# ImageProof: Enabling Authentication for Large-Scale Image Retrieval

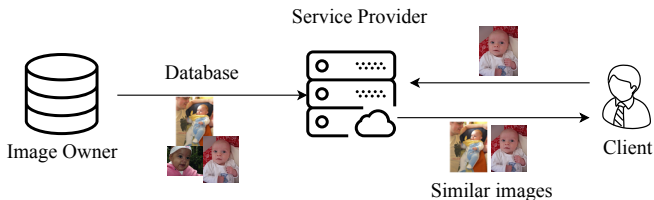Shangwei Guo[1] Jianliang Xu[1] Ce Zhang[1] Cheng Xu[1] Tao Xiang[2]

[1]Department of Computer Science, Hong Kong Baptist University

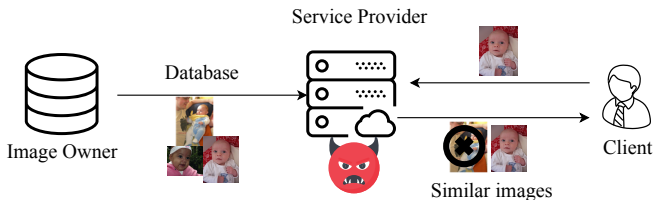[2]College of Computer Science, Chongqing University

{csswguo,xujl,cezhang,chengxu}@comp.hkbu.edu.hk, txiang@cqu.edu.cn

ICDE 2019

# Background

- Content-based image retrieval (CBIR) has been widely used in business
- Data-as-a-Service (DaaS) enables companies to build and then outsource image retrieval systems to cloud platforms
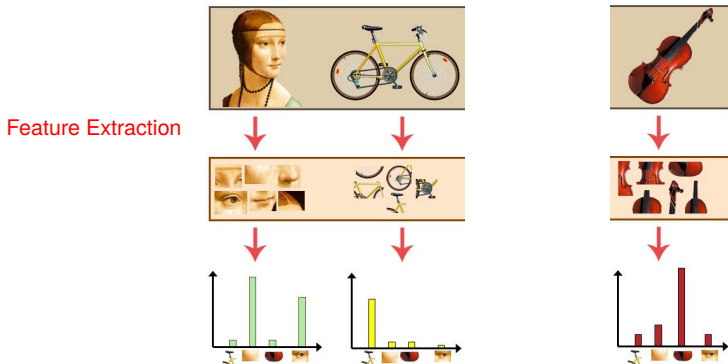
# Background

- Content-based image retrieval (CBIR) has been widely used in business
- Data-as-a-Service (DaaS) enables companies to build and then outsource image retrieval systems to cloud platforms



- **Security Threat**:
  - Query result integrity not guaranteed due to software/hardware malfunctions, hack attacks
  - Examples
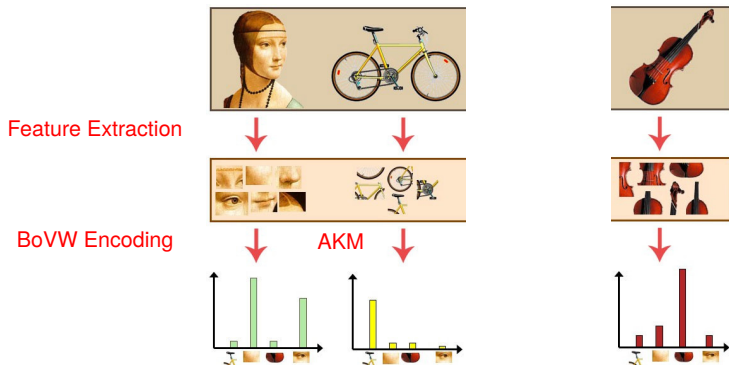    - Product image search
    - Medical image search

# SIFT-Based Image Retrieval

- Detect and extract local features using scale invariant feature transform (SIFT) and its variants



Feature Extraction

# SIFT-Based Image Retrieval

- Detect and extract local features using scale invariant feature transform (SIFT) and its variants
- **Twe Steps**
  - Bag-of-visual-words (BoVW) encoding
    - Approximate $k$-means (AKM) using randomized k-d trees



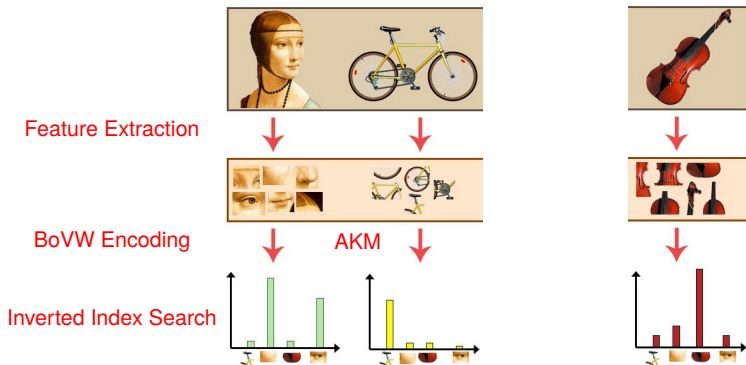Feature Extraction

BoVW Encoding — AKM

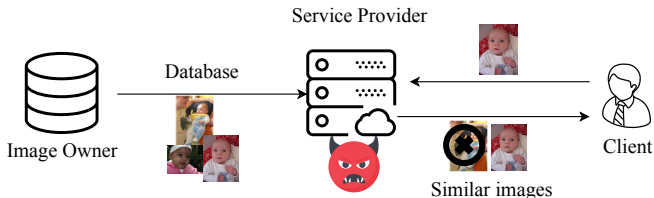# SIFT-Based Image Retrieval

- Detect and extract local features using scale invariant feature transform (SIFT) and its variants
- **Twe Steps**
  - Bag-of-visual-words (BoVW) encoding
    - Approximate $k$-means (AKM) using randomized k-d trees
  - Inverted index search: search similar images with impact-ordered inverted index



Feature Extraction

BoVW Encoding     AKM

Inverted Index Search

# Problem Model

- Malicious threat model
- The service provider (SP) could return **incorrect** results (e.g., faked or low-ranked images)
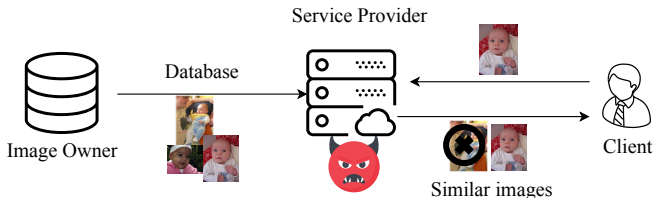
# Problem Model

- Malicious threat model

- The service provider (SP) could return **incorrect** results (e.g., faked or low-ranked images)



- Query authentication for SIFT-based image retrieval and top-$k$ query

# Problem Model

- Malicious threat model

- The service provider (SP) could return **incorrect** results (e.g., faked or low-ranked images)



Service Provider

Database

Image Owner

Client

Similar images

- Query authentication for SIFT-based image retrieval and top-$k$ query

- **Challenges**
    - Designing a query authentication scheme for a large and complex retrieval system is a big challenge in itself
    - The client usually has only limited storage, communication, and computation resources

# Problem Model

Service Provider

Database and ADS

Image Owner

Similar images & VO

Client

- **Our Solution**:
  - Taking the advantage of the authenticated data structures (ADSs), the SP returns a verification object (VO) to prove
    - Soundness: The results must be the images which have not been tampered with
    - Completeness: The results include the $k$ most similar images

## Our Contributions

- Propose an efficient authentication scheme, ImageProof, for SIFT-based image retrieval with large or medium-sized codebooks

# Our Contributions

- Propose an efficient authentication scheme, ImageProof, for SIFT-based image retrieval with large or medium-sized codebooks
- Two novel ADS components:
  - Merkle randomized k-d tree
  - Merkle inverted index with cuckoo filters

# Our Contributions

- Propose an efficient authentication scheme, ImageProof, for SIFT-based image retrieval with large or medium-sized codebooks
- Two novel ADS components:
  - Merkle randomized k-d tree
  - Merkle inverted index with cuckoo filters
- Develop several optimization techniques to further reduce the costs of both the SP and the client

- **Merkle Hash Tree**
  - An authenticated binary tree, enabling users to verify individual data objects without retrieving the entire database



Figure 1: An example of a Merkle hash tree.

# Preliminaries

- **Merkle Hash Tree**
  - An authenticated binary tree, enabling users to verify individual data objects without retrieving the entire database
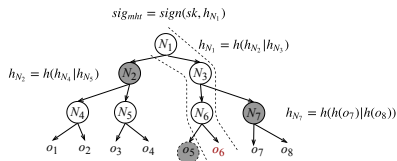- **Cuckoo Filter**
  - An efficient data structure for approximate set membership tests
    - Two hash values per item
    - Support delete operation



Figure 1: An example of a Merkle hash tree.



Figure 2: A cuckoo filter, two hash values per item.

# Scheme Overview

- Ensure the integrity of query processing for each step



Authenticated BoVW Encoding

Authenticated Inverted Index Search

# Scheme Overview

- Ensure the integrity of query processing for each step
- Two novel ADS components:
  - Merkle randomized k-d tree
  - Merkle inverted index with cuckoo filters



Merkle Randomized k-d Tree

Merkle Inverted Index

# Merkle Randomlized k-d Tree (MRKD-tree)

- **ADS**
  - Internal nodes and leaf nodes



Figure 3: An example of the MRKD-tree and VO generation for query $q_1$, $q_2$.

# Merkle Randomlized k-d Tree (MRKD-tree)

- **ADS**
  - Internal nodes and leaf nodes



Figure 3: An example of the MRKD-tree and VO generation for query $q_1$, $q_2$.

- **Authenticated Query Processing**
  - Given a set of feature vectors, calculate the BoVW vector
  - Generate a single verification object (VO) for all feature vectors by maximizing the use of shared tree nodes

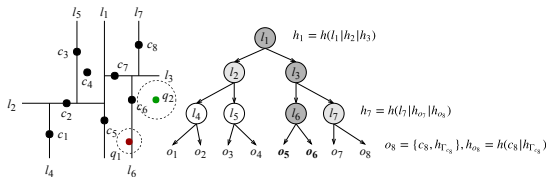# Merkle Inverted Index With Cuckoo Filters

- **ADS**
  - Each Merkle inverted list $\Gamma_{c_i}$ consists of five components, i.e., the associated cluster $c_i$, the digest $h(\Theta_{c_i})$, the cluster weight $w_{c_i}$, the cuckoo filter $\Theta_i$ and its posting list

Table 1: An example of the Merkle inverted lists.

| $c_i$ | $h_{\Gamma_{c_i}}$ | $w_{c_i}$ | $\Theta_i$ | | Posting Lists | | | |
|---|---|---|---|---|---|---|---|---|
| $c_5$ | $h(2\sqrt{2}\|h(\Theta_{c_5})\|h_{pos_{5,1}})$ | $2\sqrt{2}$ | $\Theta_{c_5}$ | $\mapsto$ | $\langle 1, 0.34, h_{pos_{5,1}}\rangle$ | $\langle 3, 0.26, h_{pos_{5,2}}\rangle$ | $\langle 4, 0.25, h_{pos_{5,3}}\rangle$ | ... |
| $c_6$ | $h(\sqrt{2}\|h(\Theta_{c_6})\|h_{pos_{6,1}})$ | $\sqrt{2}$ | $\Theta_{c_6}$ | $\mapsto$ | $\langle 5, 0.41, h_{pos_{6,1}}\rangle$ | $\langle 8, 0.32, h_{pos_{6,2}}\rangle$ | $\langle 3, 0.28, h_{pos_{6,3}}\rangle$ | ... |

# Merkle Inverted Index With Cuckoo Filters

- **ADS**
  - Each Merkle inverted list $\Gamma_{c_i}$ consists of five components, i.e., the associated cluster $c_i$, the digest $h(\Theta_{c_i})$, the cluster weight $w_{c_i}$, the cuckoo filter $\Theta_i$ and its posting list

Table 1: An example of the Merkle inverted lists.

| $c_i$ | $h_{\Gamma_{c_i}}$ | $w_{c_i}$ | $\Theta_i$ | | Posting Lists | | |
|---|---|---|---|---|---|---|---|
| $c_5$ | $h(2\sqrt{2}|h(\Theta_{c_5})|h_{pos_{5,1}})$ | $2\sqrt{2}$ | $\Theta_{c_5}$ | $\mapsto$ | $\langle 1, 0.34, h_{pos_{5,1}}\rangle$ | $\langle 3, 0.26, h_{pos_{5,2}}\rangle$ | $\langle 4, 0.25, h_{pos_{5,3}}\rangle$ ... |
| $c_6$ | $h(\sqrt{2}|h(\Theta_{c_6})|h_{pos_{6,1}})$ | $\sqrt{2}$ | $\Theta_{c_6}$ | $\mapsto$ | $\langle 5, 0.41, h_{pos_{6,1}}\rangle$ | $\langle 8, 0.32, h_{pos_{6,2}}\rangle$ | $\langle 3, 0.28, h_{pos_{6,3}}\rangle$ ... |

- **Authenticated Query Processing**
  - Find top-$k$ most similar images and generate the VO of inverted index search
  - Ensure the integrity of top-$k$ search with fewer postings with the help of cuckoo filters

# Merkle Inverted Index With Cuckoo Filters

- **Main Idea**
  - Termination conditions:
    1. $s_k^L \geq S^U(Q, I)$, the upper bound of the similarity scores of the images popped, where $s_k^L$ is the lower bound of the k-th similar score
    2. $s_k^L \geq$ the upper bound of the similarity scores of the images not popped

| $c_i$ | $h_{\Gamma_{c_i}}$ | $w_{c_i}$ | $\Theta_i$ | | **Posting Lists** | | | | | |
|-------|--------------------|-----------|------------|---|---|---|---|---|---|---|
| $c_5$ | $h(2\sqrt{2}|h(\Theta_{c_5})|h_{pos_{5,1}})$ | $2\sqrt{2}$ | $\Theta_{c_5}$ | $\mapsto$ | $\langle 1, 0.34, h_{pos_{5,1}} \rangle$ | $\langle 3, 0.26, h_{pos_{5,2}} \rangle$ | $\langle 4, 0.25, h_{pos_{5,3}} \rangle$ | $\langle 10, 0.17, h_{pos_{5,4}} \rangle$ | $\langle 7, 0.11, h_{pos5,5} \rangle$ | ... |
| $c_6$ | $h(\sqrt{2}|h(\Theta_{c_6})|h_{pos_{6,1}})$ | $\sqrt{2}$ | $\Theta_{c_6}$ | $\mapsto$ | $\langle 5, 0.41, h_{pos_{6,1}} \rangle$ | $\langle 8, 0.32, h_{pos_{6,2}} \rangle$ | $\langle 3, 0.28, h_{pos_{6,3}} \rangle$ | $\langle 6, 0.25, h_{pos_{6,4}} \rangle$ | $\langle 4, 0.10, h_{pos_{6,5}} \rangle$ | ... |

# Merkle Inverted Index With Cuckoo Filters

- **Main Idea**
  - Termination conditions:
    1. $s_k^L \geq S^U(Q, I)$, the upper bound of the similarity scores of the images popped, where $s_k^L$ is the lower bound of the k-th similar score
    2. $s_k^L \geq$ the upper bound of the similarity scores of the images not popped

| $c_i$ | $h_{\Gamma_{c_i}}$ | $w_{c_i}$ | $\Theta_i$ | | **Posting Lists** | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| $c_5$ | $h(2\sqrt{2}|h(\Theta_{c_5})|h_{pos_{5,1}})$ | $2\sqrt{2}$ | $\Theta_{c_5}$ | $\mapsto$ | $\langle 1, 0.34, h_{pos_{,1}} \rangle$ | $\langle 3, 0.26, h_{pos_{,2}} \rangle$ | $\langle 4, 0.25, h_{pos_{5,3}} \rangle$ | $\langle 10, 0.17, h_{pos_{5,4}} \rangle$ | $\langle 7, 0.11, h_{pos5,5} \rangle$ | ... |
| $c_6$ | $h(\sqrt{2}|h(\Theta_{c_6})|h_{pos_{6,1}})$ | $\sqrt{2}$ | $\Theta_{c_6}$ | $\mapsto$ | $\langle 5, 0.41, h_{pos_{6,1}} \rangle$ | $\langle 8, 0.32, h_{pos_{6,2}} \rangle$ | $\langle 3, 0.28, h_{pos_{6,3}} \rangle$ | $\langle 6, 0.25, h_{pos_{6,4}} \rangle$ | $\langle 4, 0.10, h_{pos_{6,5}} \rangle$ | ... |

- Estimate the similarity bounds using the cuckoo filters

Table 2: Example: the postings for $S(Q, 5)$.

$$\text{Without cuckoo filter:} \quad \begin{aligned} S^U(Q, 5) &\mapsto \langle 5, 0.41, h_{pos_{6,1}} \rangle, \langle 4, 0.25, h_{pos_{5,3}} \rangle \\ S^L(Q, 5) &\mapsto \langle 5, 0.41, h_{pos_{6,1}} \rangle \end{aligned}$$

# Merkle Inverted Index With Cuckoo Filters

- **Main Idea**
  - Termination conditions:
    1. $s_k^L \geq S^U(Q, I)$, the upper bound of the similarity scores of the images popped, where $s_k^L$ is the lower bound of the k-th similar score
    2. $s_k^L \geq$ the upper bound of the similarity scores of the images not popped

| $c_i$ | $h_{\Gamma_{c_i}}$ | $w_{c_i}$ | $\Theta_i$ | | Posting Lists | | | | |
|-------|-------------------|-----------|------------|---|---|---|---|---|---|
| $c_5$ | $h(2\sqrt{2}|h(\Theta_{c_5})|h_{pos_{5,1}})$ | $2\sqrt{2}$ | $\Theta_{c_5}$ | $\mapsto$ | $\langle 1, 0.34, h_{pos_{5,1}} \rangle$ | $\langle 3, 0.26, h_{pos_{5,2}} \rangle$ | $\langle 4, 0.25, h_{pos_{5,3}} \rangle$ | $\langle 10, 0.17, h_{pos_{5,4}} \rangle$ | $\langle 7, 0.11, h_{pos5,5} \rangle$ ... |
| $c_6$ | $h(\sqrt{2}|h(\Theta_{c_6})|h_{pos_{6,1}})$ | $\sqrt{2}$ | $\Theta_{c_6}$ | $\mapsto$ | $\langle 5, 0.41, h_{pos_{6,1}} \rangle$ | $\langle 8, 0.32, h_{pos_{6,2}} \rangle$ | $\langle 3, 0.28, h_{pos_{6,3}} \rangle$ | $\langle 6, 0.25, h_{pos_{6,4}} \rangle$ | $\langle 4, 0.10, h_{pos_{6,5}} \rangle$ ... |

  - Estimate the similarity bounds using the cuckoo filters

Table 2: Example: the postings for $S(Q, 5)$.

Without cuckoo filter: $\quad S^U(Q, 5) \quad \mapsto \quad \langle 5, 0.41, h_{pos_{6,1}} \rangle, \langle 4, 0.25, h_{pos_{5,3}} \rangle$
$\qquad\qquad\qquad\qquad S^L(Q, 5) \quad \mapsto \quad \langle 5, 0.41, h_{pos_{6,1}} \rangle$
With cuckoo filter: $\qquad S^U(Q, 5) \quad \mapsto \quad \langle 5, 0.41, h_{pos_{6,1}} \rangle$
$\qquad\qquad\qquad\qquad S^L(Q, 5) \quad \mapsto \quad \langle 5, 0.41, h_{pos_{6,1}} \rangle$

# ImageProof

- **ADS Generation**
  - Build Merkle inverted lists $\{\Gamma_{c_i}\}$ and MRKD-trees $\{\mathcal{T}_i\}$
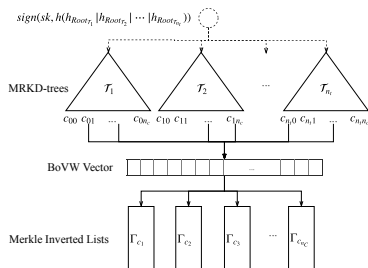


Figure 4: An overview of ADSs for ImageProof.

# ImageProof

- **ADS Generation**
  - Build Merkle inverted lists $\{\Gamma_{c_i}\}$ and MRKD-trees $\{\mathcal{T}_i\}$
- **Authenticated Query Processing**
  - Search the top-$k$ images and generate the VOs for both the BoVW encoding and the inverted index search
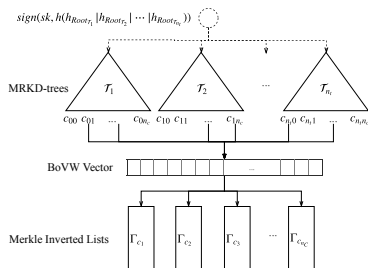  - Send the VOs, together with the top-$k$ results the client



Figure 4: An overview of ADSs for ImageProof.

# ImageProof

- **ADS Generation**
  - Build Merkle inverted lists $\{\Gamma_{c_i}\}$ and MRKD-trees $\{\mathcal{T}_i\}$
- **Authenticated Query Processing**
  - Search the top-$k$ images and generate the VOs for both the BoVW encoding and the inverted index search
  - Send the VOs, together with the top-$k$ results the client
- **Result Verification**
  - Check the integrity of image retrieval
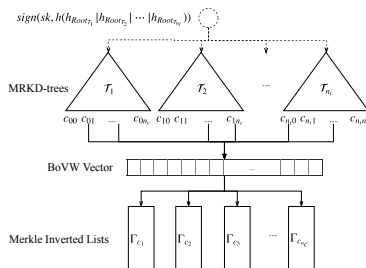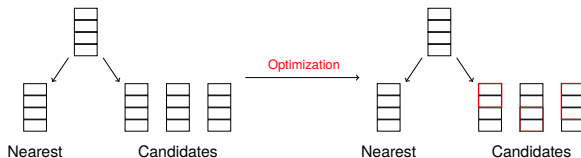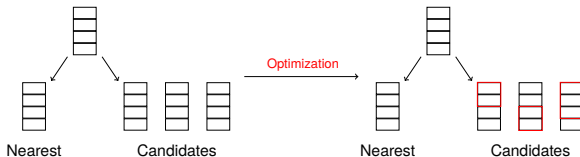  - Verify the integrity of raw image data



$sign(sk, h(h_{Root_{\mathcal{T}_1}}|h_{Root_{\mathcal{T}_2}}|\cdots|h_{Root_{\mathcal{T}_{n_i}}}))$

MRKD-trees $\mathcal{T}_1$ $\mathcal{T}_2$ $\cdots$ $\mathcal{T}_{n_i}$

$c_{00}\ c_{01}\ \cdots\ c_{0n_c}$ $c_{10}\ c_{11}\ \cdots\ c_{1n_c}$ $c_{n_i0}\ c_{n_i1}\ \cdots\ c_{n_in_c}$

BoVW Vector

Merkle Inverted Lists $\Gamma_{c_1}$ $\Gamma_{c_2}$ $\Gamma_{c_3}$ $\cdots$ $\Gamma_{c_C}$

Figure 4: An overview of ADSs for ImageProof.

# Optimization

- Compressing nearest neighbor candidates

- Compressing nearest neighbor candidates



- Frequency-grouped inverted index

- **Experimental Setup**
  - Dataset: MirFlickr1M
  - Algorithms
    - **Baseline**: The scheme that combines the proposed MRKD-trees without sharing nodes and the authenticated inverted index search in PVLDB2008
    - **ImageProof**: The proposed scheme
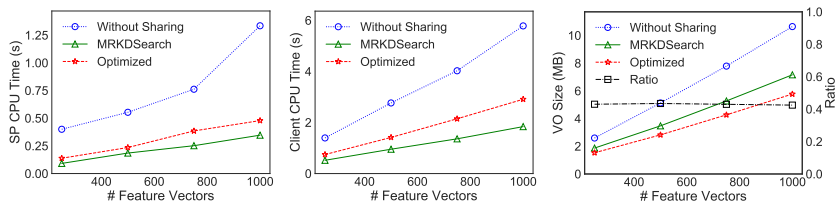    - **Optimized**: The optimized ImageProof

# BoVW Performance

Figure 5: BoVW performance as the number of feature vectors increases.
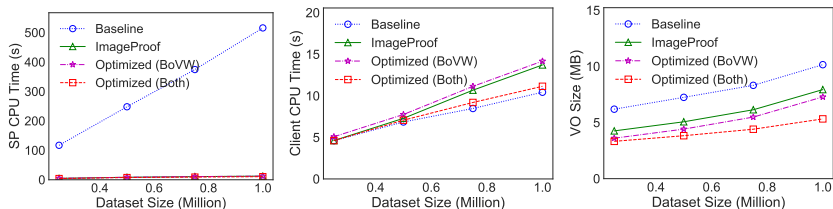
# Overall Performance



Figure 6: Overall performance as dataset size increases.

# Summary

- Focus on the query authentication problem in SIFT-based image retrieval
- Two authenticated data structures (ADSs) for both BoVW encoding and inverted index search
- Extensive experiments on real-world image dataset

# Thanks

# Q&A