# Image-Based Place Recognition Using Deep Feature Embeddings and Similarity Search

**Mehmet Arda Uçar**
**February 2026**

## 1. Introduction

Place recognition is a fundamental problem in computer vision where a query image's location is determined by comparing it against a reference gallery. This project presents an end-to-end system that relies exclusively on visual features extracted from image pixels, with no metadata or auxiliary information. The system is modular, scalable, and academically rigorous, incorporating proper evaluation metrics (Recall@K, mAP), duplicate detection, and open-set recognition handling.

## 2. System Design

The system follows a standard retrieval pipeline: preprocessing → feature extraction → similarity search → evaluation.Preprocessing: The dataset is scanned and split into gallery (67 images) and query (34 images) sets. MD5 hashing detects exact duplicates, while dHash with Hamming distance identifies near-duplicates. No critical cross-split leakage was detected (1 intra-gallery duplicate, 1 same-place near-duplicate identified).Feature Extraction: A pre-trained ResNet50 (ImageNet) backbone extracts 2048-dimensional embeddings via global average pooling. All embeddings are L2-normalized for cosine similarity computation.Retrieval: Two backends are implemented: NumPy brute-force (baseline) and FAISS IndexFlatIP (scalable). For 67 gallery images, brute-force is sufficient (~0.03–0.10ms/query).Evaluation: Recall@K and mAP are computed with multi-positive query support, and failure analysis identifies confusion pairs.

## 2.1 Technology Stack:

Python 3.11 with type safety; PyTorch for model inference (no training); Pillow for image I/O and preprocessing with torchvision transforms; NumPy for similarity computation and metric calculation; FAISS (optional) for approximate nearest neighbor search; imagehash for perceptual duplicate detection; pytest for unit testing framework; black/ruff/mypy for code quality enforcement.

## 3. Experimental Setup

Dataset: 67 gallery images, 34 query images across 6 landmark categories (Grand Canyon, Eiffel Tower, Palace of Westminster, Buckingham Palace, Galata Tower, Anıtkabir). Gallery expanded from 37 to 67 images for better balance, though some locations remain underrepresented (3 locations have 0 gallery images).Hardware: CPU inference (Intel processor), batch size 8, 4 parallel workers for image loading.Evaluation: K ∈ {1, 5, 10, 20}, cosine similarity ranking, L2-normalized 2048-dim embeddings.

# 4. Results

## 4.1 Baseline Performance (NumPy, No Filtering)

| Metric | Value | Correct Queries |
|---|---|---|
| Recall@1 | 85.29% | 29/34 |
| Recall@5 | 91.18% | 31/34 |
| Recall@10 | 91.18% | 31/34 |
| Recall@20 | 94.12% | 32/34 |
| mAP | 74.09% | – |
| Search Time | ~0.03–0.10 ms/query | – |

Expanding the gallery from 37 to 67 images improved Recall@1 from 65.52% to 85.29% (+19.77%). This confirms that dataset balance has a stronger impact on performance than architectural changes.

## 4.2 Open-Set Confidence Filtering (Max-Similarity)

| Threshold | UNKNOWN Rate | Recall@1 | mAP | FP Removed | TP Removed |
|---|---|---|---|---|---|
| 0.50 | 0% | 85.29% | 74.09% | 0 | 0 |
| 0.75 | 5.88% | 82.35% | 73.42% | 1 | 1 |
| 0.80 | 23.53% | 73.53% | 71.50% | 4 | 4 |

A threshold of 0.75 provides a reasonable precision–recall trade-off with limited performance loss. A threshold of 0.80 is overly conservative, removing a significant number of true positives. Since most predictions already have high similarity scores (>0.75), aggressive filtering offers limited benefit.

## 4.3 FAISS vs NumPy

| Backend | Recall@1 | Recall@5 | mAP | Search Time |
|---|---|---|---|---|
| NumPy | 85.29% | 91.18% | 74.09% | ~0.03–0.10 ms |
| FAISS | 85.29% | 91.18% | 74.09% | ~0.47 ms |

FAISS (IndexFlatIP) produces identical accuracy to NumPy, as both perform exact similarity search. For small datasets (<1,000 images), NumPy is faster due to lower overhead. FAISS becomes advantageous in large-scale scenarios (10k+ images) when approximate indexing methods are used.

# 5. Failure Analysis

Failures: 3/34 queries fail at top-5 (8.82% failure rate). All failures involve Eiffel Tower queries confused with Anıtkabir or Palace of Westminster.Confusion Pairs:

- Eiffel Tower → Palace of Westminster (3x)
- Eiffel Tower → Anıtkabir (1x)
- Buckingham Palace → Eiffel Tower (1x)

Root Causes: Architectural similarity between European landmarks, limited gallery diversity (Palace of Westminster has 0 gallery images), generic ImageNet features

lacking fine-grained discriminative power, small similarity margins (mean=0.0125) indicating model uncertainty.

## 6. Key Observations

Dataset balance has stronger impact than architectural changes: +19.77% Recall@1 from gallery expansion. Most failures are semantically plausible confusions (similar structures) rather than random errors. Small similarity margins ($\approx$0.0125) indicate model uncertainty, suggesting moderate representation improvements could yield significant gains. Open-set thresholding demonstrates precision-recall trade-off, but most predictions already have high confidence, limiting filtering benefits. FAISS provides identical accuracy to NumPy for small datasets, validating correct implementation and future scalability.

## 7. Proposed Improvements

Short-term (High Impact): Hard negative mining targeting Eiffel Tower/Palace/Anıtkabir confusion (expected +5–8% Recall@1), balanced gallery expansion for underrepresented locations, adaptive threshold calibration.Medium-term: Fine-tuning with landmark-specific metric learning (Triplet Loss, ArcFace) for +10–15% Recall@1, advanced pooling (GeM, NetVLAD) for spatial sensitivity, ensemble fusion (ResNet50 + EfficientNet) for +3–5% mAP.Long-term: Supervised metric learning for viewpoint-invariant features (+20–25% Recall@1 to reach 95%+), multi-scale feature integration (global + local descriptors), approximate nearest neighbor optimization for 10k+ galleries.

## 8. Conclusion

This system demonstrates that proper engineering and dataset balancing achieve strong performance (85.29% Recall@1, 74.09% mAP) using standard pre-trained models. The implementation satisfies academic requirements through correct metric computation, comprehensive testing (34 unit tests), duplicate detection preventing data leakage, and modular architecture. Primary limitations stem from gallery imbalance and generic ImageNet features. With domain-specific fine-tuning and balanced gallery expansion, the system can reach 95%+ Recall@1.