

# Rasterized Image Databases for Image Compression

Zoya Bylinskii\*, Maria Shugrina\*, Andrew Spielberg\*, Wei Zhao\*

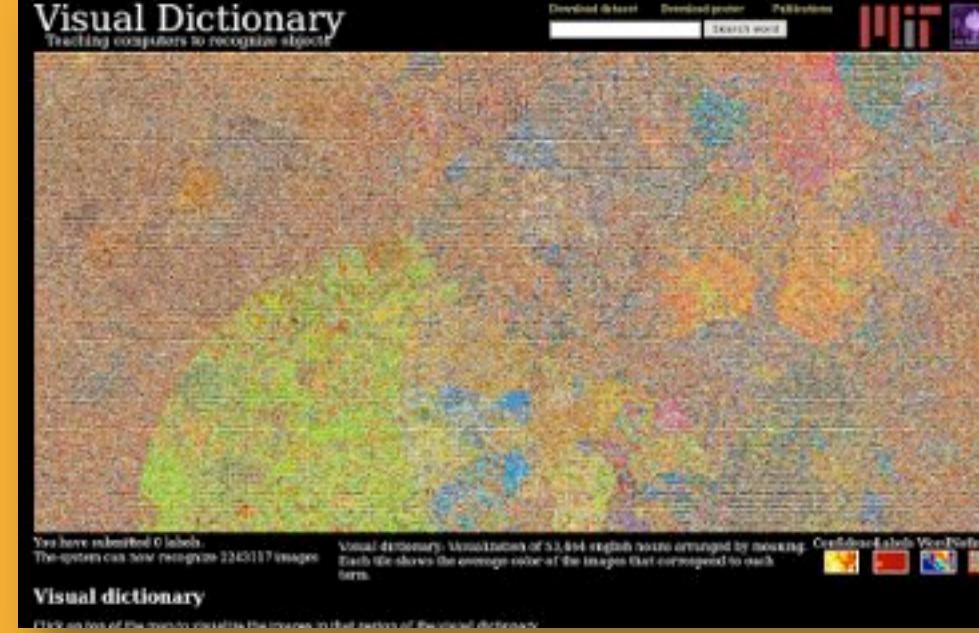


bigdata@CSAIL  
MIT BIG DATA INITIATIVE

6.830  
Database Systems  
Final Project

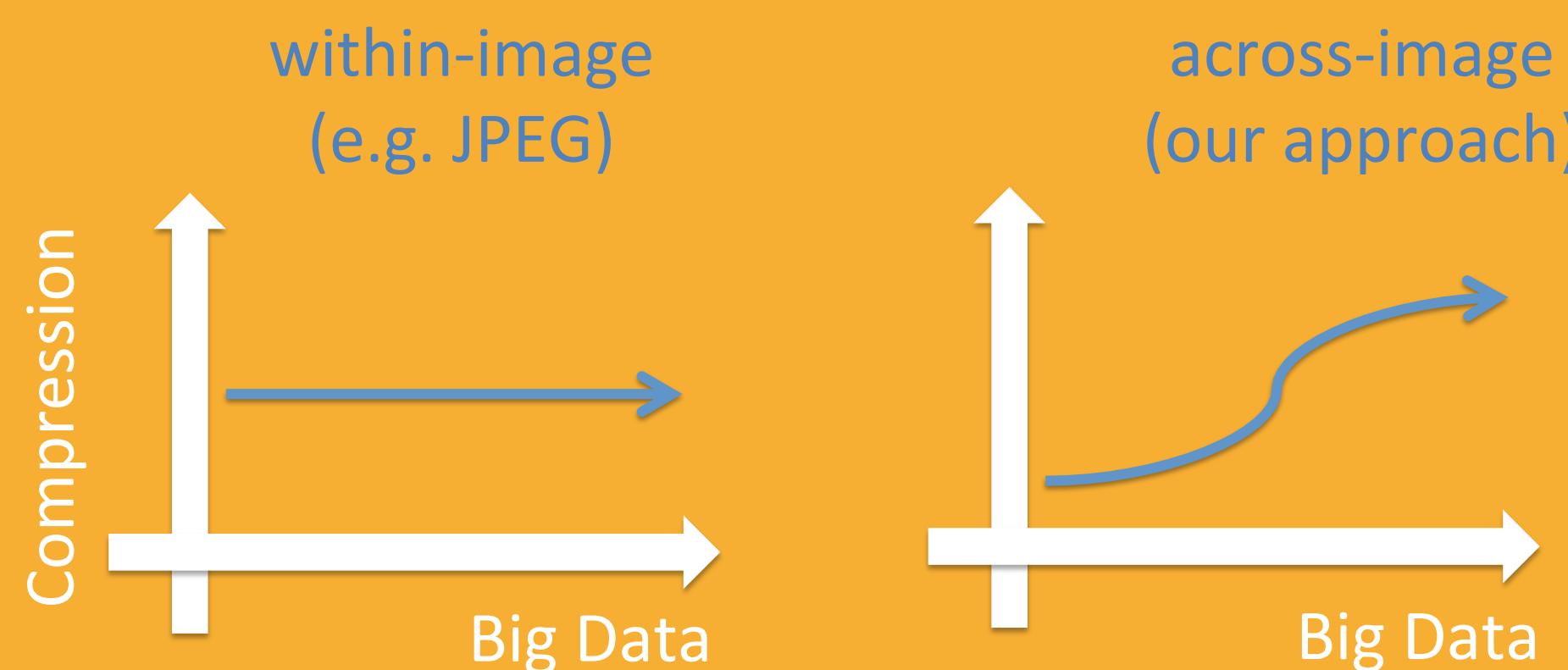
## MOTIVATION

- More than 1.8 billion images uploaded to the internet every day
- Redundancy in large image collections can lead to more efficient storage

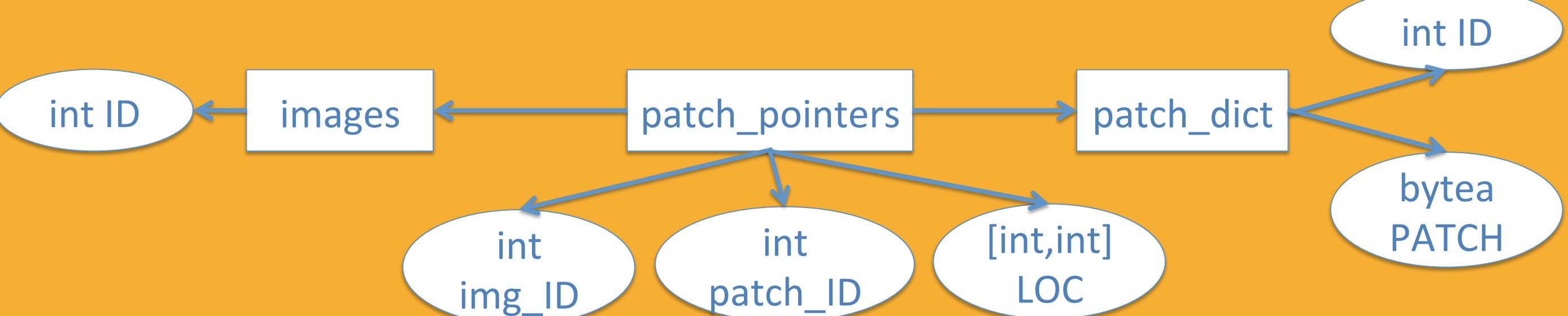


## PATCH-BASED IMAGE CODING

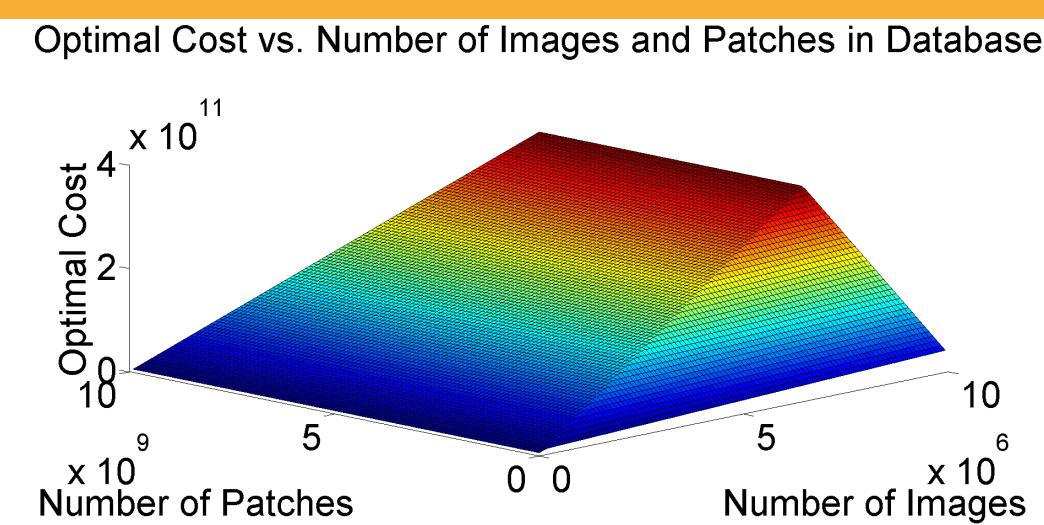
- Lossy compression
- Compression savings INCREASE as the database size INCREASES
- Rasterization is a simple operation – easily deconstruct and reconstruct images using patches



## DATABASE SCHEMA



## COST CALCULATIONS



$$\begin{aligned} & \text{minimize}_{\text{patch\_dict}, M} c(k, d, m, n) \\ & \text{subject to } S(P_j, M(P_j)) \leq T, j = 1, \dots, k \left(\frac{m}{n}\right)^2 \end{aligned}$$

Cost without patch-based scheme:

$$c'(k, m) = 3km^2$$

Cost with patch-based scheme:

$$c(k, d, m, n) = 8k \left(\frac{m}{n}\right)^2 + 3dn^2$$

Constraint for cost savings:

$$d < \frac{m^2(3n^2 - 8)k}{3n^4}$$

## METHOD OVERVIEW

### SCHEMA:

```
images(id int PRIMARY KEY);
patch_dict(id int PRIMARY KEY, patch bytea);
patch_pointers(img_id int REFERENCES images(id),
               patch_id int REFERENCES patches(id),
               x int, y int);
```

### ALGORITHM:

**Algorithm 1** Insert Image  $I$  into database

- $Patches \leftarrow \text{Patchify}(I, n)$
- for**  $P_j$  in  $Patches$  **do**
- $\quad SimPat \leftarrow \text{FindLikelySimilarPatches}(P_j, patch\_dict)$
- $\quad P_{NN} \leftarrow \text{argmin}_{P_i \in SimPat} \{S(P_i, P_j)\}$
- if**  $S(P_{NN}, P_j) > T$  **then**
- insert**  $P_j$  into patches

**Algorithm 2** Modification of alg. 1, with ANN

- $Patches \leftarrow \text{Patchify}(I, n)$
- for**  $P_j$  in  $Patches$  **do**
- $\quad SimPat \leftarrow \text{FindLikelySimilarPatches}(P_j, patch\_dict)$
- $\quad P_{ANN} \leftarrow M(P_j)$
- if**  $S(P_{ANN}, P_j) > T$  **then**
- insert**  $P_j$  into patches

**Algorithm 3** Optimization of alg. 2, with ANN

- $Patches \leftarrow \text{Patchify}(I, n)$
- $UniquePatches \leftarrow \text{FindUniquePatches}(I)$
- $PatchesToStoreInDB \leftarrow []$
- for**  $P_j$  in  $UniquePatches$  **do**
- $\quad SimPat \leftarrow \text{FindLikelySimilarPatches}(P_j, patch\_dict)$
- $\quad P_{ANN} \leftarrow M(P_j)$
- if**  $S(P_{ANN}, P_j) > T$  **then**
- PatchesToStoreInDB.Add**( $P_j$ )
- BatchInsert**( $PatchesToStoreInDB$ )

### IMPLEMENTATION:

PostgreSQL database with Java API.  
Image loading, patching, searching, hashing, Near-Neighbor search, reconstruction implemented in Java.

[https://github.com/shumash/db\\_project](https://github.com/shumash/db_project)

Distance per color channel in CIE (LUV) color space:

$$S(P_i, P_j, u) = \frac{\|P_i(u) - P_j(u)\|^2}{n^2}$$

Thresholds  $T_1, T_2, T_3$  per channel.

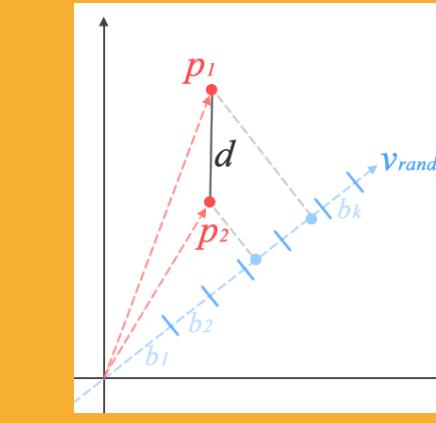
Image size  $m \times n$ , patch size  $n \times n$ .

### KEY OPTIMIZATIONS:

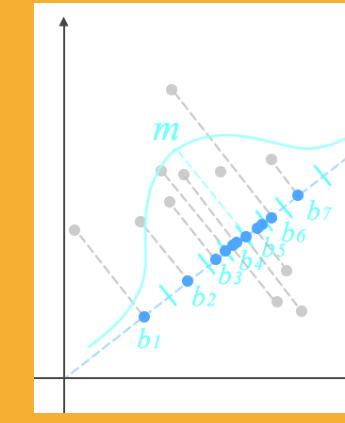
- Minimize database queries (batch query per image)
- Near-neighbor search
- Patch BufferPool for patches (LRU)
- Minimize re-computation of image vectors
- Exploit image self-similarity

## NEAR NEIGHBOR SEARCH

LSH: AND and OR hashes for high precision/recall.  
Can we find one AND hash for reasonable  $P_1$ , low  $P_2$ ?



Random-projection hash  
(we use AND of 10 in 32-bit int):  
poor  $P_1, P_2$

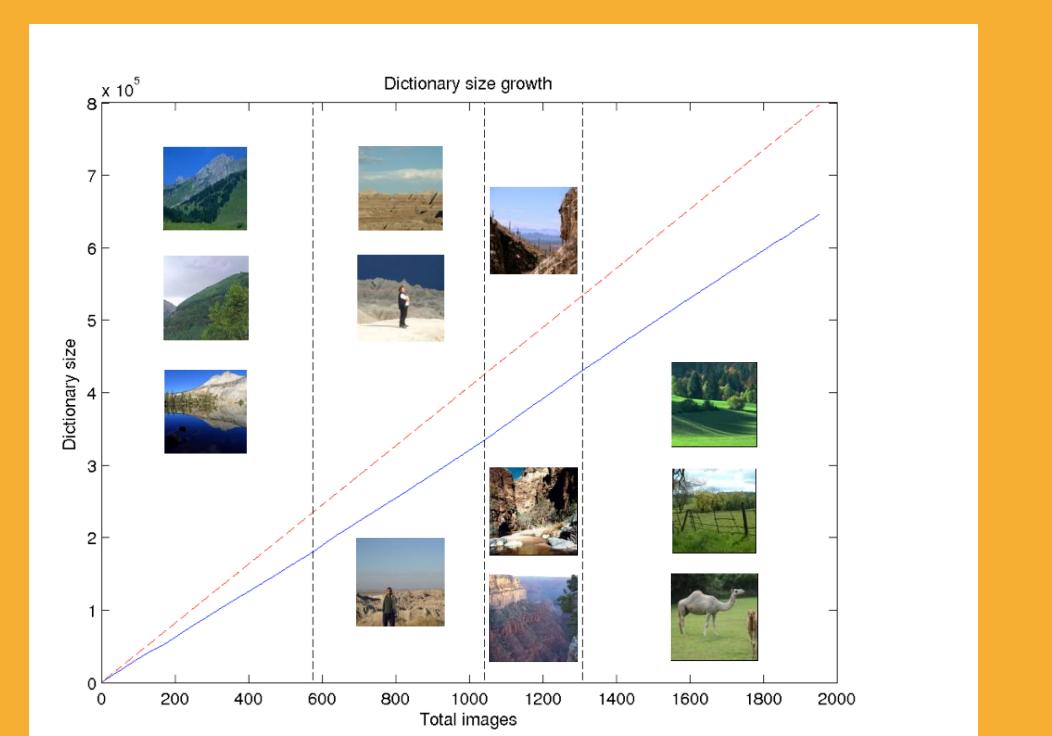
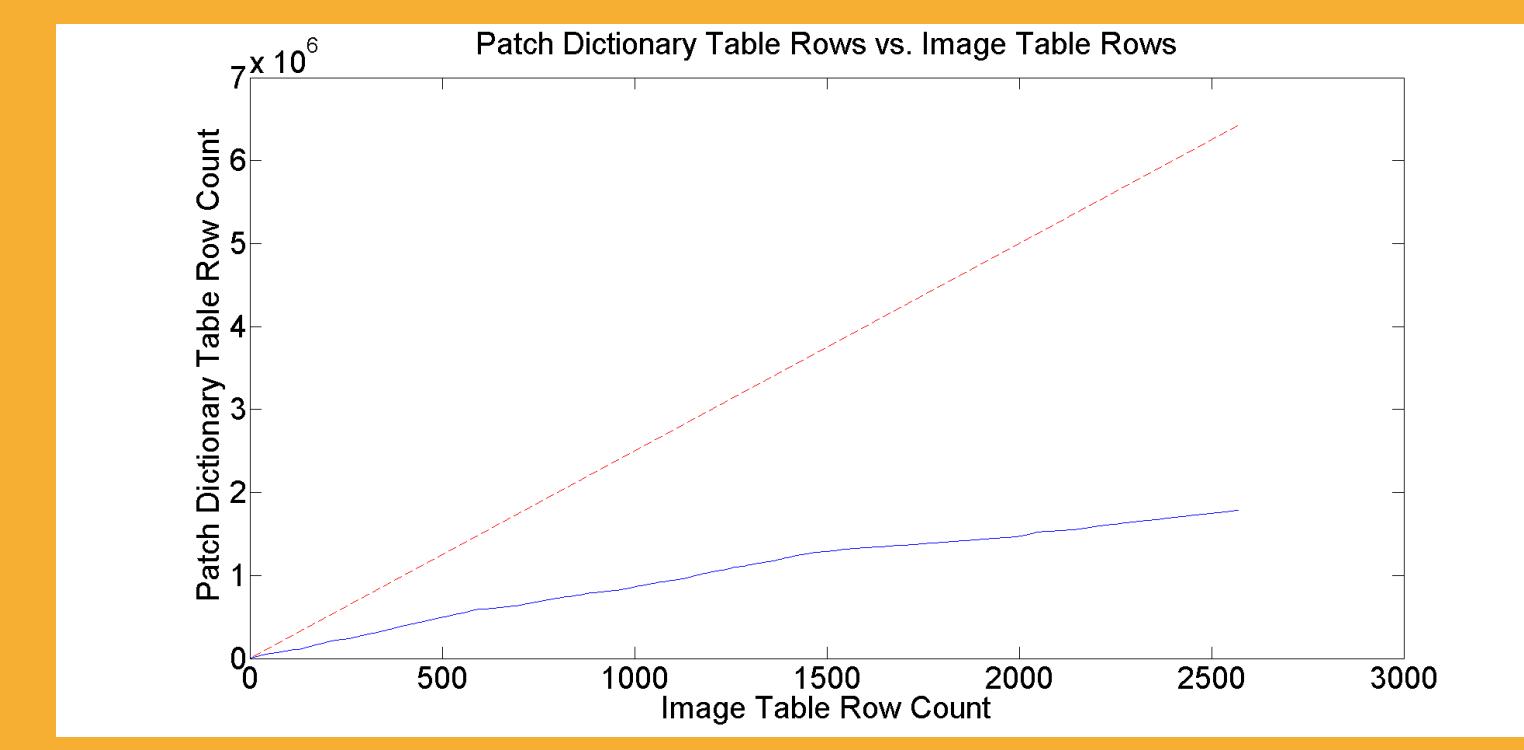


Using PCA directions as projection vectors to improve  $P_1, P_2$

$$\begin{aligned} P[b(P_i) = b(P_j)] &= S(P_i, P_j) < T \rangle > P_1 \\ P[b(P_i) = b(P_j)] &= S(P_i, P_j) > cT \rangle < P_2 \end{aligned}$$

**ADDITIONAL INSIGHT:**  
Natural images contain many nearly uniform patches, which are easy to index and search:  
`uniform_patch_dict`

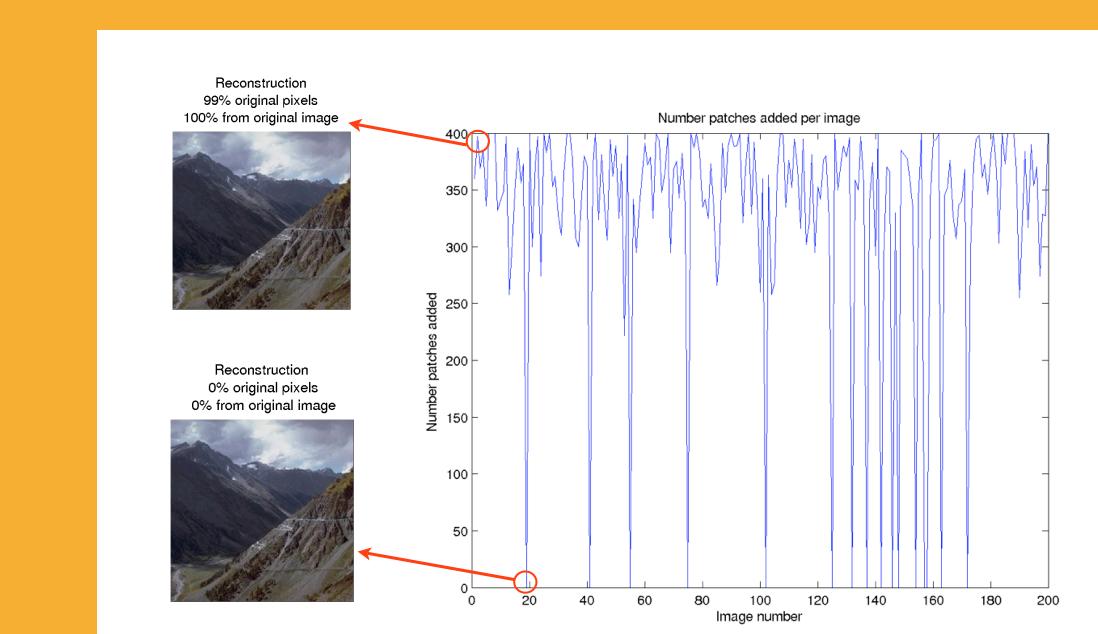
## GENERALIZABILITY



Generalizability across image types:

## APPLICATIONS

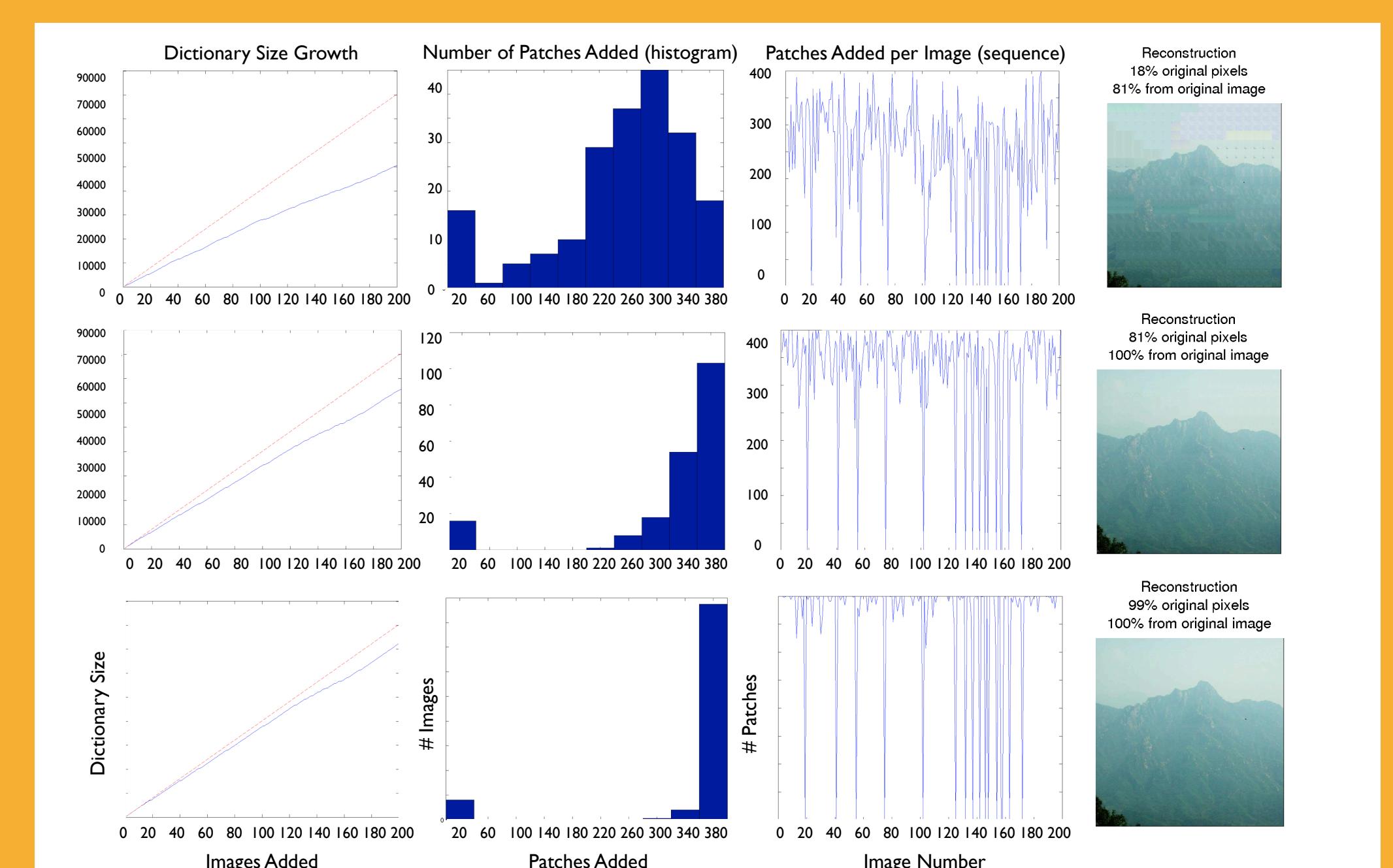
- duplicate detection
- similar image retrieval



## EXTENSIONS

- context-aware patch sampling
- task-aware threshold selection
- patch hierarchies

Weighting storage savings and reconstruction speed against compression quality:



The importance of the right patch sampling technique:

