

CATIRI: An Efficient Method for Content-and-Text Based Image Retrieval

Mengqi Zeng¹, Bin Yao¹, Member, CCF, ACM, IEEE, Zhi-Jie Wang^{2,5,6}, Member, CCF, ACM, Yanyan Shen¹, Feifei Li³, Senior Member, IEEE, Member, ACM, Jianfeng Zhang⁴, Hao Lin⁴, and Minyi Guo¹, Fellow, IEEE, Member, CCF, ACM

¹*Department of Computer Science and Engineering, Shanghai Jiao Tong University, Shanghai 200240, China*

²*School of Data and Computer Science, Sun Yat-Sen University, Guangzhou 510006, China*

³*School of Computing, University of Utah, Salt Lake City 84112, USA*

⁴*Alibaba Group, Hangzhou 311121, China*

⁵*Guangdong Key Laboratory of Big Data Analysis and Processing, Guangzhou 510006, China*

⁶*National Engineering Laboratory for Big Data Analysis and Applications, Beijing 100871, China*

E-mail: zmq19950714@126.com, yaobin@cs.sjtu.edu.cn, wangzhij5@mail.sysu.edu.cn, shen-yy@cs.sjtu.edu.cn

lifeifei@cs.utah.edu, xingdian@taobao.com, bixuan@taobao.com, guo-my@cs.sjtu.edu.cn

Received July 09, 2018; revised December 12, 2018.

Abstract The combination of visual and textual information in image retrieval remarkably alleviates the semantic gap of traditional image retrieval methods, and thus it has attracted much attention recently. Image retrieval based on such a combination is usually called the *content-and-text based image retrieval* (CTBIR). Nevertheless, existing works in CTBIR mainly make efforts on improving the retrieval quality. To the best of knowledge, little attention has been focused on how to enhance the retrieval efficiency. Nowadays, image data is widespread and expanding rapidly in our daily life. Obviously, it is important and interesting to investigate the retrieval efficiency. To this end, this paper presents an efficient image retrieval method named *CATIRI* (Content-And-Text based Image Retrieval using Indexing). *CATIRI* follows a three-phase solution framework that develops a new indexing structure called MHIM-tree. The MHIM-tree seamlessly integrates several elements including Manhattan Hashing, Inverted file, and M-tree. To use our MHIM-tree wisely in the query, we present a set of important metrics and reveal their inherent properties. Based on them, we develop a top- k query algorithm for CTBIR. Experimental results based on benchmark image datasets demonstrate that *CATIRI* outperforms by an order of magnitude a competitor.

Keywords image retrieval, text-and-visual feature, indexing, top- k

1 Introduction

With the booming development of electronic information technology and the network communication technology, massive image data are generated every day, mainly from social networks (e.g., Facebook, Twitter and Flickr) and general image databases (e.g., Google Images and Baidu Images) [3, 4]. The scale of image databases is now in the hundreds of millions, and is expanding rapidly [3, 5, 6]. Image retrieval as a

fundamental operation in image database plays a significant role [7, 8].

One of representative branches is the text-based image retrieval (TBIR)[9, 10]. It is based on annotations associated with images, and the relevant images are retrieved from the databases by matching the textual query with the textual annotations of images [1, 10]. TBIR is simple and fast, and has been widely used on the Internet (e.g., many social image search engines).

Yet, TBIR bears some limitations due to manual labelling, e.g., expensive labor and time cost, covering a limited part of semantic interpretations on the image content, inconsistent and/or biased annotations [3]. To alleviate these issues some automatic image annotation techniques have been proposed [1].

Another branch is the content-based image retrieval (CBIR) [11, 12]. It is based on low-level visual properties (e.g., color, texture, sketch) extracted from the images, and in CBIR systems the similarity of two images is measured by these visual features [11, 13]. CBIR thoroughly overcomes the disadvantages of TBIR, since it does not use annotations. Hence, much attention has been attracted to CBIR, some researches are devoted to the retrieval quality (see e.g., [14, 11, 15]), while others are devoted to the retrieval efficiency (see e.g., [16, 5, 17, 12, 8, 7]), since CBIR is much more time-consuming than TBIR. Nevertheless, there exists a critical issue [3]: the semantic gap between low-level visual features and high-level semantic features. Typically, the semantic gap refers to the mismatch between the information extracted from an image and the interpretation for the users [14]. Two major reasons incur semantic gap [3, 14]: one is the impossibility for users to describe the query image exactly and adequately, and the other is the incomprehension about users' intention behind the query.

The researches [18, 19] have shown that combining textual and visual information even with simple fusion strategies may improve the quality of retrieval results. Thereby, it has attracted a lot of attention in recent years [20, 4, 19, 21, 22], and typically it is called the content-and-text based image retrieval (CTBIR). In existing literature, efforts for CTBIR are mostly on improving the retrieval quality by exploiting various techniques such as latent semantic kernels [21], semantic combination [22], composite correlation quantization model [4], etc. To the best of our knowledge, few attention has focused on the retrieval efficiency in CTBIR. One can easily understand that CTBIR is also time-consuming, since we need to consider both textual and visual information, let alone the increasing scale of image databases. In view of the facts above, we make an attempt to investigate the retrieval efficiency of CTBIR.

Specifically, this paper proposes an efficient CTBIR method dubbed as *CATIRI*, which can preserve the quality of retrieval results and is designed for high efficiency. The key principle of *CATIRI* is to model the top- k image retrieval problem as a problem similar to spatial-keyword query, and solve the reduced problem by exploiting (i) a new indexing structure called MHIM-tree that can be viewed as a seamlessly blend of several existing techniques such as Manhattan hashing [23], inverted file [24] and M-tree [25], and (ii) a

set of important properties that allow us to prune most of nodes in the query processing. As we demonstrate with benchmark datasets, *CATIRI* is significantly superior to the baseline. To summarize, the novelty and contributions of this paper are as follows.

- We reveal that the retrieval efficiency for CBTIR is also urgently important yet it is widely ignored. This paper makes the first deep investigation on the retrieval efficiency of CBTIR.
- We show the top- k image retrieval problem can be modelled as a problem similar to the spatial-keyword query. This finding could open a hopeful direction for image retrieval.
- We present a top- k image retrieval method, *CATIRI*, that yields a new indexing structure and reveals a set of important properties.
- We give rigorous theoretical analysis for our method, and conduct a comprehensive experiments to evaluate its performance.

The rest of the paper is organized as follows. Section 2 reviews previous work most related to ours. Section 3 formulates our problem, and Section 4 presents our indexing structure MHIM-tree. In Section 5 several important metrics and their inherent properties are examined, and Section 6 presents our three-phase solution. Section 7 presents the theoretical analysis for our proposed indexing as well as the retrieval algorithm. In Section 8 we experimentally evaluate our proposed solution, and finally Section 9 concludes the paper.

2 Related work

Image retrieval is a classic and hot topic in tens of years. We review existing image retrieval methods most related to ours. Generally, they can be classified into three categories: (i) text-based image retrieval (TBIR); (ii) content-based image retrieval (CBIR); and (iii) content-and-text based image retrieval (CTBIR).

TBIR. In TBIR, images are annotated with text that represents high-level semantics, and image retrieval is performed by text retrieval techniques. For example, Zhang *et al.* [9] suggested a user-term feedback based technique for text-based image retrieval. Li *et al.* [10] proposed an approach to learn a robust classifier for text-based image retrieval. TBIR is simple and fast yet manual labelling is labor intensive process [38], and so many researches focus on automatic image annotation techniques [1].

CBIR. Correspondingly, in CBIR, image content is used to measure similarity, which is described by visual features such as color, shape and texture. CBIR overcomes the disadvantages of manual labelling, but

it is not so satisfactory [14] due to the semantic gap between low-level visual features and high-level semantics. Many researches have been devoted to narrowing or bridging the semantic gap, and various techniques and algorithms are proposed. For example, Tong *et al.* [11] proposed a support vector machine active learning algorithm for content-based image retrieval. Yang *et al.* [39] suggested multiple-instance learning techniques for image retrieval. Deng *et al.* [15] incorporated prior human knowledge in the form of a hierarchical structure to perform content-based image retrieval.

Compared to TBIR, CBIR is much more time-consuming. A lot of works have been devoted to the efficiency in CBIR. For example, Xia *et al.* [17] developed Multi-Kernel Locality Sensitive Hashing schema, significantly improving the retrieval performance of KLSH [13] by utilizing multiple kernels. Natsev *et al.* [40] proposed WALRUS retrieval algorithm, which extracts features of region and compute their signature as index to measure similarity. Shen *et al.* [7] addressed how to speed up *interactive image retrieval*, where they assume a query is interactively refined towards the optimal answers by exploiting user feedback. Falchi *et al.* [8] studied the possibility of caching the answers to content-based image retrieval queries in metric space, with the aim of reducing the average cost of query.

Although these line of works are related to our work, it can be seen that they are clearly different from ours, since this paper mainly focuses on the content-and-text based image retrieval, instead of CBIR.

CTBIR. In CTBIR, it can take advantage of both visual and textual information, and bridge the semantic gap. Thereby, it has attracted a lot of attention in recent years [20, 4, 19, 21, 22]. For example, Chatzichristofis *et al.* [20] proposed a two-stage approach. They first use a text modality to rank the collection and then perform CBIR only on the top-k items. Long *et al.* [4] proposed a composite correlation quantization model for image retrieval. Their model jointly finds correlation-maximal mappings, and learns composite quantizers that convert the isomorphic latent features into compact binary codes. Zhou *et al.* [19] proposed a seamless joint querying and relevance feedback scheme, based on both keywords and low-level visual contents. Caicedo *et al.* [21] proposed a strategy to fuse visual features and unstructured-text data in a medical image retrieval system. Clinchant *et al.* [22] proposed a set of techniques called semantic combination that better manages the complementarities between text and image search systems.

This line of works are mainly devoted to improving the retrieval quality by developing various fusion techniques. To the best of our knowledge, less attention has been focused on the efficiency in CTBIR. This paper

makes an attempt to address the efficiency issue in CTBIR. Remark that, the demo paper [37] could be closest to our work, since they “mentioned” the retrieval efficiency by directly exploiting an existing technique “inverted index”. Nevertheless, their focus is still on the retrieval quality, as indicated in their experiments. Essentially, we compare a similar version of [37]. Our solution *CATIRI* achieves excellent query performance, and is significantly superior to the baseline.

Others. There are also other works that close to our, but clearly different from our work. For example, Rabbitti *et al.* [41] presented an approach to image retrieval that allows to take into account the imprecision assigned to the different parts of the query, and to rank the images on the basis of this imprecision. Chu *et al.* [42] introduced a semantic data model to capture the hierarchical, spatial, temporal, and evolutionary semantics of images in pictorial databases. Brown *et al.* [43] proposed a prototype content-based image retrieve system that aims to save space cost. Chen *et al.* [44] developed a distributed retrieval and recommendation system for geo-textual images.

3 Preliminaries

3.1 Problem Formulation

Let D be an image database, and each image object I in D is defined as a pair (I_v, I_t) , where I_v is a vector of visual features (e.g., GIST [26] and bag-of-words [27]), and I_t is a document consisting of textual annotations. Usually, document I_t can be represented by a vector, in which each dimension is the weight of a distinct term in the document [28, 29].

Let D_c be the collection of all documents in database D . The weight of a term t in I_t can be computed as follows [29]:

$$w(I_t, t) = p(t|I_t) = (1 - \lambda) \frac{tf(t, I_t)}{|I_t|} + \lambda \frac{tf(t, D_c)}{|D_c|} \quad (1)$$

where $tf(t, I_t)$ (resp., $tf(t, D_c)$) is the term frequency of term t in I_t (resp., D_c), $|I_t|$ (resp., $|D_c|$) is the total number of tokens in document I_t (resp., all documents of D_c), $\lambda \in [0, 1]$ is a smoothing coefficient of the Jelinek-Mercer method [29]. Notice that, here $tf(t, I_t)/|I_t|$ (resp., $tf(t, D_c)/|D_c|$) is essentially the maximum likelihood estimate of term t in I_t (resp., D_c).

Let Q be a query defined as a pair (Q_v, Q_k) , where Q_v is a vector of visual features, and Q_k is a set of query keywords. The text relevancy, between a given query Q and an image object I , can be regarded as the probability that query keywords Q_k appear in the document

Table 1. Notations

Notation	Meaning
D/D_c	image database; collection of documents in D
$I/I_v/I_t$	image object; vector of visual features of I ; document of I
$w(I_t, t)$	weight of term t in document I_t
$Q/Q_v/Q_k$	top- k query; vector of visual features of Q ; keywords of Q
$P(Q_k I_t)$	text relevancy between Q_k and I_t
$vDist(Q_v, I_v)$	distance between Q_v and I_v
$Dist(Q_v, I_v)$	Manhattan distance between the hashcodes of Q_v and I_v
$S_v(\cdot)/S_t(\cdot)$	visual similarity; text relevancy
$S(\cdot)/sDiff(\cdot)$	similarity score; change to insert an entry
$R/rDiff(\cdot)$	minimum bounding sphere of entry; enlargement of the covering radius of R
$\mathbb{V}/vDiff(\cdot)$	document vector of image object or node; difference between vectors
E/C	entry in a node; category label(s)

I_t . And it can be computed as follows:

$$P(Q_k|I_t) = \prod_{t \in Q_k} p(t|I_t) \quad (2)$$

Correspondingly, the visual similarity, between a given query Q and an image object I , can be usually obtained by computing the distance between their visual vectors. That is,

$$vDist(Q_v, I_v) = \|Q_v - I_v\| \quad (3)$$

Furthermore, the similarity score between an image object I and the query Q can be computed as follows [30]:

$$S(Q, I) = \alpha \mathbf{N}(S_v(Q, I)) + (1 - \alpha) \mathbf{N}(S_t(Q, I)) \quad (4)$$

where $S_v(\cdot)$ represents the visual similarity between Q and I , $S_t(\cdot)$ represents the text relevancy between Q and I ; $\mathbf{N}(\cdot)$ is a normalization operator that transforms them into a range $[0, 1]$, and $\alpha \in [0, 1]$ is a balance parameter between the visual similarity and text relevancy. Note that, the similarity score described above is a linear combination of text relevance and visual similarity. Nevertheless, our solution can also work for other functions, since it is independent of the ‘‘fusion’’ functions. Formally, in this paper we focus on the following problem.

Definition 1 (Top- k image retrieval). *Given an image database D and a query Q , the top- k image retrieval is asked to return k image objects that are ranked highest in terms of the similarity scores.*

3.2 Problem Analysis

It is widely accepted [3, 13, 17] that feature hashing can map visual features of image into the compact binary code. Most of existing hashing methods adopt

Hamming distance to measure the similarity between points in the hashcode space. Yet, these hashing methods may destroy the neighborhood structure in the original feature space, violating the essential goal of feature hashing [31].

Recently, a new feature hashing technique called *Manhattan Hashing* [23] shows that it can effectively preserve the neighborhood structure in the original feature space. Nevertheless, their work mainly focuses on the retrieval quality for the content-based image retrieval (CBIR). Note that, in this paper we are interested in the content-and-text based image retrieval (CTBIR), and our goal is to achieve a remarkable improvement on the retrieval efficiency.

Nevertheless, Manhattan hashing is still a powerful tool for our work. It enlightens us to re-examine our problem, and particularly, we realize that, by using Manhattan hashing, it may allow us to transform our retrieval problem into a new version that still preserve the inherent equivalence. The followings show the reduction.

Denote by $Dist(Q_v, I_v)$ the Manhattan distance between Q_v and I_v in the hashcode space. One can rewrite Equation 4 as follows:

$$S(Q, I) = \alpha \left(1 - \frac{Dist(Q_v, I_v)}{maxD}\right) + (1 - \alpha) \frac{P(Q_k|I_t)}{maxP} \quad (5)$$

where $maxD$ is the upper bound on $Dist(Q_v, I_v)$ and is used to normalize $Dist(Q_v, I_v)$ into $[0, 1]$, while $maxP$ is the upper bound on $P(Q_k|I_t)$ and is used to normalize the $P(Q_k|I_t)$ into $[0, 1]$.

This reduction is important, since it is the cornerstone of our proposed method. In particular, we observe that, to some extent the reduced problem is similar to spatial-keyword query problems [32], since they also involve text and have two components in the query input. This mightily motivates us to develop solutions by in-

tegrating Manhattan hashing [23] with the techniques widely used in keyword query such as inverted file [24], and used in similarity search such as M-tree [25]. To the best of our knowledge, to date (1) none of existing spatial-keyword query methods can solve our problem, and (2) this is the first work to expose the top- k image retrieval problem can be modelled as a problem similar to the spatial-keyword query. In what follows, we present our solution including indexing structure, important metrics and specific algorithms. For ease of reference, the reader can refer to Table 1 for main symbols used in this paper.

4 The MHIM-tree

In this section, we present a novel indexing structure, named MHIM-tree. Similar to many existing indexing structures such as *sequenced multi-attribute tree* [16], our MHIM-tree is also a blend of several existing techniques (with some specific adaptations). In a nutshell, our MHIM-tree integrates several elements including Manhattan hashing [23], inverted file [24], M-tree [25]. For ease of presentation, we next shortly review these three existing elements. And then we examine the details of the MHIM-tree.

4.1 Three Existing Elements

Manhattan hashing. As mentioned earlier, it is mainly used to transform the visual features of images into binary Manhattan hashcodes that preserve the neighborhood structure in the original feature space. More specifically, it first uses the iterative quantization (ITQ) [2] to project the original data with higher dimension to the lower dimensional space, and then uses Manhattan quantization (MQ) [23] to get the corresponding hashcode, such as 110011 (a 3-dimensional hashcode).

Inverted file. It has a vocabulary of all terms, and each term is associated with an inverted list. Each inverted list comprises a sequence of postings, each of which normally contains the identifier of an object whose description contains the term and the frequency of the term in the description. In general, the postings in each inverted list are sorted by object’s identifier.

M-tree. Similar to B-tree, the M-tree is also a balanced search tree, mainly used for indexing generic multi-dimensional “metric space”, where the distance function satisfies the symmetry, non-negativity and triangle inequality postulates. The key idea is to group nearby objects and represent them with their minimum bounding sphere in the next higher level of the tree. Since all objects lie within this bounding sphere, a query that does not intersect the bounding sphere also cannot intersect any of the contained objects. At

the leaf level, each indexed object is represented by a sphere, where sphere center is the features of object and radius is zero; at higher levels, the aggregation of an increasing number of objects is called *routing object*, represented by a bounding sphere of contained objects.

4.2 Data structure of MHIM-tree

Based on M-tree, the MHIM-tree is a hight-balanced tree with the balancing factor B . For a **leaf node**, it contains entries in the form of (id, R, \mathbb{V}, C) , and each entry represents an image object.

- id : the identifier of an image object I .
- $R(O, r)$: the minimum bounding sphere of image object I , whose center O is the hashcode of I and radius is zero. Remark that, the sphere for a single image object is essentially a *point*. For example, assume the 3-dimensional hashcode of I is 001101, then $R.O = (0, 3, 1)$ and $R.r = 0$.
- \mathbb{V} : a vector that represents the document I_t . For example, assume that there are n distinct terms t_1, t_2, \dots, t_n in D_c , then $\mathbb{V} = \{v_1, v_2, \dots, v_n\}$, where $v_i = w(I_t, t_i)$ (i.e., $p(t_i|I_t)$, recall Equation 1) if $t_i \in I_t$; otherwise $v_i = 0$.
- C : a label for the category containing I . Note that, in general all images in the database D are divided into many categories.

For a **non-leaf node**, it contains entries in the form of (p, R, \mathbb{V}, C) , and each entry represents a routing object corresponding to a child node. Note that, here we slightly abuse the notations R, \mathbb{V}, C but their meanings shall be clear in the context.

- p : a pointer to the corresponding child node.
- $R(O, r)$: the minimum bounding sphere of the routing object that bounds the spheres of all entries in the child node.
- \mathbb{V} : the representative vector of all the document vectors in the subtree rooted at the child node. For example, assume $\mathbb{V} = (v_1, v_2, \dots, v_n)$, then $v_j = \max\{E_i.\mathbb{V}.v_j\}$, where E_i is an entry contained in the child node.
- C : a set of all labels for categories in the child node. That said, $C = \bigcup\{E_i.C\}$.

In addition, each node \mathcal{N} (regardless of leaf or non-leaf node) is attached with an inverted file, which is used to index all the documents in the subtree rooted at \mathcal{N} . The inverted file is composed of two main components:

- A vocabulary of all distinct terms in documents that are in the subtree rooted at \mathcal{N} .
- A collection of inverted lists, each of which is linked to a term t in the vocabulary. Note that, here each inverted list is a set of tuples in the form

of $\langle E, c, w((E, c), t) \rangle$, where E refers to an entry in the node containing term t , c is one of *category labels* in entry E , and $w((E, c), t)$ is the weight of term t in both entry E and category c .

Remark that, when \mathcal{N} is a leaf node, $c = E.C$ and $w((E, c), t) = w(E.id, t)$, since entry E represents an image object. Otherwise, c is a *category label* in the set $E.C$, and $w((E, c), t)$ is the maximum weight of term t in all documents of image objects that belong to the subtree rooted at E and category c , since entry E represents a child node of \mathcal{N} . Here $w((E, c), t)$ is computed as follows:

$$w((E, c), t) = \max\{w((E_i, c), t) | E_i \in E.p\} \quad (6)$$

where E_i is an entry contained in the node represented by E .

Example 1. To further understand the general picture of the MHIM-tree, Figure 1 shows a general picture of our MHIM-tree. Here I_1, I_2, \dots, I_9 are nine image objects, which are grouped into four different categories C_1, C_2, C_3, C_4 .

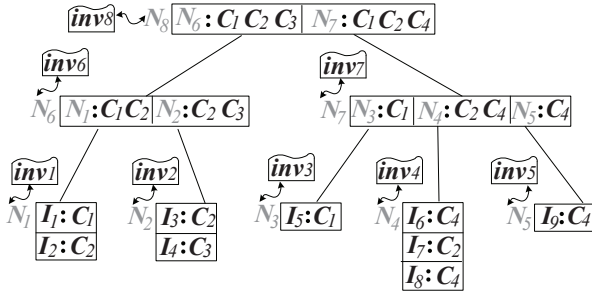


Fig.1. Index structure of the MHIM-tree

5 Important Metrics and Properties

In order to use the MHIM-tree wisely in the retrieval, this section defines a set of important metrics, and meanwhile examines their inherent properties.

As we know, for two points $\mathbf{x} = [x_1, x_2, \dots, x_d]^T$, $\mathbf{y} = [y_1, y_2, \dots, y_d]^T$ in Manhattan space, the Manhattan distance between \mathbf{x} and \mathbf{y} is the sum of differences on all dimensions. The specific calculation formula is as follows [23]:

$$Dist(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^d |x_i - y_i| \quad (7)$$

Obviously, Manhattan distance satisfies the symmetry, non-negativity and triangle inequality postulates.

Consider, in Manhattan hashcode space, a point $P = [p_1, p_2, \dots, p_d]$ and a sphere $R = (O, r)$, where

$O = [o_1, o_2, \dots, o_d]$ is center of R and r is radius. The first metric is as follows:

Definition 2. Given P and $R = (O, r)$, the minimum Manhattan distance between P and R in Manhattan hashcode space, is defined as follows:

$$minDist(P, R) = \max(Dist(P, O) - r, 0) \quad (8)$$

The metric above has the following property:

Lemma 1. The distance $minDist(P, R)$ is no larger than the distance between P and any object enclosed in R .

Proof. Let S be any image object enclosed in R , and the hashcode of which is $[s_1, s_2, \dots, s_d]$. It is obvious that $\forall i \in [1, d]$, we have

$$|p_i - s_i| \geq |p_i - o_i| - |s_i - o_i|.$$

By Equation 7, we can obtain

$$Dist(P, S) \geq Dist(P, O) - Dist(S, O).$$

And $Dist(S, O) \leq r$, so

$$Dist(P, S) \geq Dist(P, O) - r.$$

Then, by Definition 2 and non-negativity of Manhattan distance, we have

$$minDist(P, R) \leq Dist(P, S).$$

□

The above lemma offers a lower bound on the Manhattan distance between a query and all image objects enclosed in the Sphere of a node, and it shall be used in our proofs later. The second metric we present is the maximum text relevancy between a query and a node, which takes the category into consideration.

Definition 3. Given a query Q and a node \mathcal{N} , the maximum text relevancy between Q and \mathcal{N} is defined as follows:

$$maxRel(Q_k, \mathcal{N}) = \max_{c \in \mathcal{N}.C} P(Q_k | (\mathcal{N}, c)) \quad (9)$$

$P(Q_k | (\mathcal{N}, c))$ is an upper bound on the text relevancy between Q_k and documents of all image objects that are contained in the subtree rooted at \mathcal{N} and belong to category c . It is computed as:

$$\begin{aligned} P(Q_k | (\mathcal{N}, c)) &= \prod_{t \in Q_k} p(t | (\mathcal{N}, c)) \\ &= \prod_{t \in Q_k} w((\mathcal{N}, c), t) \end{aligned}$$

where for $t \in (\mathcal{N}, c)$, $w((\mathcal{N}, c), t)$ is stored in the inverted file corresponding to the parent node of \mathcal{N} ; for $t \notin (\mathcal{N}, c)$, $w((\mathcal{N}, c), t) = \lambda \frac{tf(t, D_c)}{|D_c|}$.

Lemma 2. *The text relevancy $\maxRel(Q_k, \mathcal{N})$ is no less than the text relevancy between query Q and any image object contained in the subtree rooted at node \mathcal{N} .*

Proof. Let I be any image object belonging to the subtree rooted at \mathcal{N} , the text relevancy between q and I is:

$$P(Q_k|I_t) = \prod_{t \in Q_k} w(I_t, t).$$

According to Equation 6, we have

$$w((\mathcal{N}, I.C), t) \geq w(I_t, t).$$

By combining the above two results, we can obtain:

$$P(Q_k|(\mathcal{N}, I.C)) \geq P(Q_k|I_t).$$

Further, according to Equation 9 and the fact $I.C \in \mathcal{N}.C$, this yields

$$\maxRel(Q_k, \mathcal{N}) \geq P(Q_k|I_t).$$

□

The lemma above offers an upper bound on the text relevancy between a query and all image objects contained in the subtree rooted at a node, and it shall be used in the proof of Lemma 3. Furthermore, based on Equation 5 we propose the third metric, the maximum similarity score, which combines \minDist and \maxRel . Specifically, we have

Definition 4. *Given a query Q and a node \mathcal{N} , the maximum similarity score between Q and \mathcal{N} is defined as follows:*

$$\begin{aligned} \text{MaxScore}(Q, \mathcal{N}) = & \alpha \left(1 - \frac{\minDist(Q_v, \mathcal{N}.R)}{\max D} \right) + \\ & (1 - \alpha) \frac{\maxRel(Q_k, \mathcal{N})}{\max P} \end{aligned} \quad (10)$$

where α , $\max D$ and $\max P$ are same to that in Equation 5.

Lemma 3. *The similarity score $\text{MaxScore}(Q, \mathcal{N})$ is no less than the similarity score between any image object and query Q , which is contained in the subtree rooted at node \mathcal{N} .*

Proof. Let I be an image object belonging to the subtree rooted at node \mathcal{N} , the similarity score of I for query Q is $S(Q, I)$, and can be computed by Equation 5. By Lemma 1, we have

$$\minDist(Q_v, \mathcal{N}.R) \leq \text{Dist}(Q_v, I_v).$$

Further, by Lemma 2, we can obtain:

$$\maxRel(Q_k, \mathcal{N}) \geq P(Q_k|I_t).$$

By Equations 5 and 10, this yields

$$\text{MaxScore}(Q, \mathcal{N}) \geq S(Q, I).$$

□

For a query, the lemma above offers an upper bound on the similarity scores of all image objects contained in the subtree rooted at node \mathcal{N} . The last metric is the maximum Manhattan distance between a point and a Sphere.

Definition 5. *Given a point P and a sphere $R = (O, r)$, the maximum Manhattan distance between P and R in Manhattan hashcode space is defined as:*

$$\maxDist(P, R) = \text{Dist}(P, O) + r \quad (11)$$

Lemma 4. *The distance $\maxDist(P, R)$ is no smaller than the distance between P and any object enclosed in R .*

Proof. Let S be any object enclosed in R , the hashcode of which is $[s_1, s_2, \dots, s_d]$, it is obvious that for $\forall i \in [1, d]$, we have

$$|p_i - s_i| \leq |p_i - o_i| + |s_i - o_i|.$$

By Equation 7, we can obtain

$$\text{Dist}(P, S) \leq \text{Dist}(P, O) + \text{Dist}(S, O).$$

And $\text{Dist}(S, O) \leq r$, so

$$\text{Dist}(P, S) \leq \text{Dist}(P, O) + r.$$

By Definition 5, this yields

$$\maxDist(P, R) \geq \text{Dist}(P, S).$$

□

Besides, one can easily obtain the extra results below. Given a query Q and the root node, $root$, of MHIM-tree, one can see that, for any image object I in the tree, it is enclosed in $root.R$. Thus, by Lemma 4 and Lemma 2, one can obtain the following corollary:

Corollary 1. $\maxDist(Q_v, root.R) \geq \text{Dist}(Q_v, I_v)$.

Corollary 2. $\maxRel(Q_k, root) \geq P(Q_k|I_t)$.

The above two facts imply that, in the MHIM-tree we can essentially view $\maxDist(Q_v, root.R)$ as $\max D$, and $\maxRel(Q_k, root)$ as $\max P$, respectively.

6 The CATIRI Retrieval Algorithm

CATIRI comprises several phases: (1) text-and-visual feature preprocessing; (2) building MHIM-tree; (3) image retrieval via MHIM-tree.

6.1 Preprocessing

In this phase, we preprocess (i) the documents of textual annotations, and (ii) the visual features of image objects.

The text preprocessing is mainly for calculating the text relevancy. Specifically, for each document I_t , one can calculate, according to Equation 1, the weight $w(I_t, t)$ of each distinct term t . This way, one can obtain a set of term-weight pairs $(t, w(I_t, t))$. In addition, one may need to calculate $\lambda \frac{tf(t, D_c)}{|D_c|}$ for each distinct term t in the database D . This value represents the weight of t when term t is not in document I_t . After text processing, one can easily get the text relevancy according to Equation 2.

The visual feature preprocessing is for transforming the visual features of image objects into hashcodes. As mentioned earlier, we use Manhattan hashing algorithm [23] to get the hashcodes of visual features.

6.2 Building MHIM-tree

As mentioned in Section 4, our MHIM-tree is based on the classic M-tree. This allows us to adapt the construction algorithm of M-tree to achieve our goal. Generally, the construction of MHIM-tree is similar to that of M-tree. To save space, we only present the parts that are different from M-tree.

Algorithm 1: Insert($I, category$)

Input: Image object I and its category label $category$

- 1 $R \leftarrow \text{createSphere}(I_v, 0)$;
 - 2 $E \leftarrow \text{createEntry}(I.id, R, I_t, category)$;
 - 3 $\mathcal{N} \leftarrow \text{ChooseLeaf}(E)$;
 - 4 Add E to node \mathcal{N} , add I_t and $category$ to inverted file corresponding to \mathcal{N} ;
 - 5 $\text{AdjustTree}(\mathcal{N})$;
-

The MHIM-tree is built by **Insert** operation, whose pseudo-codes are described in Algorithm 1. In this algorithm, **ChooseLeaf** is responsible for choosing the best leaf node to place a new image object. This operation starts from the root, and recursively selects the best subtree to enclose the image object, until a leaf node is reached. Algorithm 2 covers more details about the **ChooseLeaf** operation. Note that, the selection criteria is to make the change for “inserted” subtree as small as possible. Here the change includes the *enlargement of the covering radius*, and also the *difference of the “document” vector*. The followings show how the

change is measured.

Let E be the entry of the image object to be inserted, E_i be one of entries in the current node. Without loss of generality, assume that one wants to insert E into the subtree rooted at E_i , then the enlargement of the covering radius is:

$$rDif(E_i, E) = \mathbf{CMB}(E_i.R, E.R).r - E_i.r \quad (12)$$

where $\mathbf{CMB}(E_i.R, E.R)$ is minimum bounding sphere enclosing $E_i.R$ and $E.R$, so its radius is:

$$\mathbf{CMB}(E_i.R, E.R).r = \max(\text{Dist}(E_i.O, E.O) + E.r, E_i.r). \quad (13)$$

Correspondingly, the change of the document vector is:

$$vDif(E_i, E) = 1 - \text{cosSim}(E_i.V, E.V) \quad (14)$$

where cosSim is the cosine similarity between two vectors. With the above two concepts in mind, the “overall” change is measured by the following:

$$sDif(E_i, E) = \beta \frac{rDif(E_i, E)}{\max R} + (1 - \beta)vDif(E_i, E) \quad (15)$$

where $\max R$ is the maximum $rDif$, which is used to normalize $rDif$ into $[0, 1]$. In addition, $\beta \in [0, 1]$ is a parameter used to balance visual similarity and text relevancy. In general, one can set $\beta = \alpha$ in order to maintain the consistency.

Algorithm 2: ChooseLeaf(E)

Input: An entry E representing an image object

Output: The leaf node that E should be inserted into

- 1 $\mathcal{N} \leftarrow \text{root}$;
 - 2 **while** \mathcal{N} is a non-leaf node **do**
 - 3 **for each** entry E_i in \mathcal{N} **do**
 - 4 \lfloor change $\leftarrow sDif(E_i, E)$;
 - 5 Select the entry E_i with the smallest change;
 - 6 $\mathcal{N} \leftarrow$ the node pointed by $E_i.p$;
 - 7 **return** \mathcal{N} ;
-

Another important operation in Algorithm 1 is the **AdjustTree** operation, which is responsible for adjusting the MHIM-tree after inserting a new image object. This operation starts from the leaf node, recursively adjusts the entry of each visited node and updates the inverted file in its parent node, until the root is reached. Algorithm 3 covers more details about this operation.

A key operation in Algorithm 3 is **SplitNode**, used to split the node into two new nodes when the number of entries in a node exceeds the *maximum limit*. The division criteria is to make more similar entries in the

same node, and less similar entries in different nodes. By incorporating the above criteria (based on $sDif$ in Equation 15) and the Split policy in [25], Algorithm 4 covers the details about **SplitNode**.

Algorithm 3: AdjustTree(\mathcal{N})

Input: A node \mathcal{N}

- 1 $\mathcal{N}_p \leftarrow$ parent node of \mathcal{N} ;
- 2 **if** \mathcal{N} needs to be split **then**
- 3 $\{\mathcal{N}_1, \mathcal{N}_2\} \leftarrow$ SplitNode(\mathcal{N});
- 4 **if** \mathcal{N} is root **then**
- 5 Create a new node \mathcal{N}' ;
- 6 Add $\mathcal{N}_1, \mathcal{N}_2$ into \mathcal{N}' as child nodes and update the inverted file corresponding to $\mathcal{N}_1, \mathcal{N}_2$ and \mathcal{N}' ;
- 7 root $\leftarrow \mathcal{N}'$;
- 8 **else**
- 9 Delete \mathcal{N} from the parent node \mathcal{N}_p , and add $\mathcal{N}_1, \mathcal{N}_2$ to \mathcal{N}_p as child nodes;
- 10 **else**
- 11 $E_{\mathcal{N}} \leftarrow$ entry of \mathcal{N} ;
- 12 adjustEntry($E_{\mathcal{N}}$);
- 13 Update the inverted file corresponding to \mathcal{N}_p ;
- 14 **if** \mathcal{N} is not root **then**
- 15 AdjustTree(\mathcal{N}_p);

Algorithm 4: SplitNode(\mathcal{N})

Input: A node \mathcal{N} that needs to be split
Output: Two new nodes obtained by splitting \mathcal{N}

- 1 $m \leftarrow$ the minimum number of entries in a node;
- 2 $\mathcal{N}_1, \mathcal{N}_2 \leftarrow$ empty node;
- 3 Insert E_1 in node \mathcal{N} into \mathcal{N}_1 ;
- 4 **for** each remaining entry E_i in node \mathcal{N} **do**
- 5 $rD \leftarrow Dist(E_i.O, \mathcal{N}.O)$;
- 6 Select the entry with the largest rD , and insert it into \mathcal{N}_2 ;
- 7 **while** there are unassigned entries in node \mathcal{N} **do**
- 8 **if** one node has to contain all the rest to reach m **then**
- 9 Insert all the rest of entries into it;
- 10 **break**;
- 11 **for** each remaining entry E_i **do**
- 12 $d_1 \leftarrow sDif(\mathcal{N}_1, E_i)$;
- 13 $d_2 \leftarrow sDif(\mathcal{N}_2, E_i)$;
- 14 Select the entry with the largest $|d_1 - d_2|$ and insert it into the node with smaller $sDif()$;
- 15 **return** $\{\mathcal{N}_1, \mathcal{N}_2\}$;

6.3 Retrieval via MHIM-tree

Our query algorithm follows the best-first traversal paradigm [33]. It fully exploits the proposed indexing structure as well as the metric properties. Generally speaking, we employ a priority for nodes and image objects to determine the order in which they are to be visited. For an image object, its priority is the similarity score computed by Equation 5, and for a node its priority is the $MaxScore$ computed by Equation 10. In the algorithm, we use the common rule: “the larger the score is, the higher the priority is”. In the retrieval process, our algorithm repeatedly chooses the next node or image object from the priority queue, which keeps track of unvisited nodes and image objects, until k image objects have been obtained. Algorithm 5 illustrates the details for top- k image retrieval.

Algorithm 5: Query($Q, root, k$)

Input: A query Q , the root node $root$, and the number of desired images k
Output: A list of desired images

- 1 $que \leftarrow$ createPriorityQueue();
- 2 $result \leftarrow$ createList();
- 3 $count \leftarrow 0$;
- 4 $que.Enqueue(root, 1)$;
- 5 **while** que is not empty **do**
- 6 $elem \leftarrow que.Dequeue()$;
- 7 **if** $elem$ is an image object **then**
- 8 $result.add(elem)$;
- 9 $count++$;
- 10 **if** $count = k$ **then**
- 11 **break**;
- 12 **else if** $elem$ is a leaf node **then**
- 13 **for** each entry (i.e., image object I) in $elem$ **do**
- 14 $priority \leftarrow S(Q, I)$;
- 15 $que.Enqueue(I, priority)$;
- 16 **else**
- 17 **for** each entry (i.e., node \mathcal{N}) in $elem$ **do**
- 18 $priority \leftarrow MaxScore(Q, \mathcal{N})$;
- 19 $que.Enqueue(\mathcal{N}, priority)$;
- 20 **return** $result$;

7 Theoretical Analysis

In this section, we cover main theoretical results related to our method. In analysis, we use m/M to denote the minimum/maximum number of entries in a node, and d to denote the dimension of hashcode.

Theorem 1. *Our top- k query processing algorithm is correct.*

Proof. To prove our query algorithm is correct, the key of point is to show that all image objects are dequeued in the descending order of their similarity scores. The followings validate this fact. Given a query Q , two image object I_1 and I_2 , without loss of generality, assume I_1 is dequeued before I_2 (in the query processing). It is easily know that, after I_1 is dequeued, there are two cases:

- (i) I_2 is in the queue;
- (ii) a node \mathcal{N} containing I_2 in the queue.

If it is the former case, then it is not hard to see that $S(Q, I_1) \geq S(Q, I_2)$, since I_1 has a higher priority than I_2 . In contrast, if it is the latter case, we can get $S(Q, I_1) \geq \text{MaxScore}(Q, \mathcal{N}) \geq S(Q, I_2)$, according to Lemma 3. Therefore, image objects are dequeued in the descending order of their similarity scores in Algorithm 5. This completes the proof. \square

Theorem 2. *The upper bound on time cost of constructing an MHIM-tree containing N image objects is*

$$O(MN \lceil \log_m N \rceil (d + |D_c|)).$$

Proof. The height of an MHIM-tree containing N image objects is at most $\lceil \log_m N \rceil - 1$, since the branch number of each node is between m and M . And the maximum number of nodes in the tree can be estimated as

$$\text{maxNode} = \lceil \frac{N}{m} \rceil + \lceil \frac{N}{m^2} \rceil + \dots + 1 \approx \frac{N-1}{m-1}.$$

To insert an image object, one can follow a path from the root node to the leaf node for **ChooseLeaf** operation. For each visited node in **ChooseLeaf** operation, we scan no more than M entries to find the most appropriate one, and the cost per entry is $O(d + |D_c|)$. Clearly, the cost of **ChooseLeaf** operation should be $O((\lceil \log_m N \rceil - 1) * M * (d + |D_c|))$.

In addition, we need to adjust the tree by **AdjustTree** operation, since the new node may incur some updates for the upper level nodes. The updates can be done by following the reversed path from the chosen leaf to the root. For each visited node in **AdjustTree** operation, we need to adjust the attributes and the inverted file in its parent node. Hence, the cost for **AdjustTree** operation is $O((\lceil \log_m N \rceil - 1) * (d + |D_c|))$.

Based on the above facts, it easily get that the cost for inserting an image object (without **SplitNode** operations) is $O((\lceil \log_m N \rceil - 1) * (M + 1) * (d + |D_c|))$.

On the other hand, it is not hard to see that the cost for each **SplitNode** operation is $O(\frac{M^2+M}{2} * (d + |D_c|))$, and the number of **SplitNode** operations is no more than $\text{maxNode} - 1$.

Thus, to sum up, the time cost of building an MHIM-tree is at most

$$\begin{aligned} O((M + 1) * N * (\lceil \log_m N \rceil + \frac{M}{2(m-1)}) * (d + |D_c|)) \\ \approx O(MN \lceil \log_m N \rceil (d + |D_c|)). \end{aligned}$$

\square

Theorem 3. *For an MHIM-tree containing N image objects with c (≥ 1) categories, the upper bound on space cost is $O(Nd + (\frac{c}{m} + 2)N|D_c|)$.*

Proof. Without loss of generality, assume that the number of all nodes is nodeSum . It is easily verified that in the MHIM-tree there are N entries that represent image objects, and nodeSum entries that represent nodes. The space cost per entry is $O(d + |D_c|)$, in which $O(d)$ is for R in the entry and $O(|D_c|)$ is for \mathbb{V} in the entry. So the space cost of for all entries is $O((\text{nodeSum} + N) * (d + |D_c|))$.

The other part of space cost is the inverted files corresponding to nodes. Note that, in the inverted lists, one term is at most related to one tuple for each entry that represents an image object, and is at most related to c tuples for each entry that represents a node. In addition, the total number of terms is no more than $|D_c|$. Thus, the total space cost of all inverted files is $O((\text{nodeSum} * c + N) * |D_c|)$.

To sum up, the space cost of an MHIM-tree is

$$O((\text{nodeSum} + N) * d + [\text{nodeSum} * (c + 1) + 2N] * |D_c|).$$

Note that, since the maximum number of nodes is $\frac{N-1}{m-1}$ (stated in the proof of Theorem 2), we can relax the above result, obtaining the following (by replacing nodeSum): $O(Nd + (\frac{c}{m} + 2)N|D_c|)$. This completes the proof. \square

Theorem 4. *For a top- k query Q on an MHIM-tree containing N image objects, the time complexity in the best-case is $O((k + \lceil \log_M \frac{N}{k} \rceil)(d + |Q_k|))$, and in the worst-case is $O(N(d + |Q_k|))$.*

Proof. The time cost for visiting a node or an image object is $O(d + |Q_k|)$, according to Equation 5 and 10. In the best-case, all image objects of the query result are close to each other in the MHIM-tree. So in the retrieval process, the number of visited nodes and image objects is at least

$$\begin{aligned} k + (\lceil \frac{k}{M} \rceil + \lceil \frac{k}{M^2} \rceil + \dots + 1) + (\lceil \log_M N \rceil - \lceil \log_M k \rceil) \\ = k + \frac{k-1}{M-1} + \lceil \log_M \frac{N}{k} \rceil \approx O(k + \lceil \log_M \frac{N}{k} \rceil). \end{aligned}$$

Thus, the time complexity is $O((k + \lceil \log_M \frac{N}{k} \rceil)(d + |Q_k|))$. In the worst-case, we may need to visit the

whole tree in the retrieval, so the number of visited nodes and image objects is at most

$$N + \max Node = N + \frac{N-1}{m-1} \approx O(N).$$

Thus, the time complexity is $O(N(d + |Q_k|))$. \square

8 Experiments

8.1 Experimental Settings

Datasets. In our experiments, we use three benchmark datasets: IAPR [34], LabelMe [35] and NUS-WIDE [36]. NUS-WIDE is a large social image dataset to test the scalability of our method, used to investigate if our method is promising for processing large-scale image retrieval. For ease of reference, Table 2 summarizes the main properties of these three datasets. In order to enhance the reliability and comparability of experiment results, for IAPR and LabelMe we scale all images to the same size, and extract 512-dimensional GIST descriptor as visual vector. For NUS-WIDE, we use 500-dimensional bag-of-words (included in the dataset) as visual vector.

Table 2. Properties of Dataset

Property	IAPR	LabelMe	NUS-WIDE
images	20000	73000	269648
distinct words	7873	19291	425000
words	348630	442215	4949317
words per image	1-55	1-317	1-632

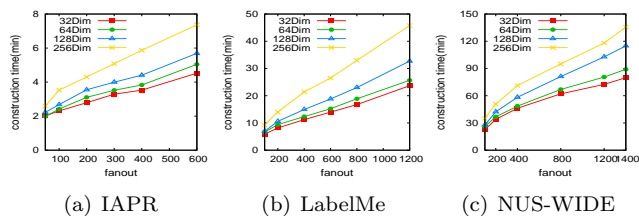
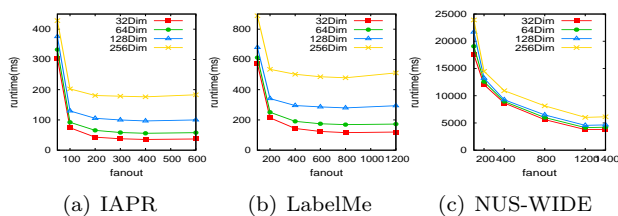
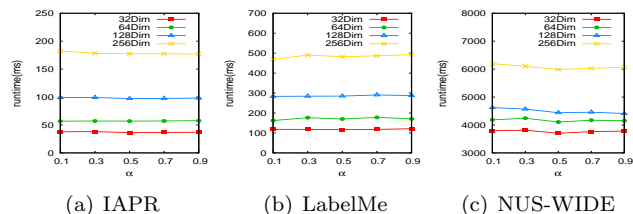
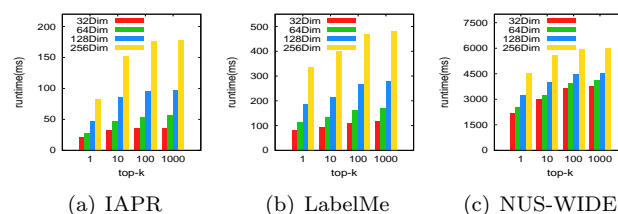
Evaluation metrics. In our experiments, four evaluation metrics are used: (i) the construction time of our MHIM-tree; (ii) the runtime of executing top- k image retrieval; (3) the I/O cost per query; and (iv) the mean average precision (MAP), which is to measure the retrieval quality. Note that, we adopt 60 query topics and the *ground truth* for top-1000 retrieval task used in ImageCLEF2007 [30] to evaluate the accuracy for the IAPR dataset. For LabelMe and NUS-WIDE, there is no ground truth for top- k image retrieval, and so we randomly choose 1000 image objects (with image and textual caption) to form the query set, and only evaluate the running efficiency.

Parameter settings. Following [29], we set $\lambda = 0.2$ in our experiments. In addition, in our experiments the parameter k is set to [1, 10, 100, 1000], where 1000 is the default value. To evaluate the influence of the dimension in the hashing process, we set $d = [32, 64, 128, 256]$, in which $d = 128$ is the default setting. To investigate the impact of the fanout (i.e., the balance factor B) of our MHIM-tree, we set $B = [50, 100, 200, 300, 400, 600]$

for IAPR dataset, in which $B = 400$ is the default setting. For the LabelMe dataset, the values we used are [100, 200, 400, 600, 800, 1200], and $B = 800$ is set to the default value. For the NUS-WIDE dataset, we use $B = 1200$ as the default value. A major reason we use different B for these datasets is that, the sizes of these datasets are different, a more micromesh settings allow us to find exactly the impact trend and help us find appropriate default values for other tests. Also, we study the impact of the parameter α , which is used to adjust the weigh between textual and visual parts when calculating the similarity score. Specifically, we set $\alpha = [0.1, 0.3, 0.5, 0.7, 0.9]$, and the default value is 0.5. Experiments were executed on an Intel(R) Core(TM) i5-5200U CPU @2.20HZ and 4GB RAM.

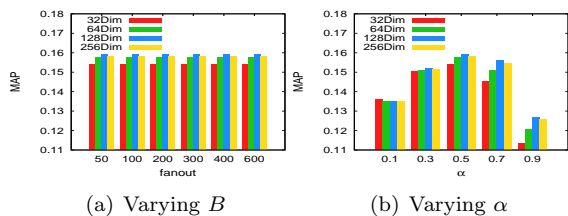
Compared method. For our problem, it is challenging to find an appropriate baseline for comparison, although there are already some representative works [20, 4, 19, 21, 22, 37] that investigate the content-and-text based image retrieval (CTBIR). Among these works, most of them are devoted to designing better fusion techniques (e.g., [22, 19, 21, 20]). That said, they focused on how to *effectively* combine text and content information to improve the retrieval quality. Our paper directly uses their fusion technique as mentioned in Section 3. Note that, for their papers it is unclear on how to perform top- k image retrieval. Their works are essentially orthogonal to our work. Besides, the work in [4] focused on developing advanced *hashing techniques* to improve the retrieval quality; for their paper it is also unclear on how to perform top- k image retrieval. Our paper directly uses the existing hashing technique. Note that, as we discussed in Section 6.1, designing more effective hashing techniques is not the focus of this paper. Thus, our work is also orthogonal to [4]. Among all these works, the paper [37] could be closest to our work, since they “mentioned” the *retrieval efficiency* by directly exploiting an existing technique “inverted index”. Thus, we would like to choose their method as the Baseline1. Note that, their paper uses visual words to compute the similarity score, it is different from our similarity score measurement; it could incur some deviations in evaluation. For the fair comparison, we adapt their method by using also the inverted file and the same similarity measurement, getting a version similar to their method. The general steps of the baseline are briefly described as follows.

It first creates a inverted file for all image objects, which records $p(t|I_t)$. For a top- k query, it uses the inverted file to obtain the text relevancy between the query and all image objects by Equation 2, and then sorts all image objects in the descending order to form an array. After that, it traverses the array, and calculates the similarity score for each visited image ob-

Fig.2. Construction time vs. B Fig.3. Runtime vs. B Fig.4. Runtime vs. α Fig.5. Runtime vs. k

ject. During the traversal, it keeps a priority queue to record k objects with highest similarity scores. When visiting an image object I whose text relevancy is lower than a threshold τ , it stops traversing and returns all image objects in the queue as the query result, where $\tau = \frac{S_{min}-\alpha}{1-\alpha}$, S_{min} refers to the minimum similarity score in the priority queue.

Furthermore, we adopt Manhattan hashing to optimize Baseline1 according to paper [23], and then get Baseline2 method. we will compare our method with Baseline1 and Baseline2, to enhance the comparability and reliability of our experiment.

Fig.6. MAP vs. B and α

8.2 Experimental Results of CATIRI

Indexing Construction Cost. Figure 2 shows the construction time of MHIM-tree, it is easy to find that the construction time is almost proportional to the fanout for all datasets. In addition, the higher the dimension is, the longer the construction time is. This is consistent to our theoretical analysis in Section 7.

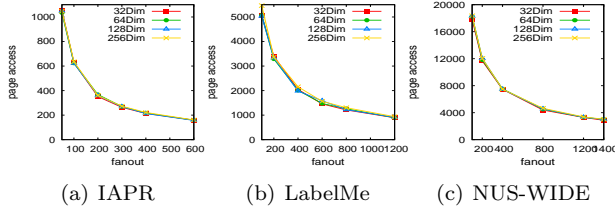
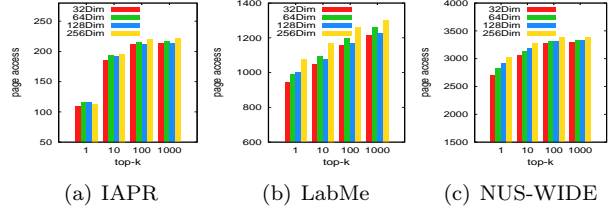
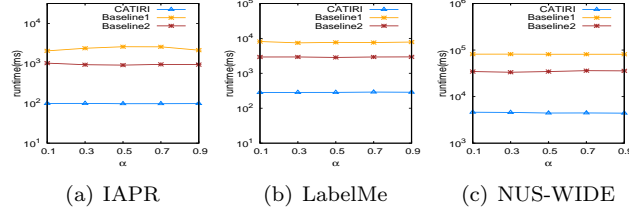
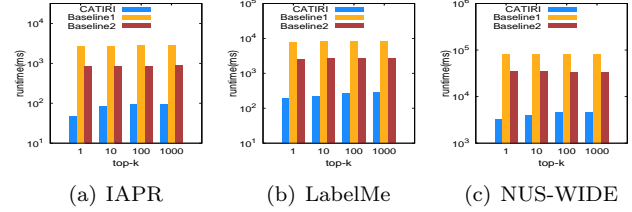
Query Time. Figure 3 shows the query performance of the top- k image retrieval with different fanout. It can be seen that the runtime is almost proportional to the dimension d . The smaller d is, the less the query time is. Besides, one can see that, when the fanout B is very small, the runtime is long and decreases rapidly when

B increases. This implies that a too small fanout value could hurt the performance. When the fanout reaches a certain value, the runtime performance enters a steady stage. Nevertheless, when the fanout increases to a too large value, the performance starts to degrade. This implies that a too larger fanout value could be not very helpful. This is because, by limiting the capacity of each node, the fanout directly determines the height of MHIM-tree. Clearly, when the fanout is too small, the index tree is so high that the retrieval needs a long time to travel. On the other hand, when the fanout is too large, the height of the MHIM-tree is so small that it offers little help for the retrieval. This essentially shows us that choosing an appropriate fanout is important.

Figure 4 shows the runtime when we vary α . It can be seen that there is no obvious fluctuation in terms of the runtime. This implies that α has little effect on the runtime. On the other hand, as we expected, the runtime is still almost proportional to the dimension d .

The query runtime performance for different k is shown in Figure 5. On one hand, the runtime increases when k is relatively small, while the growth speed turns slow when k reaches a large enough value. The reason is same to our analysis for I/O performance. That said, the pruning ability of our MHIM-tree is powerful and scalable. On the other hand, one can see that dimension d can make impact on the performance, and usually the larger the dimension is, the larger the runtime is. This is because computing the distances in the retrieval needs to consider each dimension.

I/O Cost. Figure 7 shows the I/O cost per query for different B . One can see that the I/O cost decreases when the fanout increases. This is mainly because both the number of nodes and the height of the MHIM-tree decrease when the fanout B increases. In addition, as we expected, the dimension d rarely influences the I/O

Fig.7. I/O cost vs. B Fig.8. I/O cost vs. k Fig.9. Runtime vs. α Fig.10. Runtime vs. k

cost.

Figure 8 shows the impact of k on the I/O cost. We can see that when k is large enough, the I/O cost will not increase when k increases. The main reason is that, the pruning power of the MHIM-tree plays a greater role in the retrieval when k is large. In addition, we can also see that, for different dimensions, the I/O cost is almost equivalent. This result is consistent to our previous discussion.

Accuracy. Furthermore, we also study the impact of B on the retrieval quality. From Figure 6(a) we can see that, the fanout has no influence on the MAP. This further verifies that the fanout B is essentially irrelevant with the similarity score. The MAP values for different α are shown in Figure 6(b). We can see that the MAP first ascends and then descends with the increase of α , and reaches the peak at $\alpha = 0.5$. These results validate that the fusion of visual and text information can essentially improve the accuracy of query, and the best performance can be achieved at the balance point of the visual and text information. Another phenomenon is that 128 dimensions achieve the best MAP for all cases except $\alpha = 0.1$. The MAP keeps nearly unchanged for $\alpha = 0.1$, because α is so small that the visual similarity makes almost no contribution on the query result. In addition, we can see that the most appropriate dimension for our experiments could be 128.

8.3 Comparing with Baseline

Query quality. Table 3 shows the MAP values for our method and for the baselines. It can be seen that the MAP of our method is close to that of the baselines. This means that our method can successfully preserve the retrieval accuracy. The major reason is that re-

trieval accuracy mainly depends on the fusion technique and hashing method, CATIRI directly uses existing fusion technique and hashing technique, obtaining similar retrieval quality.

Table 3. MAP on IAPR

α	0.1	0.3	0.5	0.7	0.9
CATIRI	0.1351	0.1517	0.1589	0.1561	0.1265
Baseline1	0.1321	0.1398	0.1499	0.1558	0.1159
Baseline2	0.1374	0.1533	0.1559	0.1526	0.1264

Runtime. Figure 9 and 10 compare the runtime time of our method against the baselines, it is easy to find that our MHIM-tree significantly outperforms the baseline algorithms under all settings of α and k . And our method almost outperforms both baselines by an order of magnitude. This essentially demonstrates the effectiveness and efficiency of our proposed method.

Construction Cost. Table 4 shows the index construction cost for our method and baselines. It can be seen that the construction time of our method is larger than that of the baselines, but it is within the acceptable range. Considering the huge improvement in retrieval efficiency of our method, it is worth sacrificing a little construction time to get better retrieval efficiency.

Table 4. Construction Time (minute) on Datasets

Dataset	IAPR	LabelMe	NUS-WIDE
CATIRI	4.41	23.00	102.77
Baseline1	0.42	0.70	5.62
Baseline2	1.62	6.00	25.25

9 Conclusion

In this paper we proposed *CATIRI* for the content-and-text based image retrieval. *CATIRI* consists of three phases. It first transforms the image query problem into the one similar to spatial keyword query by feature hashing, and then builds an MHIM-tree to manage the textual and visual information, and finally performs the retrieval based on our top-*k* query algorithm that fully exploits the MHIM-tree and a set of important properties. We presented theoretical analysis for *CATIRI*, and experimentally demonstrated that it can remarkably improve retrieval efficiency while preserving the quality of retrieval results.

Acknowledgement(s) This work was supported by the National Basic Research 973 Program of China under Grant No.2015CB352403, the National Key Research and Development Program of China under Grant Nos.2018YFC1504504, 2016YFB0700502 and 2018YFB1004400, the National Natural Science Foundation of China under Grant Nos.61872235, 61729202, 61832017, U1636210, 61832013, 61672351 and 61472453, 61702320, U1401256, U1501252, U1611264, U1711261, U1711262, U61811264, and Guangdong Province Key Laboratory of Popular High Performance Computers of Shenzhen University under Grant No, SZU-GDPHPCL2017.

References

- [1] Wu L, Jin R, Jain A K. Tag Completion for Image Retrieval. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 2013, 35(3): 716-727.
- [2] Gong Y, Lazebnik S, Gordo A, Perronnin F. Iterative quantization: A procrustean approach to learning binary codes. In *Proc. Computer Vision and Pattern Recognition*, Jun. 2011, pp.817-824.
- [3] Datta R, Joshi D, Li J, Wang J Z. Image retrieval: Ideas, influences, and trends of the new age. *Acm Computing Surveys*, 2008, 40(2): 1-60.
- [4] Long M, Cao Y, Wang J, Yu P S. Composite Correlation Quantization for Efficient Multimodal Retrieval. In *Proc. the 39th Int. ACM SIGIR Conf. on Research and Development in Information Retrieval*, Jul. 2016, pp.579-588.
- [5] Zhu L, Shen J, Xie L, Cheng Z. Unsupervised visual hashing with semantic assistant for content-based image retrieval. *IEEE Trans. on Knowledge and Data Engineering*, 2017, 29(2): 472-486.
- [6] Xu B, Bu J, Chen C, Cai D, He X. EMR: A scalable graph-based ranking model for content-based image retrieval. *IEEE Trans. on Knowledge and Data Engineering*, 2015, 27(1): 102-114.
- [7] Shen H T, Jiang S, Tan K L, Huang Z, Zhou X. Speed up interactive image retrieval. *The Int. Journal on Very Large Data Bases*, 2009, 18(1): 329-343.
- [8] Falchi F, Lucchese C, Orlando S, Perego R, Rabitti F. Caching content-based queries for robust and efficient image retrieval. In *Proc. the 12th Int. Conf. on Extending Database Technology: Advances in Database Technology*, Mar. 2009, pp.780-790.
- [9] Zhang C, Chai J Y, Jin R. User term feedback in interactive text-based image retrieval. In *Proc. the 28th Int. ACM SIGIR Conf. on Research and Development in Information Retrieval*, Aug. 2005, pp.51-58.
- [10] Li W, Duan L, Xu D, Tsang I W. Text-based image retrieval using progressive multi-instance learning. In *Proc. Int. Conf. on Computer Vision*, Nov. 2011, pp.2049-2055.
- [11] Tong S, Chang E. Support vector machine active learning for image retrieval. In *Proc. the 9th ACM Int. Conf. on Multimedia*, Sep. 2001, pp.107-118.
- [12] Liu D, Hua K A, Vu K. Fast query point movement techniques with relevance feedback for content-based image retrieval. In *Proc. Int. Conf. on Extending Database Technology*, Aug. 2006, pp.700-717.
- [13] Kulis B, Grauman K. Kernelized locality-sensitive hashing for scalable image search. In *Proc. IEEE Int. Conf. on Computer Vision*, Sep. 2009, pp.2130-2137.
- [14] Smeulders A W M, Worring M, Santini S, Gupta A, Jain R C. Content-based image retrieval at the end of the early years. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 2000, 22(12): 1349-1380.
- [15] Deng J, Berg A C, Li F F. Hierarchical semantic indexing for large scale image retrieval. In *Proc. Computer Vision and Pattern Recognition*, Jun. 2011, pp.785-792.
- [16] Ooi B C, Tan K L, Chua T S, Hsu W. Fast image retrieval using color-spatial information. *The Int. Journal on Very Large Data Bases*, 1998, 7(2): 115-128.
- [17] Xia H, Wu P, Hoi S C H, Jin R. Boosting multi-kernel locality-sensitive hashing for scalable image retrieval. In *Proc. Int. ACM SIGIR Conf. on Research and Development in Information Retrieval*, Aug. 2012, pp.55-64.
- [18] Christel M G. Examining user interactions with video retrieval systems. In *Proc. Multimedia Content Access: Algorithms and Systems*, Oct. 2007, pp.1-15.
- [19] Zhou X S, Huang T S. Unifying keywords and visual contents in image retrieval. *IEEE Multimedia*, 2002, 9(2): 23-33.
- [20] Zagoris K, Chatzichristofis S A, Arampatzis A. Bag-of-visual-words vs global image descriptors on two-stage multimodal retrieval. In *Proc. the 34th Int. ACM SIGIR Conf. on Research and Development in Information Retrieval*, Dec. 2011, pp.1251-1252.
- [21] Caicedo J C, Moreno J G, Niuo E A, Gonzalez F A. Combining visual features and text data for medical image retrieval using latent semantic kernels. In *Proc. the Int. Conf. on Multimedia Information Retrieval*, Mar. 2010, pp.359-366.
- [22] Clinchant S, Ah-Pine J, Csurka G. Semantic combination of textual and visual information in multimedia retrieval. In *Proc. the 1st ACM Int. Conf. on Multimedia Retrieval*, Oct. 2011, pp.44.
- [23] Kong W, Li W J, Guo M. Manhattan hashing for large-scale image retrieval. In *Proc. the 35th Int. ACM SIGIR Conf. on Research and Development in Information Retrieval*, Aug. 2012, pp.45-54.
- [24] Zobel J, Moffat A. Inverted files for text search engines. *ACM Computing Surveys*, 2006, 38(2): 6.
- [25] Ciaccia P, Patella M, Zezula P. M-tree: An efficient access method for similarity search in metric spaces. In *Proc. 23th Int. Conf. on Very Large Data Bases*, Aug. 1997, pp.357-368.
- [26] Oliva A, Torralba A. Modeling the shape of the scene: A holistic representation of the spatial envelope. *Int. Journal of Computer Vision*, 2001, 42(3): 145-175.
- [27] Sivic J, Zisserman A. Video Google: A text retrieval approach to object matching in videos. In *Proc. IEEE Int. Conf. on Computer Vision*, Oct. 2003, pp.1470-1477.
- [28] Ponte J M, Croft W B. A language modeling approach to information retrieval. In *Proc. the 21th Int. ACM SIGIR Conf. on Research and Development in Information Retrieval*, Aug. 1998, pp.275-281.
- [29] Zhai C, Lafferty J. A study of smoothing methods for language models applied to information retrieval. *ACM Trans. on Information Systems*, 2004, 22(2): 179-214.
- [30] Depeursinge A, Muller H. Fusion techniques for combining textual and visual information retrieval. In *ImageCLEF*, 2010, pp.95-114.
- [31] Wang J, Liu W, Kumar S, Chang S. Learning to hash for indexing big data: a survey. *Proceedings of the IEEE*, 2016, 104(1): 34-57.
- [32] Cao X, Chen L, Cong G, Jensen C S, Qu Q, Skovsgaard A. Spatial keyword querying. In *Proc. Int. Conf. on Conceptual Modeling*, Nov. 2012, pp.16-29.

- [33] Hjaltason G R, Samet H. Distance browsing in spatial databases. *ACM Trans. on Database Systems*, 1999, 24(2): 265-318.
- [34] Grubinger M, Clough P, Muller H, Deselaers T. The iapr tc-12 benchmark: A new evaluation resource for visual information systems. In *Proc. Int. Workshop OntoImage*, May. 2006, pp.5.
- [35] Russell B C, Torralba A, Murphy K P, Freeman W T. LabelMe: a database and web-based tool for image annotation. *Int. Journal of Computer Vision*, 2008, 77(1-3): 157-173.
- [36] Chua T S, Tang J, Hong R, Li H, Luo Z, Zheng T Y. NUS-WIDE: a real-world web image database from National University of Singapore. In *Proc. the ACM Int. Conf. on image and video retrieval*, Oct. 2009, pp.1-9.
- [37] Mamou J, Mass Y, Shmueli-Scheuer M, Sznajder B. A unified inverted index for an efficient image and text retrieval. In *Proc. the 32th Int. ACM SIGIR Conf. on Research and Development in Information Retrieval*, Jul. 2009, pp.814-815.
- [38] Rasiwasia N, Costa Pereira J, Coviello E, Doyle G, Lanckriet G R G, Levy R, Vasconcelos N. A new approach to cross-modal multimedia retrieval. In *Proc. the 18th ACM Int. Conf. on Multimedia*, Oct. 2010, pp.251-260.
- [39] Yang C, Lozano-Perez T. Image database retrieval with multiple-instance learning techniques. In *Proc. Int. Conf. on Data Engineering*, Feb. 2000, pp.233-243.
- [40] Natsev A, Rastogi R, Shim K. WALRUS: A similarity retrieval algorithm for image databases. In *Proc. ACM SIGMOD Record*, Jun. 1999, pp.395-406.
- [41] Rabitti F, Savino P. An information retrieval approach for image databases. In *Proc. 18th Int. Conf. on Very Large Data Bases*, Aug. 1992, pp.574-584.
- [42] Chu W W, Jeong I T, Taira R K. A semantic modeling approach for image retrieval by content. *The VLDB Journal*, 1994, 3(4): 445-477.
- [43] Brown L, Gruenwald L. A Prototype Content-Based Retrieval System that Uses Virtual Images to Save Space. In *Proc. 27th Int. Conf. on Very Large Data Bases*, Sep. 2001, pp.693-694.
- [44] Chen L, Gao Y, Xing Z, Jensen C S, Chen G. I2RS: a distributed geo-textual image retrieval and recommendation system. *Proc. the Very Large Data Bases Endowment*, 2015, 8(12): 1884-1887.