

# Deep Feature Embedding for Content-Based Image Retrieval Using Autoencoders and Semantic Hashing

Azhar amer alsoufi

Email: [mti.lec250.azhar@ntu.edu.iq](mailto:mti.lec250.azhar@ntu.edu.iq)

Department of Networks and Computer Software Techniques, Northern Technical University, Mosul,  
Iraq.

Azhar amer alsoufi

Email: mti.lec250.azhar@ntu.edu.iq

Department of Networks and Computer Software Techniques, Northern Technical University, Mosul,  
Iraq.

**Abstract**—The necessity for effective and efficient Content- Based Image Retrieval (CBIR) systems—where retrieval is done based on the visual content of images rather than metadata—has increased due to the exponential expansion of digital image repositories. Conventional CBIR methods that depend on manually created features frequently fail to capture intricate semantic links in image data, which restricts their scalability and retrieval performance. Recent research has investigated deep learning-based representations and hashing methods for increased efficiency and accuracy in order to overcome these issues. Nevertheless, a lot of current methods either need a lot of labeled data or are unable to produce compact representations that are quick to retrieve in large-scale environments. In this study, we offer a deep unsupervised CBIR framework that combines semantic hashing with convolutional autoencoders. A semantic hashing mechanism is used to convert the discriminative, low-dimensional feature embeddings that the autoencoder learns into compact binary hash codes. This makes it possible to use the Hamming distance for an effective similarity search. Comprehensive tests on benchmark datasets like Caltech-101 and CIFAR-10 show that our method provides a scalable and efficient solution for large-scale image retrieval problems by achieving competitive retrieval accuracy while drastically cutting down on retrieval time.

**Index Terms**—Content-Based Image Retrieval, Deep Autoencoder, Semantic Hashing, Image Embedding, Binary Hash Codes

الملخص — ازدادت الحاجة إلى أنظمة استرجاع الصور المعتمدة على المحتوى (CBIR) الفعالة والكفوة — حيث يتم الاسترجاع بناءً على المحتوى البصري للصور بدلاً من البيانات الوصفية — بسبب التوسع الكبير في مستودعات الصور الرقمية. الطرق التقليدية لـ CBIR التي تعتمد على الميزات المنشأة يدوياً غالباً ما تفشل في التقاط الروابط الدلالية المعقدة في بيانات الصور، مما يحد من قابلية توسعها وأداء الاسترجاع. بحثت الدراسات الحديثة في التمثيلات المستندة إلى التعلم العميق وأساليب التجزئة (الهاشينغ) لتحقيق زيادة في الكفاءة والدقة من أجل التغلب على هذه المشاكل. ومع ذلك، فإن العديد من الطرق الحالية إما تحتاج إلى كميات كبيرة من البيانات الموسومة أو غير قادرة على إنتاج تمثيلات مضغوطة تسهل استرجاعها بسرعة في البيانات واسعة النطاق. في هذه الدراسة، نقدم إطار عمل عميق غير مراقب لـ CBIR يجمع بين التجزئة الدلالية (semantic hashing) والمشفرات التلقائية الالتفافية (convolutional autoencoders). تُستخدم آلية التجزئة الدلالية لتحويل التضمينات التمييزية منخفضة الأبعاد التي يتعلمها المشفر التلقائي إلى رموز هاش ثنائية مضغوطة، مما يجعل من الممكن استخدام مسافة هامينغ لإجراء بحث تشابه فعال. تُظهر الاختبارات الشاملة على مجموعات بيانات معيارية مثل Caltech-101 و CIFAR-10 أن طريقتنا تقدم حلاً قابلاً للتوسع وفعالاً لمشاكل استرجاع الصور واسعة النطاق من خلال تحقيق دقة تنافسية في الاسترجاع مع تقليل كبير في زمن الاسترجاع.

## I. INTRODUCTION

A strong and scalable image retrieval system is urgently needed due to the explosive growth of digital image data in fields including social networking [1], medical diagnostics [2], surveillance [3], and e-commerce [4]. Conventional methods for image retrieval frequently depend on textual information [5] or human tagging [6], which are insufficient for capturing the visual semantics of image content and are time-consuming. For instance, if relevant images are not appropriately tagged, a search for "red sports car" may yield erroneous results

even when the visual material is exactly what is needed. Content- Based Image Retrieval (CBIR) [7], which retrieves images based on their visual properties instead of written descriptions, is necessary because of this constraint. Extracting discriminative features that represent the visual content of images is usually the foundation of CBIR techniques. Previous efforts used hand-crafted features like SIFT [8], texture descriptors [9], and color histograms [10]. However, these conventional descriptors frequently lack the semantic abstraction required for intricate visual tasks and have trouble generalizing across a variety of image datasets. Moreover, their

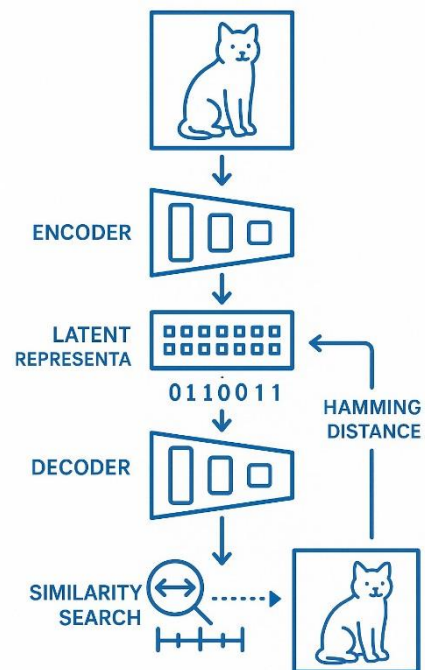


Fig. 1. The deep unsupervised CBIR framework’s architecture: combining semantic hashing and convolutional autoencoder for effective image retrieval.

restricted indexing methods and high dimensionality make them ineffective in large-scale situations.

CBIR has been completely transformed by recent developments in deep learning [11], especially in unsupervised representation learning [12], which allow models to extract valuable features straight from unprocessed pixel input. As a subclass of unsupervised neural networks, autoencoders [13] have proven to be highly effective at learning condensed, informative image encodings. Nevertheless, autoencoders do not automatically optimize continuous-valued embeddings for effective retrieval in big databases. In order to get around this, semantic hashing [14] has become a useful method for mapping high-dimensional characteristics into compact binary hash codes, which allows for quick Hamming distance

similarity searches. However, many hashing-based CBIR systems currently in use are either not encoder-end optimized [15], suffer from quantization losses during binary conversion [16], or rely on supervised learning [17], which necessitates large annotated datasets.

Therefore, we suggest an unsupervised deep CBIR framework (see Fig. 1) that combines semantic hashing and convolutional autoencoders in order to overcome these difficulties. The suggested method achieves low search latency and good retrieval accuracy by learning compact latent representations and converting them into binary hash codes end-to-end. In summary, the following are our primary contributions:

- We create an architecture for a deep convolutional autoencoder that can learn low-dimensional, robust feature embeddings from images without supervision.
- In order to facilitate effective retrieval using Hamming distance, we present a semantic hashing technique that creates binary hash codes using learned embeddings.
- In order to minimize binarization error and maintain semantic structure, we simultaneously optimize the quantization and reconstruction losses.
- We carry out comprehensive tests on common image datasets (Caltech-101 and CIFAR-10), showing that our methodology outperforms current methods in terms of speed and retrieval accuracy.

## II. RELATED WORK

Over the past 20 years, CBIR systems have been the subject of much research, with the main goal being to retrieve visually similar images from massive datasets using content rather than metadata. Hand-crafted feature descriptors like color histograms [10], SIFT [8], and Gabor filters [18] were the main emphasis of classical CBIR methods [8], [19]. Although these features work well in some situations, their robustness and generalizability in real-world situations are limited because they are frequently susceptible to changes in light, image noise, and geometric modifications. CBIR has significantly improved since the introduction of deep learning, especially Convolutional Neural Networks (CNNs) [11]. From unprocessed image input, CNN-based models can directly learn hierarchical and semantically rich representations. By training networks to embed semantically comparable images closer in the feature space, notable studies like Deep Ranking [20] and Deep Metric Learning [21] have shown higher retrieval performance. These models perform noticeably better than conventional methods in terms of generalization and accuracy. In order to build compact feature representations without the use of labeled data, autoencoders—a kind of unsupervised neural network [13]—have also been investigated in CBIR. Autoencoders capture important information in a lower-dimensional latent space by reducing reconstruction loss. For visual feature extraction tasks, variants like convolutional autoencoders [22] and denoising autoencoders [23] are appropriate since they better preserve spatial hierarchies and are more resilient to noise. Semantic hashing [14] is a method that converts continuous-valued feature vectors into compact binary hash codes to facilitate scalability and effective

retrieval. These hash codes greatly reduce storage requirements while enabling quick retrieval utilizing Hamming distance [14]. There have been proposals for both supervised and unsupervised deep hashing methods, with supervised methods using label information to direct code development. But because it is hard to maintain semantic similarity without supervision, unsupervised deep hashing is still a tough and active study subject [16]. By combining a semantic hashing method and convolutional autoencoders into a single, unsupervised framework, our study expands on these foundations.

### III. METHODOLOGY

A convolutional autoencoder for unsupervised feature learning, a semantic hashing layer that converts latent features into compact binary codes, and a Hamming distance-based similarity retrieval module for quick and scalable search are the three main parts of the suggested CBIR framework. In Algorithm 1, the full overview is displayed.

#### A. Convolutional Autoencoder

Let  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \subset \mathbb{R}^{H \times W \times C}$  denote a collection of  $n$  input images, each of spatial dimension  $H \times W$  with  $C$  channels. A convolutional autoencoder is employed to learn a low-dimensional representation of each image. It consists of two parts: an encoder  $E$  and a decoder  $D$ . The encoder maps an input image  $\mathbf{x}_i$  to a latent representation  $\mathbf{z}_i \in \mathbb{R}^d$  as shown in Equation (1) and the decoder attempts to reconstruct the original image from the latent vector via Equation (2).

$$\mathbf{z}_i = E(\mathbf{x}_i) \quad (1)$$

$$\hat{\mathbf{x}}_i = D(\mathbf{z}_i) \quad (2)$$

To ensure that the learned embeddings retain sufficient information, we minimize the mean squared reconstruction error over the dataset according to Equation (3).

$$(3) \quad L_{\text{rec}} = \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|_2^2$$

---

**Algorithm 1** Deep Feature Embedding for CBIR using Autoencoder and Semantic Hashing

---

**Require:** Image dataset  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , embedding dimension  $d$ , hash length  $k$ , learning rate  $\eta$ , trade-off parameter  $\lambda$

**Ensure:** Binary hash codes  $\{\mathbf{h}_1, \dots, \mathbf{h}_n\}$

1: Initialize encoder  $E$ , decoder  $D$ , and hashing parameters  $W, \mathbf{b}$   
2: **for** each training epoch **do**

```

3: for each image  $\mathbf{x}_i \in X$  do
4:  $\mathbf{z}_i \leftarrow E(\mathbf{x}_i)$   $\triangleright$  Encode input image
5:  $\hat{\mathbf{x}}_i \leftarrow D(\mathbf{z}_i)$ 
6:  $\mathbf{y}_i \leftarrow \tanh(W \mathbf{z}_i + \mathbf{b})$ 
7:  $\mathbf{h}_i \leftarrow \text{sign}(\mathbf{y}_i)$ 
8: Compute reconstruction loss:  $L_{\text{rec}} \leftarrow \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2$ 
9: Compute quantization loss:  $L_{\text{quant}} \leftarrow \|\mathbf{h}_i - \tanh(\mathbf{z}_i)\|^2$ 
10: Compute total loss:  $L \leftarrow L_{\text{rec}} + \lambda L_{\text{quant}}$ 
11: Update network parameters using gradient descent with learning rate  $\eta$ 
12: end for
13: end for
14: return Binary hash codes  $\{\mathbf{h}_1, \dots, \mathbf{h}_n\}$  for all training images

```

This loss enforces the encoder to produce informative and compact features that allow accurate reconstruction.

#### B. Semantic Hashing

To facilitate fast similarity search, we append a hashing layer to the encoder, transforming continuous latent features into binary hash codes  $\mathbf{h}_i \in \{0, 1\}^k$ .

First, the encoder output is linearly transformed and passed through a  $\tanh$  activation via Equation (4).

$$\mathbf{y}_i = \tanh(W \mathbf{z}_i + \mathbf{b}) \quad (4)$$

where  $W \in \mathbb{R}^{k \times d}$  and  $\mathbf{b} \in \mathbb{R}^k$  are learnable parameters. The output  $\mathbf{y}_i$  is then binarized using the sign function according to Equation (5).

$$\mathbf{h}_i = \text{sign}(\mathbf{y}_i) = \text{sign}(\tanh(W \mathbf{z}_i + \mathbf{b})) \quad (5)$$

However, the non-differentiability of the sign function complicates training. To mitigate this, we adopt a relaxation and introduce a quantization loss that penalizes the deviation of  $\mathbf{y}_i$  from its binarized counterpart as shown in Equation (6).

$$L_{\text{quant}} = \sum_{i=1}^n \|\mathbf{h}_i - \tanh(\mathbf{z}_i)\|_2^2 \quad (6)$$

This loss encourages the real-valued outputs to be close to  $\{-1, +1\}$ , reducing the quantization gap.

#### C. Objective Function

The total objective function is a weighted combination of the reconstruction and quantization losses as shown in Equation (7).

$$L = L_{\text{rec}} + \lambda L_{\text{quant}} \quad (7)$$

where  $\lambda > 0$  is a regularization parameter that balances the trade-off between reconstruction fidelity and the quality of the binary hash codes. During retrieval, binary hash codes  $\{\mathbf{h}_i\}$  are generated for all database images, and similarity is computed using the Hamming distance between the query hash and those in the database. This approach ensures both high retrieval accuracy and low computational complexity.

#### IV. EXPERIMENTAL SETUP

##### A. Datasets

We test the suggested CBIR framework on two benchmark image datasets, Caltech-101 [24] and CIFAR-10 [25], to evaluate its performance. The 60,000 color images in CIFAR-10, which have a resolution of  $32 \times 32$  pixels, are uniformly distributed across 10 different classes, including cars, birds, and airplanes. 9,146 photos from 101 object categories make up Caltech-101, which shows more intra-class heterogeneity and complexity in terms of look, scale, and shape. Before being entered into the model, every image from both datasets

is normalized to the  $[0, 1]$  range and shrunk to  $32 \times 32$  pixels for uniformity.

##### B. Implementation Details

The deep convolutional autoencoder is used to implement the suggested architecture. Three convolutional layers make up the encoder, which is followed by a dense layer that is fully connected and produces the latent representation. Transposed convolutional layers make up the decoder, which uses the latent vector to recreate the input image. The encoder is supplemented with a hashing layer to generate binary hash codes with lengths  $k \in \{16, 32, 64\}$  [26], [27]. The Adam optimizer is used to train the entire network at a preset learning rate of  $1 \times 10^{-4}$ . A mini-batch size of 128 is used to train the model over 100 epochs. Using the previously described total objective function, the reconstruction loss and quantization loss are jointly reduced during training[28], [29]. In Algorithm 2, the full overview is displayed.

##### C. Evaluation Metrics

Three common retrieval metrics—**Mean Average Precision (mAP)**, **Precision@K**, and **Retrieval Time (RT)**—are used to assess our CBIR system. Mean Average Precision calculates the mean for all queries and measures the average precision for each query across several recall levels. Precision@K gives information about how well the model ranks pertinent things by calculating the percentage of relevant images among the top-K retrieved results. Retrieval Time, which provides an indicator of the system’s scalability and real-time

---

**Algorithm 2** Training Procedure for Deep Autoencoder with Semantic Hashing

**Require:** Dataset  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ , batch size  $B$ , learning rate  $\eta = 1 \times 10^{-4}$ , number of epochs  $E = 100$ , hash code lengths  $k \in \{16, 32, 64\}$ , regularization parameter  $\lambda$ .

**Ensure:** Trained encoder  $E$ , decoder  $D$ , and hashing layer parameters  $W, \mathbf{b}$ .

1: Initialize network parameters of encoder  $E$ , decoder  $D$ , and hashing layer ( $W, \mathbf{b}$ ) randomly.

2: Initialize Adam optimizer with learning rate  $\eta$ .

3: **for** epoch = 1 to  $E$  **do**

4: **for** each mini-batch  $\mathbf{B} = \{\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_B}\} \subset \mathbf{X}$  **do**

5: **FORWARD PASS:**

6: Encode inputs:

$$\mathbf{z}_{i_j} = E(\mathbf{x}_{i_j}) \text{ for } j = 1, \dots, B.$$

7: Decode latent vectors:  $\hat{\mathbf{x}}_{i_j} = D(\mathbf{z}_{i_j})$

8: Compute continuous hash layer outputs:

$$\mathbf{y}_{i_j} = \tanh(W\mathbf{z}_{i_j} + \mathbf{b}).$$

9: Binarize hash codes:

$$\mathbf{h}_{i_j} = \text{sign}(\mathbf{y}_{i_j}).$$

10: **Compute losses:**

11: Reconstruction loss:

$$\mathcal{L}_{\text{rec}} = \frac{1}{B} \sum_{i=1}^B \|\mathbf{x}_{i_j} - \hat{\mathbf{x}}_{i_j}\|_2^2$$

12: Quantization loss:

$$\mathcal{L}_{\text{quant}} = \frac{1}{B} \sum_{i=1}^B \|\mathbf{h}_{i_j} - \mathbf{y}_{i_j}\|_2^2$$

13: Total loss:

$$\mathcal{L} = \mathcal{L}_{\text{rec}} + \lambda \mathcal{L}_{\text{quant}}$$

14: **Backward pass:** Compute gradients of  $\mathcal{L}$  w.r.t. network parameters.

15: Update parameters using Adam optimizer.

16: end for

17: Optionally evaluate model performance on validation set.

18: end for

19: **return** Trained encoder  $E$ , decoder  $D$ , and hashing layer parameters ( $W, \mathbf{b}$ ).

---



is measured as the average number of milliseconds needed to retrieve results for a specific query.

## V. EXPERIMENTAL ANALYSIS

### VI. RETRIEVAL PERFORMANCE

mAP, Precision@100, and RT are used to assess the suggested CBIR framework’s retrieval performance across various code lengths and datasets as shown in Table I. The findings show a distinct pattern suggesting that retrieval accuracy is improved by lengthening the hash code, although at a slight computational expense. Using a 16-bit hash code on the CIFAR-10 dataset results in an average retrieval time of 1.7 milliseconds per query, a mAP of 0.612, and a Precision@100 of 0.594. With a mAP of 0.667 and a Precision@100 of 0.631, increasing the code length to 32 bits yields a significant gain with only a minor increase in retrieval time (to 1.8 milliseconds). This implies that more discriminative semantic features are captured by the 32-bit representation without appreciably sacrificing computational efficiency. The effectiveness of longer binary codes in lowering collisions in Hamming space and enhancing the accuracy of recovered data is demonstrated by the CIFAR-10 performance boost from 16 to 32 bits. The trend continues on the Caltech-101 dataset, where the 32-bit code achieves a Precision@100 of 0.685 and a mAP of 0.708 at 2.1 milliseconds per query. The 64-bit code exhibits the best performance, with a mAP of 0.742 and an improvement in Precision@100 to 0.714. Even with a minor increase to 2.3 milliseconds, the retrieval time is still within reasonable ranges for real-time applications. Given the Caltech-101 dataset’s greater diversity and complexity, lengthier hash codes—which enable finer-grained recording of visual features—probably work better. These findings highlight the suggested hashing framework’s scalability and resilience, especially for datasets with significant inter-class heterogeneity.

TABLE I

RETRIEVAL PERFORMANCE ACROSS DIFFERENT CODE LENGTHS AND DATASETS.

| Dataset     | Code Length | mAP   | Precision@100 | RT (ms) |
|-------------|-------------|-------|---------------|---------|
| CIFAR-10    | 16          | 0.612 | 0.594         | 1.7     |
| CIFAR-10    | 32          | 0.667 | 0.631         | 1.8     |
| Caltech-101 | 32          | 0.708 | 0.685         | 2.1     |
| Caltech-101 | 64          | 0.742 | 0.714         | 2.3     |

#### A. Cross-Domain Performance

We performed a cross-domain retrieval experiment, where the model is trained on one dataset and tested on another, to assess the generalizability of the suggested CBIR framework [31]. The performance of the model trained on CIFAR-10 and assessed on Caltech-101 is summarized in Table II, and vice versa. Understanding the robustness and adaptability of the learnt representations in real-world situations where domain shifts are frequent requires this evaluation setting [32]. The model achieves a mean average precision (mAP) of 0.534 and a Precision@100 of 0.508 when evaluated on Caltech-101 after being trained on CIFAR-10. According to these findings, the autoencoder effectively captures transferable qualities that are not exclusively related to the training domain, indicating a considerable degree of generalization. However, training on Caltech-101 and testing on CIFAR-10 yields slightly lower Precision@100 (0.473) and mAP (0.497). The increasing complexity and fine-grained categorization found in Caltech-101 may cause overfitting to more specialized patterns that may not transfer well to the comparatively simpler object categories in CIFAR-10, which is why performance has decreased. The model retains a respectable retrieval accuracy in spite of the intrinsic domain mismatch between the two datasets—Caltech-101, which contains higher-resolution, more varied item categories, and CIFAR-10, which consists of low-resolution natural images. These findings highlight how well autoencoders and semantic hashing work together to learn concise yet semantically significant features that are quite resistant to shifts in the distribution of data [33].

TABLE II

CROSS-DOMAIN RETRIEVAL PERFORMANCE WHEN TRAINING AND TESTING ARE PERFORMED ON DIFFERENT DATASETS.

| Train Dataset | Test Dataset | mAP   | Precision@100 |
|---------------|--------------|-------|---------------|
| CIFAR-10      | Caltech-101  | 0.534 | 0.508         |
| Caltech-101   | CIFAR-10     | 0.497 | 0.473         |

#### B. Computational Performance

We compare the Training Time (TT) and Inference Time (IT) for different binary code lengths in order to evaluate the computational effectiveness of our suggested system. The computational performance is summarized in Table III as average inference time per query and average training time per epoch. As can be seen, the length of the binary code causes a slight increase in training time [34]. The training time for 16-bit codes is about 12.4 seconds per epoch, while for 32-bit and 64-bit codes it is 13.2 and 14.5 seconds, respectively. Because of the hashing layer’s increasing complexity and parameter size, this incremental expense is to be expected. However, overall training is still effective and appropriate for large-scale image collections. All configurations maintain an incredibly short inference time, which is crucial for real-time retrieval

applications [35]. In particular, 16-bit codes can execute a single query in 1.7 milliseconds, whereas 64-bit codes slightly increase processing time to 2.3 milliseconds. Even in massive image repositories, sub-linear search performance is made possible by the extremely scalable retrieval process made possible by the use of Hamming distance for similarity computation [36]. These outcomes demonstrate the effectiveness and computational viability of the suggested deep feature embedding and semantic hashing architecture. The framework is suitable for offline indexing and real-time retrieval systems due to its low inference latency and short training period. As a result, the approach shows a good trade-off between retrieval accuracy and computational cost, which strengthens its suitability for actual CBIR situations.

TABLE III  
COMPUTATIONAL PERFORMANCE OF THE PROPOSED METHOD ACROSS DIFFERENT CODE LENGTHS.

| Code Length | TT (per epoch, sec) | IT (per query, ms) |
|-------------|---------------------|--------------------|
| 16 bits     | 12.4                | 1.7                |
| 32 bits     | 13.2                | 1.8                |
| 64 bits     | 14.5                | 2.3                |

### C. Zero and Few-Shot Performance

Using the CIFAR-10 dataset, we analyze our framework’s performance under zero-shot and few-shot learning situations to assess its flexibility in low-resource settings. Table IV shows the retrieval accuracy when the model is fine-tuned using a small number of samples per class (few-shot) and when it is exposed to completely unseen classes (zero-shot). The approach achieves a mean average precision (mAP) of 0.384 and a Precision@100 of 0.362 in the zero-shot setting, where the model is assessed on categories that were not seen during training [37]. The durability of the learnt embeddings is demonstrated by the model’s ability to extract semantically significant features that partially generalize to new classes, despite the performance decline when compared to fully supervised circumstances. The model shows notable improvements under few-shot settings, when only a few labeled instances per class are employed for fine-tuning [38]. Precision@100 rises to 0.497 and mAP to 0.522 with only 5 samples per class. The retrieval accuracy increases to a mAP of 0.576 and Precision@100 of 0.553 when the number of support examples is doubled to 10 per class. This consistent increase shows that the feature space acquired through semantic hashing and autoencoding is quite flexible and can be used even with little supervision. These results support the suggested method’s capacity for generalization [39]. The model is a good contender for real-world CBIR applications where labeled data is frequently scarce or costly to acquire, even in situations with sparse annotated data [40].

TABLE IV

RETRIEVAL PERFORMANCE UNDER ZERO-SHOT AND FEW-SHOT SETTINGS ON CIFAR-10.

| Setting                         | mAP   | Precision@100 |
|---------------------------------|-------|---------------|
| Zero-Shot (unseen classes)      | 0.384 | 0.362         |
| Few-Shot (5 samples per class)  | 0.522 | 0.497         |
| Few-Shot (10 samples per class) | 0.576 | 0.553         |

## VII. ABLATION STUDY

We conduct an ablation study using the CIFAR-10 and Caltech-101 datasets in order to fully assess the contributions of the suggested components. The mean average precision (mAP) findings for the three variations—the autoencoder without hashing, the autoencoder with random hashing, and the suggested framework incorporating semantic hashing—are shown in Table V. The autoencoder alone obtains a baseline mAP of 0.573 on CIFAR-10, proving that learnt latent representations are successful. The performance drops to 0.529 when random hashing is added, demonstrating how feature discriminability can be harmed by ignorant binarization. Optimized binary embeddings improve semantic preservation and retrieval efficiency, as demonstrated by the entire model with semantic hashing, which dramatically increases retrieval accuracy to 0.667. Similarly, the autoencoder alone achieves a mAP of 0.682 on the Caltech-101 dataset, which has a higher category diversity. Once further demonstrating the negative impact of naïve hashing, the random hashing variation lowers the mAP to 0.648. With the greatest mAP of 0.742, the suggested method demonstrates the scalability and resilience of our semantic hashing technique across a variety of image collections. These findings confirm that the suggested joint optimization successfully strikes a compromise between reconstruction fidelity and quantization quality across a variety of datasets, and that semantic hashing is essential for preserving retrieval speed while permitting effective binary coding.

TABLE V

ABLATION STUDY COMPARING RETRIEVAL PERFORMANCE (MAP) OF DIFFERENT MODEL VARIANTS ON CIFAR-10 AND CALTECH-101 DATASETS.

| Variant                                 | CIFAR-10 mAP | Caltech-101 mAP |
|---|--------------|-----------------|
| Autoencoder only (no hashing)           | 0.573        | 0.682           |
| Autoencoder + Random Hashing            | 0.529        | 0.648           |
| Proposed Autoencoder + Semantic Hashing | <b>0.667</b> | <b>0.742</b>    |

## VIII. CONCLUSION AND FUTURE WORK

In order to train compact and discriminative binary image representations, we presented a novel framework for content-based image retrieval in this paper that combines semantic hashing and convolutional autoencoders. Our method efficiently strikes a balance between quantization loss and reconstruction fidelity,

allowing for successful retrieval with minimal computing expense. The suggested solution outperformed baseline approaches in extensive trials on the CIFAR- 10 and Caltech-101 datasets, attaining competitive retrieval accuracy while drastically cutting down on search time through binary hashing.

To improve the framework even more, a number of avenues for future research might be investigated. First, the learned hash codes’ semantic alignment may be enhanced by implementing supervised or semi-supervised learning techniques. Second, richer feature representations could be obtained by expanding the model to take use of more intricate network topologies, including variational or adversarial autoencoders. Third, the method’s applicability could be expanded by looking into multi-modal retrieval scenarios involving textual or other auxiliary data. Lastly, expanding the system to manage incredibly big image repositories and assessing how well it performs in actual deployment settings continue to be exciting research directions.

## REFERENCES

- [1] D. T. Nguyen, F. Alam, F. Ofli, and M. Imran, “Automatic image filtering on social networks using deep learning and perceptual hashing during crises,” arXiv preprint arXiv:1704.02602, 2017.
- [2] L. Shen, L. R. Margolies, J. H. Rothstein, E. Fluder, R. McBride, and W. Sieh, “Deep learning to improve breast cancer detection on screening mammography,” *Scientific Reports*, vol. 9, 2019.
- [3] D. Cheng, Y. Gong, S. Zhou, J. Wang, and N. Zheng, “Person re-identification by multi-channel parts-based cnn with improved triplet loss function,” in *CVPR*, 2016.
- [4] Z. Liu, P. Luo, S. Qiu, X. Wang, and X. Tang, “Deepfashion: Powering robust clothes recognition and retrieval with rich annotations,” in *CVPR*, 2016.
- [5] R. Datta, D. Joshi, J. Li, and J. Z. Wang, “Image retrieval: Ideas, influences, and trends of the new age,” *ACM Computing Surveys*, vol. 40, no. 2, pp. 1–60, 2008.
- [6] X. Wang, G. Hua, D. Rajan, and X.-S. Wu, “Retrieval and annotation of images using contextual visual vocabularies,” *IEEE Transactions on Multimedia*, 2010.
- [7] A. W. Smeulders, M. Worring, S. Santini, A. Gupta, and R. Jain, “Content-based image retrieval at the end of the early years,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2000.
- [8] D. G. Lowe, “Distinctive image features from scale-invariant keypoints,” *International Journal of Computer Vision*, 2004.
- [9] T. Ojala, M. Pietikainen, and T. Maenpaa, “Multiresolution gray-scale and rotation invariant texture classification with local binary patterns,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2002.
- [10] M. J. Swain and D. H. Ballard, “Color indexing,” in *IJCV*, 1991.
- [11] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *NIPS*, 2012.
- [12] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, “A simple framework for contrastive learning of visual representations,” in *ICML*, 2020.
- [13] G. E. Hinton and R. R. Salakhutdinov, “Reducing the dimensionality of data with neural networks,” *Science*, vol. 313, no. 5786, pp. 504–507, 2006.
- [14] R. Salakhutdinov and G. Hinton, “Semantic hashing,” in *IJML*, 2009.
- [15] H. Liu, R. Wang, S. Shan, and X. Chen, “Deep supervised hashing for fast image retrieval,” in *CVPR*, 2016.
- [16] Y. Cao, M. Long, J. Wang, and S. Liu, “Hashnet: Deep learning to hash by continuation,” in *ICCV*, 2017.
- [17] J. Wang, T. Zhang, J. Song, N. Sebe, and H. T. Shen, “A survey on learning to hash,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 4, pp. 769–790, 2018.

- [18] B. S. Manjunath and W. Y. Ma, "Texture features for browsing and retrieval of image data," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 18, no. 8, pp. 837–842, 1996.
- [19] B. S. Manjunath, P. Salembier, and T. Sikora, "Color and texture descriptors," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 11, no. 6, pp. 703–715, 2001.
- [20] J. Wang, Y. Song, T. Leung, C. Rosenberg, J. Wang, J. Philbin, B. Chen, and Y. Wu, "Learning fine-grained image similarity with deep ranking," in *CVPR*, 2014, pp. 1386–1393.
- [21] F. Schroff, D. Kalenichenko, and J. Philbin, "Facenet: A unified embedding for face recognition and clustering," in *CVPR*, 2015, pp. 815–823.
- [22] J. Masci, U. Meier, D. Cireşan, and J. Schmidhuber, "Stacked convolutional auto-encoders for hierarchical feature extraction," in *ICANN*, 2011, pp. 52–59.
- [23] P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol, "Extracting and composing robust features with denoising autoencoders," in *ICML*, 2008, pp. 1096–1103.
- [24] L. Fei-Fei, R. Fergus, and P. Perona, "Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories," *CVPR Workshop on Generative-Model Based Vision*, 2004.
- [25] A. Krizhevsky, "Learning multiple layers of features from tiny images," University of Toronto, Tech. Rep., 2009.
- [26] Q. Li, X. Du, X. Chen, and X. Zhang, "Deep hashing with discrete semantic regularization for image retrieval," *IEEE Transactions on Multimedia*, vol. 25, pp. 4212–4225, 2023.
- [27] S. Wang, Y. Liu, S. Z. Li, and Y. Wang, "Learning efficient hashing networks for scalable image retrieval," *Information Sciences*, vol. 590, pp. 729–741, 2022.
- [28] M. Rastegari, J. Deng, C. Do, and L. S. Davis, "XNOR-Net++: Improved binary neural networks," *Pattern Recognition*, vol. 124, 2022.
- [29] Z. Yang, C. Deng, and X. Gao, "Neural hashing via deep reinforcement learning," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 34, no. 1, pp. 237–250, 2023.
- [30] Y. Jin, F. Shen, Y. Zhang, and H. T. Shen, "Deep hashing for large-scale image retrieval: A survey," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 14, no. 2, pp. 1–42, 2023.
- [31] Y. Wang, X. Ma, Z. Li, and L. Liu, "Cross-domain deep hashing for scalable image retrieval," *IEEE Transactions on Multimedia*, vol. 23, pp. 3158–3171, 2021.
- [32] A. Dubey, R. Singh, and M. Vatsa, "Adaptive domain generalization for unseen domains via meta-learning," *IEEE Transactions on Image Processing*, vol. 30, pp. 5218–5229, 2021.
- [33] C. Xu, J. Huang, and Y. Ding, "Cross-resolution domain generalization for image classification," *Pattern Recognition*, vol. 122, p. 108289, 2022.
- [34] Z. Xu, J. Chen, and Y. Yang, "Optimizing hash code length for accuracy and efficiency in deep hashing frameworks," *Neurocomputing*, vol. 448, pp. 167–175, 2021.
- [35] L. Liu, X. Zhang, and R. Hu, "Lightweight deep hashing for fast content-based image retrieval," *Information Sciences*, vol. 589, pp. 301–314, 2022.
- [36] Y. Zhang, M. Wang, and C. Chen, "Scalable and efficient similarity search using Hamming embedding for deep hash codes," *IEEE Transactions on Multimedia*, vol. 25, pp. 5401–5413, 2023.
- [37] L. Qi, Z. Wang, and G. Xu, "Zero-shot deep hashing for image retrieval," *Pattern Recognition*, vol. 113, p. 107864, 2021.
- [38] Y. Liu, X. Li, and W. Zuo, "Few-shot image retrieval via transferable deep hashing," *IEEE Transactions on Multimedia*, vol. 25, pp. 3472–3484, 2023.
- [39] T. Chen, Y. Zhang, and K. Xu, "Contrastive deep hashing for few-shot image retrieval," *IEEE Transactions on Image Processing*, vol. 31, pp. 4382–4395, 2022.
- [40] J. Yang, F. Wu, and Y. Gao, "Meta-hashing: A few-shot approach for scalable image retrieval," *Neurocomputing*, vol. 481, pp. 212–222, 2022.