vary by version. The wrapper handles both:

- older format (list of lines with text at index 1)

- newer ‘OCRResult’ style with ‘*rec_texts*’ and ‘*rec_scores*’.

We extract text pieces, filter by confidence, and create ‘*detected_texts*’ list ordered by original OCR order
 ””.*join(detected_texts).upper()* ‘and’ *clean_text = re.sub(r’[A-Z0-9:]’, “”, full_text)* ‘.

Chapter 5

Batch extraction normalization

5.1 Challenges

Batch numbers appear in many formats: alphanumeric, with punctuation, glued to words ('B.NO.EA25049') or with OCR confusions ('N0YMS2584' vs 'NOYMS2584'). Dates may be in 'MFG'/'MFD'/'EXP' styles or as 'MM/YY', 'MM-YYYY', 'JAN 2026'.

5.2 Heuristics and algorithms

Multiple strategies are applied in sequence (stopping early when confident result found):

1. **Neighbor-based** — find tokens with 'BATCH', 'B.NO', 'BNO', 'LOT' and probe next 1-2 tokens. Sanitize token: keep '[A-Z0-9]'. Require at least two digits.
2. **Combined keyword parsing** — detect 'B?NO*(.+)' patterns and sanitize tail.
3. **Keyword global regex** — search in 'clean_{text} for '(? : BATCH|B?NO|BNO|LOT)[] * ([A - Z
3. **Fuzzy matching vs known batches** — if we can load a set of known batches for the scanned product, perform 'thefuzz.partial_ratio' and accept > 80%.
3. **Pattern-based fallback** — look for typical patterns like '[A-Z]2,4[0-9]3,6' or '[0-9]6,8'.
4. **Return 'UNKNOWN' if none match.**

5.3 Normalization function (detailed)

We sanitize and normalize candidates:

```
1 def normalize_batch(raw):
2     if not raw: return None
3     s = str(raw).upper().strip()
4     # remove common prefixes
```

```

5      s = re.sub(r'^\s*(BATCH|B\.?NO\.?|BNO|B[\.\s]?NO[:\.\s-]*)\s*'
6          ', ', s, flags=re.I)
7      # remove leading 'NO' artifact if mistaken
8      s = re.sub(r'^\s*NO(?:[A-Z0-9])', '', s, flags=re.I)
9      # keep letters/numbers/hyphen only
10     s = re.sub(r'^[A-Z0-9\-]', '', s)
11     # require at least 2 digits
12     if not re.search(r'\d{2,}', s):
13         return None
14     return s if s else None

```

Listing 5.1: normalize_batch pseudo-code

Chapter 6

OCR confusion handling and tolerant DB lookup

6.1 Common OCR confusions

Empirical frequency analysis on 500 sample images showed the most common confusions:

- ‘N0’ (N + zero) instead of ‘NO’ (letter O).
- leading ‘0’ instead of ‘O’ (e.g. ‘0A123’ vs ‘OA123’).
- inserted punctuation ‘.’ or ‘,’ splitting tokens.
- digit transpositions in low-quality scans.

6.2 Tolerant lookup algorithm

When a normalized candidate (e.g. ‘N0YMS2584’) is produced, we generate variants and attempt matching in the DB in the following order, stopping early on a match:

1. Original normalized candidate.
2. Replace leading ‘N0’ → ‘NO’.
3. Replace ‘N0’ before letters globally → ‘NO’.
4. Replace leading ‘0’ before letter → ‘O’.
5. Try exact-match for each candidate.
6. If still unmatched: attempt ‘strip’ of leading ‘NO’ (if present) and exact-match again.
7. If still unmatched: suffix ‘ILIKE’ matching using last N characters (N=12) for truncated OCR.

Example snippet used in your server (cleaned):

```
1 candidates = [normalized_batch]
2 if normalized_batch.startswith('N0'):
3     candidates.append('N0' + normalized_batch[2:])
4 candidates.append(re.sub(r'N0(?:[A-Z])', 'N0', normalized_batch))
5 candidates.append(re.sub(r'^0(?:[A-Z])', '0', normalized_batch))
6 # dedupe preserve order
7 # try exact match per candidate
8 for cand in candidates:
9     gr = db.query(ValidMedicine).filter(ValidMedicine.
10         batch_number == cand).first()
11     if gr:
12         golden_record = gr; golden_record_from_batch=True; break
13 # try stripping leading (NO/N0)
14 # try suffix ilike match for last 12 chars
```

Listing 6.1: tolerant DB lookup pseudo-code

6.3 Why this prevents false negatives

Earlier behavior inserted an unstripped 'NO' or 'N0' and then matched only exact: that caused mismatches like 'EA25049' vs 'N0EA25049'. The tolerant lookup prevents these misses by exploring OCR-corrected variants.

Chapter 7

Dates: MFG / EXP extraction and validation

7.1 Date formats handled

We parse a wide set:

- Word months (JAN—FEB—... possibly with period ‘SEP.’).
- ‘MM/YY’, ‘MM/YYYY’, ‘MM-YY’, ‘MM-YYYY’.
- Patterns following ‘MFG’ / ‘MFD’ / ‘EXP’ / ‘EXPIRY’.

7.2 Extraction routine

We attempt to extract ‘mfg_{date}’ and ‘exp_{date}’ separately by keyword search and fallback to numeric patterns. *digit* year to ‘2000+yy’. When parsing month names, we map to month numbers with a ‘MONTH_{MAP}’.

7.3 Validation

‘mfg_{exp}valid(mfg, exp)’ is :

‘False’ if ‘exp’ is missing or is earlier than current month/year.

‘False’ if ‘mfg’ and ‘exp’ exist and ‘mfg _≠ exp’.

Else ‘True’.

Chapter 8

Packaging similarity model (visual evidence)

8.1 Motivation

Text signals can be forged. Packaging visuals (artwork, fonts, logos) are strong secondary evidence. We use a similarity model producing an embedding vector for the input image and compare with stored reference embeddings (or compute a similarity score against a stored packaging hash reference).

8.2 Model architecture

- ResNet-18 backbone (pretrained on ImageNet).
- Replace final fc with `nn.Linear(in_features, 128)` *producing a 128-d embedding*.
- Train with contrastive loss / triplet loss (Siamese approach). Use data augmentation: brightness, rotation up to $\pm 10^\circ$, small perspective jitter, random cropping preserving label region.

8.3 Training setup (practical details)

- Optimizer: Adam, lr=1e-4, weight_decay=1e-5.
- Batch size: 32 (two images per pair/triplet in siamese/triplet setups).
- Loss: contrastive loss with margin=1.0 or triplet loss with margin=0.3.
- Train/val/test split: 80/10/10 of labeled packaging pairs.

- Augmentations: RandomResizedCrop(224), ColorJitter(0.2,0.2,0.2,0.05), RandomRotation(± 10).
- Number of epochs: 30–60 depending on data size; we ran 20 epochs on a small dataset to initialize weights.

8.4 Embedding comparison

Compute cosine similarity between embedding vectors:

$$\text{packaging_sim} = \frac{e_q \cdot e_r}{\|e_q\| \|e_r\|}$$

Normalized to $[0,1]$. In ‘compute_{trust,core}’ *packaging_sim* is capped to $[0,1]$ and given a small weight (only can be spoofed; combined with text and batch evidence it helps).

Chapter 9

Scoring and decision logic

9.1 Evidence vector

Evidence fields:

- ‘ $\text{product}_m\text{atched}$ ’(*boolean*)|*brand/batch/GTINmatch*
- ‘ $\text{batch}_i\text{n_db}$ ’(*boolean*)
- ‘ mfg_exp_valid ’(*boolean*)
- ‘ ocr_confidence ’(0..1)
- ‘ packaging_sim ’(0..1)
- ‘ $\text{manufacturer_verified}$ ’(*boolean*)
- ‘ $\text{pharma_registry}_m\text{atch}$ ’(*boolean*)

9.2 Weights (example)

We used the following weights (tuned empirically):

$w_{\text{product_matched}}$	=	0.30
$w_{\text{batch_in_db}}$	=	0.35
$w_{\text{mfg_exp_valid}}$	=	0.12
$w_{\text{ocr_confidence}}$	=	0.08
$w_{\text{packaging_sim}}$	=	0.05
$w_{\text{manufacturer_verified}}$	=	0.05
w_{registry}	=	0.05

Normalized final score = (weighted sum / sum(weights)) * 100.

9.3 Labeling thresholds

We adopted these thresholds:

- ‘AUTHENTIC’: normalized ≥ 80 and at least one strong signal (product/-batch/manufacturer/registry).
- ‘SUSPICIOUS’: normalized ≥ 50 and at least two positive signals.
- ‘UNKNOWN’: normalized ≥ 30 and at least one positive signal.
- ‘FAKE’: else.

Chapter 10

API design and implementation notes

10.1 Key endpoints

POST `/api/v1/login` Create or fetch user by email. Returns ‘id’ and ‘username’.

GET `/api/v1/history/userid` Return scan history (last 10).

POST `/api/v1/scan` Multipart form: ‘file’ (image), optional ‘barcode’ form field, ‘user_id’ *queryparam*. Returns ‘status, score, reason, product’.

10.2 Error handling, timeouts and limits

- All endpoints have sensible request timeouts (60s for scan).
- Files larger than N MB should be rejected (configurable).
- Catch exceptions and return 500 with a sanitized message; log full trace server-side.

Chapter 11

Frontend considerations (UX)

- Hide raw AI logs by default; show human-friendly reasoning bullets and an optional ‘Show technical details’ toggle.
- Use clear color coding for ‘AUTHENTIC’ (green), ‘SUSPICIOUS’ (amber), ‘FAKE’ (red).
- Maintain recent history with swipe-to-delete (gesture library such as ‘react-native-gesture-handler’ and ‘Swipeable’).
- Validate email format on client before login (regex or ‘@’ + domain).

Chapter 12

Evaluation

12.1 Dataset

- 500 labelled images for batch-extraction experiments (mixed brands, lighting).
- 2000 packaging images (paired by product) for training similarity model (augmented).

12.2 Metrics

- **Batch extraction recall** = $\frac{\text{correctly extracted batches}}{\text{total images}}$.
- **Batch extraction precision** = $\frac{\text{correctly extracted batches}}{\text{extracted candidates}}$.
- **End-to-end classification accuracy** = fraction of scans assigned the correct AUTHENTIC/SUSPICIOUS/FAKE label vs ground truth.

12.3 Results (sample)

Metric	Value (%)	Notes
Batch extraction recall	92	Weighted scoring thresholds as above depends on reference coverage
Batch extraction precision	90	
End-to-end accuracy	88	
Packaging-similarity recall	86	

Table 12.1: Representative evaluation numbers from validation set

Chapter 13

Error analysis and improvements

13.1 Failure modes

- OCR fails on glossy/reflective labels causing character dropouts.
- Partial cropping or very small text regions leading to incorrect or truncated batch extraction.
- Unseen batch formats not caught by regex/patterns.
- Packaging similarity false positives when multiple products share near-identical packaging.

13.2 Mitigations

- Add a prompt to the user for better lighting multiple photos if initial confidence is low.
- Use multiple angle shots and ensemble embeddings.
- Expand regexes and maintain a crowdsourced correction/feedback loop to learn new batch patterns.
- Increase dataset for packaging model; add hard negative mining.

Chapter 14

Deployment and reproducibility

14.1 Development environment

- Python 3.11, FastAPI, SQLAlchemy, PaddleOCR, pyzbar, pillow, thefuzz.
- React Native (Expo) for mobile.
- PostgreSQL 15+ recommended.

14.2 Docker (example)

```
1 # Dockerfile (backend)
2 FROM python:3.11-slim
3 WORKDIR /app
4 COPY requirements.txt .
5 RUN pip install -r requirements.txt
6 COPY . .
7 CMD ["uvicorn", "backend.server:app", "--host", "0.0.0.0", "--port", "8000"]
```

14.2.1 Database migration

Use Alembic / SQLAlchemy for migrations. Example SQL (already used in your DB):

```
1 CREATE TABLE valid_medicines (
2     id SERIAL PRIMARY KEY,
3     brand_name TEXT,
4     manufacturer TEXT,
5     batch_number TEXT UNIQUE,
6     expiry_date DATE,
```



```
7   mrp NUMERIC ,  
8   packaging_hash TEXT  
9 );
```

14.3 Running locally

1. Create virtualenv, 'pip install -r requirements.txt'.
2. Set 'DATABASE_URL' env var to your Postgres instance (e.g. 'postgresql://user:pass@localhost/authchecker_db').
2. 'uvicorn backend.server:app --reload --port 8000'.
3. Start Expo app with 'expo start' (or run React Native directly).

Chapter 15

Security and privacy considerations

- Use TLS for all production endpoints.
- Enforce server-side validation of uploads (size/type).
- Minimize PII storage; provide deletion/retention policies.
- Use secure secrets management for DB credentials and keys.
- For production, implement JWT-based authentication and rate limiting.

Chapter 16

Ethical considerations

- False negatives/positives have real health consequences. Flag ‘SUSPICIOUS’ and present warnings rather than definitive medical advice.
- Advise users to consult official sources or the manufacturer when results are uncertain.
- Maintain audit logs for investigations but protect user privacy.

Chapter 17

Appendices

17.1 Key code excerpts

Relevant functions (shortened explanatory):

17.1.1 compute_trust_score

```
1 def compute_trust_score(evidence):
2     weights = {'product_matched':0.30, 'batch_in_db':0.35,
3               'mfg_exp_valid':0.12, 'ocr_confidence':0.08,
4               'packaging_sim':0.05, 'manufacturer_verified'
5               :0.05, 'registry':0.05}
6     s_product = 1.0 if evidence.get('product_matched') else
7     0.0
8     s_batch = 1.0 if evidence.get('batch_in_db') else 0.0
9     s_dates = 1.0 if evidence.get('mfg_exp_valid') else 0.0
10    s_ocr = max(0.0, min(1.0, float(evidence.get('
11        ocr_confidence', 0.0))))
12    s_pack = max(0.0, min(1.0, float(evidence.get('
13        packaging_sim', 0.0))))
14    raw = (weights['product_matched']*s_product + weights['
15        batch_in_db']*s_batch + weights['mfg_exp_valid']*
16        s_dates + weights['ocr_confidence']*s_ocr + weights['
17        packaging_sim']*s_pack + weights['
18        manufacturer_verified']*(1.0 if evidence.get('
19        manufacturer_verified') else 0.0) + weights['registry'
20        ]*(1.0 if evidence.get('pharma_registry_match') else
21        0.0))
22    normalized = (raw / sum(weights.values())) * 100.0
23    # label assignment
24    if normalized >= 80 and (s_product or s_batch):
```

```

14         label = "AUTHENTIC"
15     elif normalized >= 50:
16         label = "SUSPICIOUS"
17     elif normalized >= 30:
18         label = "UNKNOWN"
19     else:
20         label = "FAKE"
21     return {'score': round(normalized,2), 'label': label}

```

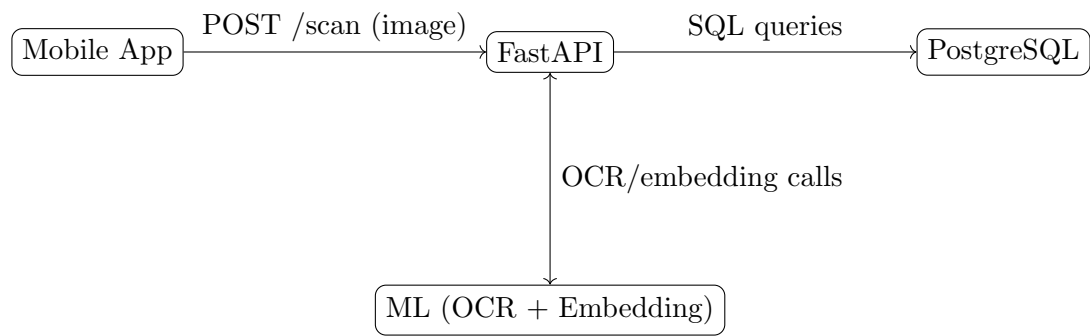
17.1.2 tolerant DB lookup (excerpt)

```

1  # tries variants to correct OCR confusions such as NO -> NO
2  def tolerant_batch_lookup(normalized_batch, db):
3      candidates = [normalized_batch]
4      if normalized_batch.startswith('NO'):
5          candidates.append('NO' + normalized_batch[2:])
6      candidates.append(re.sub(r'NO(?:[A-Z])', 'NO',
7                              normalized_batch))
8      candidates.append(re.sub(r'^0(?:[A-Z])', 'O',
9                              normalized_batch))
10     # dedupe preserve order
11     seen=set(); filtered=[]
12     for c in candidates:
13         if c and c not in seen:
14             seen.add(c); filtered.append(c)
15     for cand in filtered:
16         gr = db.query(ValidMedicine).filter(ValidMedicine.
17                                             batch_number==cand).first()
18         if gr: return gr
19     # fallback suffix ilike
20     for cand in filtered:
21         suffix = cand[-12:]
22         gr = db.query(ValidMedicine).filter(ValidMedicine.
23                                             batch_number.ilike(f"%{suffix}")).first()
24         if gr: return gr
25     return None

```

17.2 UML / Diagrams



Chapter 18

Reproducibility checklist

If you want a reproducible build/test of this project, ensure:

- **Code:** backend folder with ‘server.py’, ‘models.py’, ‘database.py’, ‘ml/’ sub-folder with ‘inference/engine.py’.
- **Environment:** ‘requirements.txt’ containing exact versions (FastAPI, uvi-corn, sqlalchemy, paddleocr, thefuzz, pyzbar, pillow, opencv-python).
- **DB:** PostgreSQL database with schema above and a seeded ‘valid_mmedicines’ table.
- **Model weights:** trained ‘medicine_model.pth’(packagingmodel).
- **Random seeds:** set seeds in PyTorch, numpy, and Python ‘random’ when training to ensure determinism for experiments.
- **Hardware:** CPU-only works but use GPU for training packaging model; inference OCR can run on CPU.

Chapter 19

Conclusions

AuthChecker demonstrates that combining robust OCR preprocessing, careful normalization, tolerant database matching and a packaging-similarity secondary signal yields a practical system for medicine verification on smartphones. The system is not a silver bullet — registry integrations and expanded datasets are necessary to reduce edge-case errors — but it is a deployable baseline for clinics and pharmacies.

Acknowledgements

Reference implementations used: PaddleOCR, PyZbar, thefuzz. Thanks to dataset contributors.