

# Probabilistic Cost Model for Nearest Neighbor Search in Image Retrieval

Kunho Kim<sup>a</sup>, Mohammad K. Hasan<sup>a</sup>, Jae-Pil Heo<sup>a</sup>, Yu-Wing Tai<sup>a</sup>, Sung-eui Yoon<sup>a,b</sup>

<sup>a</sup>*Dept. of Computer Science, KAIST*

<sup>b</sup>*Div. of Web Science and Technology, KAIST*

---

## Abstract

We present a probabilistic cost model to analyze the performance of the kd-tree for nearest neighbor search in the context of content-based image retrieval. Our cost model measures the expected number of kd-tree nodes traversed during the search query. We show that our cost model has high correlations with both the observed number of traversed nodes and the runtime performance of search queries used in image retrieval. Furthermore, we prove that, if the query points follow the distribution of data used to construct the kd-trees, the median-based partitioning method as well as PCA-based partitioning technique can produce near-optimal kd-trees in terms of minimizing our cost model. The probabilistic cost model is validated through experiments in SIFT-based image retrieval.

*Keywords:* Image retrieval, nearest neighbor search, kd-tree

---

## 1. Introduction

The nearest neighbor search [1] is one of the most widely used proximity queries with various applications such as image retrieval [2], pattern recognition [3], etc. In this paper we focus on exact or approximate nearest neighbor search queries used for content-based image retrieval (CBIR). In CBIR, various image features such as SIFT [4] are used for finding similar images in image databases. In order to find the similar images reliably, multiple features (e.g., around 1 K features) are typically extracted from each query image. The problem of CBIR then becomes finding the nearest neighbor image such that its image features match closely with those from the query image.

Since the number of images in an image database can be very large, a kd-tree [5] or kd-tree forest [6] has been widely used to accelerate the performance of nearest neighbor search. In CBIR, however, the dimension of image features is very high. For example, the dimension of a SIFT feature vector for a single feature point in an image is 128. Thus, finding the exact nearest neighbor in the image database can be very slow. A common practice in CBIR is to find an approximate nearest neighbor instead of the exact nearest neighbor, by limiting the number of nodes traversed during the traversal of kd-trees [7].

In order to construct the kd-trees for CBIR, we need to partition the image feature points contained in a node into its two child nodes. This is performed with a hyperplane, which is defined by its normal and position

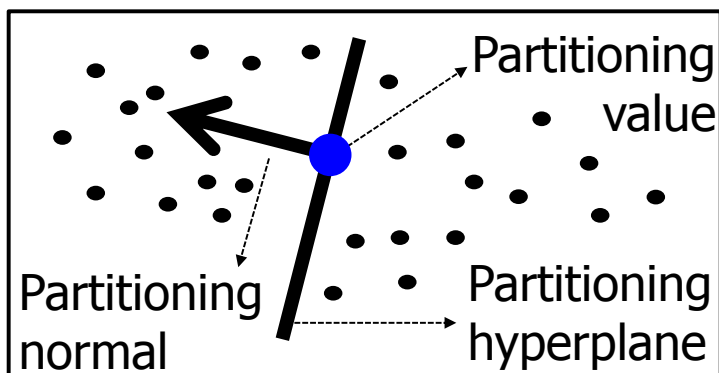


Figure 1: 2D example of a partitioning hyperplane.

24 (Fig. 1). The normal and position of the hyperplane can be called as the  
 25 partitioning normal and the partitioning value respectively.

26 Optimizing kd-trees has been actively studied in recent years [8, 9], since  
 27 constructing a high-quality kd-tree is critical for improving the performance  
 28 of many CBIR techniques. These prior techniques improve the performance  
 29 of nearest neighbor search by computing a partitioning normal that separates  
 30 feature points well and reduces computation overheads of accessing nodes at  
 31 the query time. In order to split the feature points in a node, the mean or  
 32 the median point along the chosen partitioning normals are typically selected  
 33 as the partitioning values. However, to the best of our knowledge, there has  
 34 been no prior work that quantifies the quality of kd-trees in a cost model for  
 35 CBIR and optimizes partitioning values and normals with the cost model.

36 **Contributions:** In order to quantify the quality of kd-trees and  
 37 hence to evaluate the performance of the nearest neighbor search in CBIR,

38 we propose a probabilistic cost model (Sec. 4) that measures the expected  
39 number of nodes traversed during the process of the search query. We have  
40 also proved that the conventional wisdom of partitioning data points in each  
41 node of kd-trees in the median point along a partitioning normal that is  
42 close to the principal directions of the data points can produce near-optimal  
43 kd-trees, given our proposed cost model (Sec. 5). Our cost model and as-  
44 sumptions are validated through experiments in SIFT-based image retrieval,  
45 which show strong linear correlations (i.e. up to 0.9) between our cost model  
46 and both the observed number of nodes traversed and the running time spent  
47 for queries used in CBIR applications (Sec. 6). These results can be served as  
48 a theoretic basis for other optimization problems such as updating kd-trees,  
49 to achieve a high performance for dynamic data sets.

## 50 **2. Related Work**

51 In this section we review prior techniques on content-based image retrieval  
52 (CBIR) and optimizing kd-trees. We discuss examples most relevant to our  
53 works and refer reader to [2] for the survey of recent approaches in CBIR.

### 54 *2.1. Image Retrieval*

55 Given a query image, CBIR compares the image features of images from  
56 an image database with the features from the query image, and returns the  
57 images that are most similar to the query image. Such CBIR problem can  
58 be reduced to the well-known nearest neighbor problem. Data structures

59 such as kd-tree [5] are widely used to efficiently process the nearest neighbor  
60 problem, which performs the matching between query features and image  
61 features stored in the image database [2].

62 Unique in CBIR, the dimension of image features are usually very high.  
63 Using a kd-tree for accessing such high-dimensional data is well-known to be  
64 very slow, since the nearest neighbor computation algorithms have to tra-  
65 verse most nodes of a kd-tree. Hence, approximate nearest neighbor search  
66 algorithms [7] that use a priority queue and/or the best-bin-first search  
67 technique [4] have been commonly adopted in CBIR. Recently, Muja and  
68 Lowe [10] presented an automatic tuning algorithm for various kd-tree based  
69 approximate nearest neighbor search methods.

70 Visual words [11] (or bag-of-features), vocabulary tree [12], and locality  
71 sensitive hashing [13] are some of the popular techniques for large CBIR sys-  
72 tems. A recent study by Philbin et al. [6] showed that performing nearest  
73 neighbor computation with kd-trees outperforms the technique using the vo-  
74 cabulary tree in terms of K-means clustering. It has also been shown that  
75 quality of searching results of locality sensitive hashing can be inferior to  
76 those of kd-trees [8]. While these previous works are related, we focus our  
77 analysis on evaluating kd-trees for high quality nearest neighbor computa-  
78 tion. Our results will be useful to evaluate the performance of kd-trees.

79 *2.2. Optimizing kd-Trees*

80 A few techniques have been proposed to improve the quality of kd-trees  
81 for the acceleration of tree access in CBIR. Silpa-anan et al. [8] proposed  
82 to create multiple kd-trees for an input data set and use them together as a  
83 concurrent search with a pooled priority queue. They have also demonstrated  
84 that a higher performance can be achieved by aligning the principal axis of  
85 data with the partitioning normal of a partitioning hyperplane used for data.  
86 Recently, Jia et al. [9] proposed a binary combination of axis-aligned axes for  
87 partitioning normals for kd-trees, instead of using arbitrary axes.

88 Once a partitioning normal is decided as proposed in above techniques, a  
89 mean value or a median value along the chosen partitioning normal is used to  
90 define a partitioning hyperplane. While this is a common practice, there have  
91 been no prior techniques that pay attention to optimizing the partitioning  
92 value along a chosen partitioning normal.

93 **Optimized kd-trees for ray tracing:** kd-trees have been widely used  
94 in other applications, especially ray tracing in the field of computer graphics.  
95 It has been widely known that the quality of kd-trees plays one of the key  
96 factors that govern the performance of ray tracing and, thus, optimizing  
97 kd-trees for the application has been actively studied. For kd-trees used  
98 in ray tracing, the Surface Area Heuristic (SAH) metric that measures the  
99 traversal cost during the kd-tree traversal is proposed [14]. Also, various kd-  
100 tree construction techniques [15, 16, 17] that optimize the metric have been  
101 presented. However, the usage of kd-trees for ray tracing is different from

102 that for image retrieval; the structures and ordering of polygons used for  
103 kd-trees in ray tracing are very different from the one in image retrieval. As  
104 a result, it is unclear how the metric proposed for ray tracing can be applied  
105 to image retrieval.

### 106 *2.3. Cost Models for kd-Trees*

107 The database and computational theory communities have been devel-  
108 oping various cost models for proximity queries including nearest neighbor  
109 search accelerated by spatial data structures such as kd-trees or R-trees.  
110 Bentley [18] introduced the concept of kd-trees and intuitively suggested  
111 that the median-based partitioning scheme for the kd-tree construction leads  
112 to the optimal tree in terms of a path length computed from the root node  
113 to a leaf node. This claim was an intuitive generalization of the optimality  
114 of one dimensional binary trees [19] to higher dimensional trees.

115 Friedman et al. [5] proposed a cost model for processing nearest neighbor  
116 search. Their cost model estimates the number of leaf nodes accessed in the  
117 kd-tree, while processing nearest neighbor queries. Extensions from the cost  
118 model [5] have been proposed for other spatial data structures like R-tree [20]  
119 and non-rectangular regional nodes [21]. However, these cost models have  
120 an unrealistic assumption that the number of data points is assumed to be  
121 infinite [22]. Lee and Wong [23] gave an worst-case analysis of balanced kd-  
122 trees for region queries in terms of the total number of node visits as a cost.  
123 Their worst-case analysis indicates that the dominant factor in the cost is the

124 dimensionality of data points. Sproull [24] later pointed out that these cost  
125 models are valid, when the number of points is exponential increasing as a  
126 function of the dimensionality of data points, in order to dilute the boundary  
127 effect. Arya et al. [25] developed an improved cost model that considers  
128 the boundary effect, but has an unrealistic assumption about the number of  
129 points as Sproull [24] pointed out. Berchtold et al. [26] proposed a cost metric  
130 mainly for high dimensional data, while the prior model of Friedman et al. [5]  
131 overestimates the cost by orders of magnitude for uniformly distributed high-  
132 dimensional data. Unfortunately, this method relied on the expensive Monte  
133 Carlo integration for evaluating the cost model, because the cost model is too  
134 complex to be calculated directly. Furthermore, it requires a huge number  
135 of samples for the Monte Carlo integration to achieve a reasonably accurate  
136 approximation. As a result, this model has been applied to a small scale of  
137 kd-trees.

138 These various cost models have been designed mainly for range queries  
139 and investigated in theoretical contexts. As a result, it is unclear how well  
140 these cost models are applicable to our domain, nearest neighbor search for  
141 high-dimensional image descriptors used in CBIR. Moreover, cost models  
142 proposed in theoretical communities makes assumptions inappropriate for  
143 CBIR. Some of them include that data points are uniformly distributed.



### 144 3. Access Patterns on kd-Trees and Terminologies

145 In this section we describe how a kd-tree is accessed given a query image  
146 for image retrieval. We will also define the terminologies used for the rest  
147 of the paper. We assume that the kd-tree is constructed based on image  
148 features (e.g. SIFT), and the same type of image features are extracted from  
149 the query image for searching nearest neighbor images.

150 Given an image feature,  $q_i$ , of a query image, we start to traverse the  
151 kd-tree from the root node of the kd-tree. During the kd-tree traversal, we  
152 maintain two variables: 1) the current minimum distance,  $min_d$ , and 2) the  
153 current candidate for the nearest neighbor feature that has the current min-  
154 imum distance to the query image feature  $q_i$ . The initial minimum distance  
155 is set as infinity before traversal. Since there are many methods to traverse a  
156 kd-tree, we analyze two of the most common methods, the depth-first traver-  
157 sal [1, p.517] and the best-bin-first search [27].

158 **Depth-first traversal (DFT):** The DFT is the simplest traversal  
159 method for the nearest neighbor search. In this scheme, once we visit an  
160 intermediate node, we compare the features in its left and its right child  
161 nodes to the given image feature  $q_i$ . We traverse the child node that is closer  
162 to the given image feature and store the other child node in a stack, called  
163 *traversal stack*. Once we visit a leaf node, we measure the distance between  
164  $q_i$  and each image feature stored in the leaf node. We pick the image feature,  
165  $q_k$ , stored in the leaf node that gives the shortest distance to  $q_i$ , and check  
166 whether the shortest distance is smaller than the current minimum distance

167  $min_d$ . If so, we update both the current candidate for the nearest neighbor  
168 as  $q_k$  and the current minimum distance as the distance between  $q_i$  and  $q_k$ .  
169 Then we pop a node from the traversal stack and access the node. This  
170 operation is commonly known as *back-tracking*. The process is performed  
171 recursively until the traversal stack becomes empty.

172 **Culling techniques:** The DFT is simple, but has to access all the  
173 nodes of the kd-tree, leading to a slow performance. A simple remedy to  
174 this problem is to employ various culling techniques. The most common  
175 culling method is to utilize a conservative distance bound; let us call this  
176 culling method *conservative distance culling*. For example, we can compute  
177 a conservative distance bound between the query point and the bounding box  
178 of a node of a kd-tree. If the conservative bound is bigger than the current  
179 minimum distance  $min_d$ , we cull the traversal operations on the child nodes  
180 located in the sub-tree of the node.

181 **Best-bin-first search (BBFS):** The BBFS has been known to show  
182 a higher performance than the DFT. The BBFS attempts to access nodes  
183 that are likely to have nearest neighbor points earlier than other nodes. To  
184 traverse kd-trees, the BBFS uses a priority queue instead of a stack. For each  
185 intermediate node,  $n_i$ , the BBFS computes a conservative distance bound  
186 between the query's image feature  $q_i$  and two child nodes of  $n_i$ . We traverse  
187 the child node that gives the smaller distance bound and store the other  
188 node into the priority queue with its conservative distance bound. Once we  
189 reach a leaf node, we perform the back-tracking operation, which dequeues

190 a node from the priority queue that gives the smallest distance bound and  
191 recursively traverse the node.

192 **Approximate search:** Approximate search is a common method to  
193 further speed up the kd-tree traversal by trading off the quality of search  
194 results. The approximate search is applicable for both DFS and BBFS kd-  
195 tree traversal methods. Typically, approximate search is achieved by limiting  
196 the number of traversed nodes during the kd-tree traversal, or by interrupting  
197 the search process based upon a real time clock. It has been well-known that  
198 if we limit the number of traversed nodes more, the quality of search results  
199 deteriorates.

200 **Terminologies:** We define  $T(n)$  of a node  $n$  to denote the expected  
201 number of traversed nodes under the sub-tree of the node  $n$ , when the node  $n$   
202 is accessed. We use  $n_l$  and  $n_r$  to denote the left and the right child nodes of a  
203 node  $n$  respectively. We also define  $P_{[n_l|n]}$  to denote a conditional probability  
204 that the left node  $n_l$  is traversed, given the node  $n$  is accessed;  $P_{[n_r|n]}$  is  
205 defined in a similar manner with the right node  $n_r$ . We use  $n_{pdf}(i)$  to denote  
206 the probability distribution function for values of image features that are  
207 projected onto the chosen partitioning normal given a node  $n$ , where  $i$  is  
208 in the range of the minimum,  $m$ , and the maximum values,  $M$ , among the  
209 projected values.

210 **Assumptions for mathematical derivations:** In order to simplify  
211 our derivations for the sake of the clarity, we make a few assumptions: 1)  
212 each leaf node of a kd-tree contains only a single image feature, 2) when we

213 partition image features of a node  $n$  with a hyperplane that has a partitioning  
214 value  $p$  and a chosen partitioning normal, we assign image features whose  
215 values are equal to or less than the partitioning value  $p$  into the left node of  
216 the node  $n$ , and others into the right node, and 3) we treat the values of image  
217 features as continuous. Note that all of these assumptions are introduced to  
218 simplify our derivations, and all of our theoretic results can be shown to be  
219 valid even with kd-trees that do not satisfy such assumptions, after minor  
220 modifications to our derivations.

221 In the following sections we explain how each one of these assumptions  
222 are used for deriving our cost model and its theoretical results. Note that  
223 our cost model measures the expected number of nodes traversed during the  
224 kd-tree traversal. We would also like to point out that our cost model is  
225 still useful for approximate queries that is performed with a fixed number of  
226 traversed nodes. For example, if our cost value for a kd-tree,  $T_1$ , is less than  
227 that of another tree,  $T_2$ , it also means that given a fixed number of nodes  
228 traversed during the kd-tree traversal,  $T_1$  can lead to more accurate results as  
229 compared with  $T_2$ . In this case, we show that our cost model has correlations  
230 with a quality measure on results of approximate queries. More specifically,  
231 we use a *distance error ratio* as the quality measure. This distance error  
232 ratio measures how much the distance between a nearest neighbor computed  
233 within the fixed number of traversed nodes and the query point is over the  
234 ground-truth shortest distance of the nearest neighbor to the query.

#### 235 4. Probabilistic Cost Model

236 We define our probabilistic cost model that quantifies the quality of the  
237 kd-tree by measuring the expected number of nodes traversed during the  
238 nearest neighbor search. Once we access a node, we can access its left or its  
239 right child nodes, irrespective of different traversal methods on a kd-tree for  
240 nearest neighbor search queries. Therefore, we define our probabilistic cost  
241 model for a node  $n$  that measures the expected number of nodes traversed  
242 under the sub-tree rooted at the node  $n$  in a recursive manner, as follows:

$$T(n) = 1 + P_{[n_l|n]}T(n_l) + P_{[n_r|n]}T(n_r), \quad (1)$$

243 where 1 is added after the node  $n$  is traversed. Note that we can easily extend  
244 our cost model to include costs incurred by computing distances between the  
245 query image feature and features contained in leaf nodes. This is because such  
246 cost linearly depends on the number of image features contained in each leaf  
247 node. However, in order to simplify our discussions, we intentionally ignore  
248 such costs.

249 The problem is now reduced to how to accurately define the two condi-  
250 tional probabilities,  $P_{[n_l|n]}$  and  $P_{[n_r|n]}$ , such that they can reflect the actual  
251 performance of kd-tree traversal. There are two main factors for computing  
252 these probabilities: 1) the distribution of image features of potential query  
253 images, and 2) the local geometric configurations of image features around  
254 the node  $n$  and its neighboring nodes.

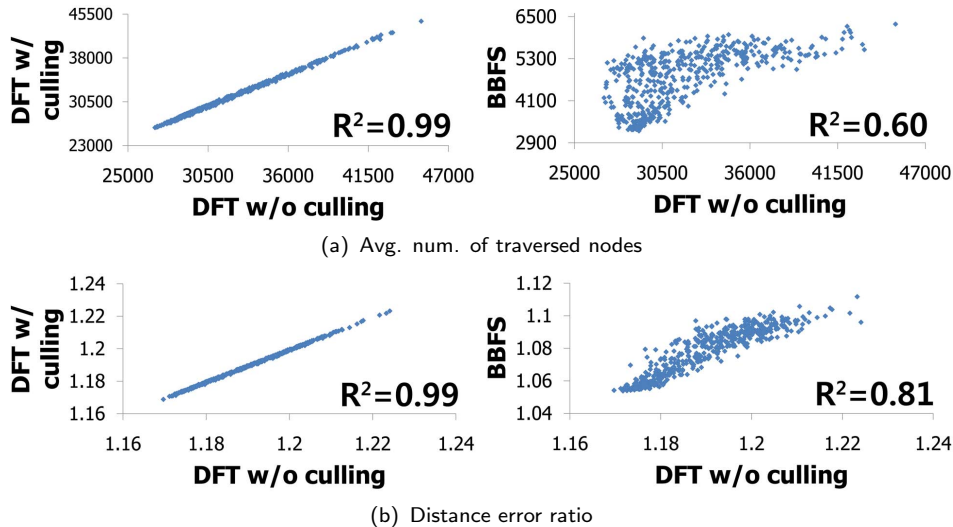


Figure 2: The left figures show correlations between the depth-first traversal (DFT) w/ and w/o culling, in terms of both the number of nodes traversed to find the exact nearest neighbor given query image features and distance error ratios given a fixed number of traversed nodes. The right figures show the correlations between the DFT and best-bin-first search (BBFS) in the same settings used for the left figures.

255     **Effects of culling and traversal methods:**     Many different culling  
256 techniques including the conservative distance culling can be used. These  
257 culling techniques as well as traversal methods can affect conditional prob-  
258 abilities of nodes. We found that it is very hard to consider effects caused  
259 by these culling and traversal methods, since these effects depend on the  
260 local geometric configurations of data stored on a node and its neighboring  
261 nodes. We, however, make an interesting observation: it is highly likely that  
262 employed culling or traversal methods do not change the relative qualities  
263 of different kd-trees. In other words, if a kd-tree,  $T_1$  has a fewer number  
264 of traversed nodes than another tree,  $T_2$ , then  $T_1$  is likely to give a fewer  
265 number of traversed nodes than  $T_2$  even with culling techniques.

266 To verify this observation, we measure correlations between the DFT  
267 with and without culling in terms of the number of traversed nodes that  
268 have been performed to find the nearest neighbor node given a query. To  
269 measure correlations, we construct 500 different kd-trees with image features  
270 from the Caltech 101 image benchmark that consists of around 10 K images;  
271 The details about how these 500 different kd-trees are constructed will be  
272 given in Sec. 6. We perform 10 K different queries used for SIFT-based image  
273 retrieval for each kd-tree. These queries are chosen from images that are in  
274 the same benchmark, but are not used for the tree construction.

275 Correlation results are plotted in Fig. 2 that contains results for both  
276 exact and approximate search queries. For Fig. 2-(a) representing results  
277 from the exact search queries, the x-axis is the number of nodes traversed  
278 by DFT without culling, and the y-axis is the number of nodes traversed by  
279 DFT with culling or BBFS. The left of Fig. 2-(a) shows correlation, 0.99,  
280 between the DFT with and without using the conservative culling. We have  
281 also measured correlations between the DFT and BBFS, and we found that  
282 the correlation is also very high (i.e. 0.6) as shown in the right of Fig. 2-(a).  
283 A simple observation from this graph is that the number of nodes visited  
284 with and without culling are linearly correlated. For a query image, if the  
285 number of nodes visited is  $K$  in DFT without culling, then the number of  
286 nodes visited in DFT with culling is  $\alpha K$ , where  $\alpha$  is a number less than 1.

287 We have also checked the correlations in approximate nearest neighbor  
288 search queries. In this case, we measure correlations in terms of distance

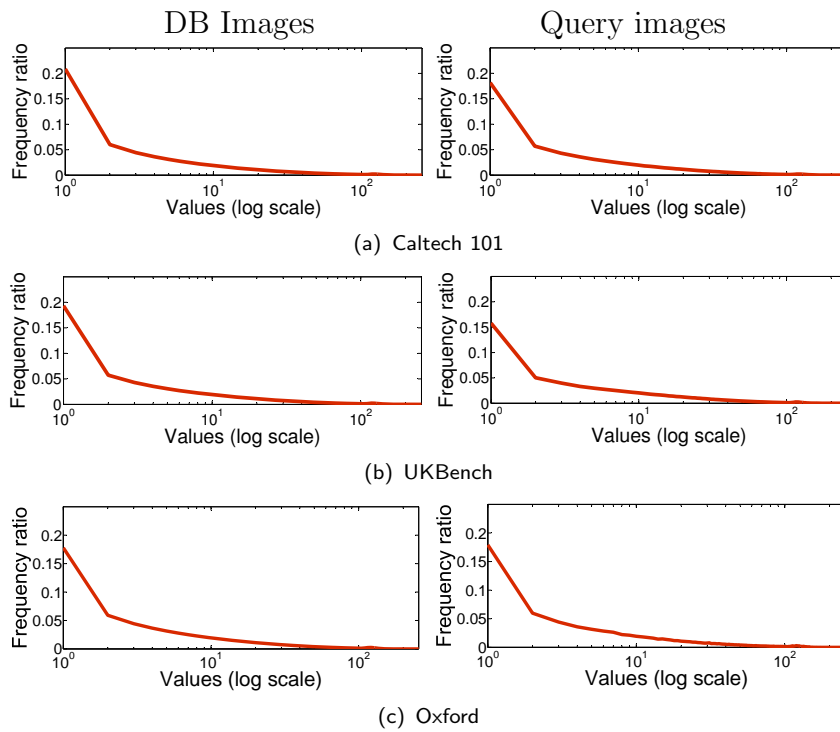


Figure 3: This figure shows distributions of values of SIFT image features; the top and bottom curves are computed from images of Caltech 101, UKBench, and Oxford benchmarks respectively. We use randomly chosen 500 images from Caltech 101 and UKBench (that create 100 K features), and 5000 images from the Oxford benchmark (that creates 1 M features) for building kd-trees (left figures under “DB Images”) and another 500 images (that create another 10 K features) for runtime query images (right figures under “Query images”).

289 error ratios. Again, we observe a high correlation, 0.99, and 0.81, in Fig. 2-  
 290 (b) for the DFT with and without culling, and for the DFT and the BBFS  
 291 respectively. Since there is such high correlation relationship, both in the  
 292 exact nearest neighbor search and the approximate nearest neighbor search,  
 293 we do not consider the effects caused by culling and traversal methods within  
 294 our cost model. Nonetheless, we believe that our model is valid for kd-tree  
 295 traversal method with culling or BBFS.



296 **Distribution of potential query images:** Fig. 3 shows the dis-  
 297 tributions of SIFT image features in our tested image benchmark datasets,  
 298 Caltech 101, UKBench, and Oxford [6]. We simply project values of all the  
 299 dimensions of these image features into one dimensional space to draw the  
 300 distributions. We draw distributions of these values of image features drawn  
 301 from different image sets of the same benchmark and drawn from different  
 302 images of two different benchmarks. Even though they are computed from  
 303 different images across different image benchmarks, these distributions have  
 304 a nearly similar tendency. This is because inside natural images, most re-  
 305 gions are smooth regions or regions with small gradients [28] and our plotting  
 306 on the image features agrees with such natural image statistics. Hence, we  
 307 can expect the distribution of potential query images follows the same distri-  
 308 bution as the distribution of image features that were used to construct the  
 309 kd-tree from an image database.

310 Now, we can define our final cost model. As discussed above, we do not  
 311 consider the effects of culling and traversal methods such as BBFS, since they  
 312 do not drastically change the relative qualities of kd-trees. Also, we assume  
 313 that the distribution of image features of query images follows that of image  
 314 features stored in the kd-tree. Then the conditional probability  $P_{[n_l|n]}$  for the  
 315 left node given a node  $n$  can be computed as  $\int_m^p n_{pdf}(i)$ , where  $p$  is the chosen  
 316 partitioning value used for defining a partitioning hyperplane. Similarly,  
 317  $P_{[n_r|n]}$  is defined as  $\int_{p+}^M n_{pdf}(i)$ . We use the range of  $(p, M]$  for  $P_{[n_r|n]}$ , since  
 318 we assumed that image features that have the partitioning value  $p$  is assigned

319 to the left node. As a result, our final cost model is defined as follows:

$$T(n) = 1 + \int_m^P n_{pdf}(i)T(n_l) + \int_{p+}^M n_{pdf}(i)T(n_r). \quad (2)$$

320 Our cost model will be validated with SIFT-based image retrieval in  
321 Sec. 6.

## 322 5. Near-Optimal Partitioning Value

323 In this section we prove that a near-optimal partitioning value exists  
324 given a chosen partitioning normal, to define a partitioning hyperplane for  
325 the high-quality kd-tree construction. Also, we show that the median-based  
326 partitioning can achieve a near-optimal result given our probabilistic cost  
327 model.

### 328 5.1. Lower Bound of Our Cost Model

329 We first show a lower bound of our cost model.

330 **Lemma 5.1.** *The following equation has the minimum value of  $-1$ , when*  
331  *$x = \frac{1}{2}$ :*

$$y = x \log x + (1 - x) \log (1 - x), \quad (3)$$

332 *where the logarithm function uses base 2 and  $0 < x < 1$ .*

333 **Proof:** We rewrite  $y$  as the following:

$$y = x \frac{\ln x}{\ln 2} + (1 - x) \frac{\ln (1 - x)}{\ln 2}.$$

334 We calculate its first derivative,  $\frac{dy}{dx}$ :

$$\begin{aligned}\frac{dy}{dx} &= \frac{1}{\ln 2}(\ln x + 1) - \frac{1}{\ln 2}[1 + \ln(1 - x)] \\ &= \frac{1}{\ln 2}[\ln x - \ln(1 - x)].\end{aligned}$$

335 The first derivative becomes 0, when  $x = \frac{1}{2}$ . In order to confirm whether  
336  $y$  has the global minima or maxima with  $x = \frac{1}{2}$ , we compute the second  
337 derivative,  $\frac{d^2y}{dx^2}$ :

$$\frac{d^2y}{dx^2} = \frac{1}{\ln 2}\left[\frac{1}{x} + \frac{1}{1 - x}\right].$$

338 Since  $0 < x < 1$ , values of the second derivative are always positive. As a  
339 result,  $y$  has the minima,  $-1$ , when  $x = \frac{1}{2}$ . □

340 **Theorem 5.1 (Lower Bound).** *Given a kd-tree with  $k$  image features, the*  
341 *cost for this kd-tree according to our probabilistic cost model (Eq. 2) is at*  
342 *least  $1 + \log k$ .*

343 **Proof:**

344 Suppose that a node  $n$  has  $k$  image features under the sub-tree rooted at  
345  $n$ . We define  $f(n)$  to be the probability of accessing the left child node  $n_l$   
346 after accessing node  $n$ ; that is  $f(n) = P_{[n_l|n]}$ . Equivalently,  $1 - f(n) = P_{[n_r|n]}$   
347 is the probability of accessing the right child node  $n_r$ . Then, the left child

348 node  $n_l$  has  $kf(n)$  image features <sup>1</sup> under its sub-tree, while the right child  
 349 node  $n_r$  has  $k(1 - f(n))$  image features. For any intermediate node, the range  
 350 of  $x$  is in  $(0, 1)$ , since some probability exists to access its child nodes.

351 Our cost model shown in Eq. 2 can be rewritten as follows:

$$T(n) = 1 + f(n)T(n_l) + (1 - f(n))T(n_r). \quad (4)$$

352 We first prove that  $T(n) \geq 1 + \log k$  by induction. The base case arises  
 353 when a node  $n$  is a leaf. In this case,  $T(n) = 1$ , since the number of traversed  
 354 node is 1 for the leaf node. Since we assumed that each leaf node,  $n_{leaf}$  has  
 355 a single image feature,  $T(n_{leaf}) \geq 1 + \log(1)$ . <sup>2</sup>

356 Let us consider cases when a node  $n$  is intermediate and thus  $k > 1$ . We  
 357 start from the assumption for two child nodes of the node  $n$ :  $T(n_l) \geq 1 +$   
 358  $\log(kf(n))$  and  $T(n_r) \geq 1 + \log(k(1 - f(n)))$ . Then,  $T(n)$  can be rewritten

---

<sup>1</sup>For the sake of clarity, we allow a continuous number of image features in our derivation. One can prove our theorem in a discrete manner by taking similar steps shown in the paper.

<sup>2</sup>If a leaf node has multiple image features, we can extend our derivation by showing  $T(n) \geq \log k + c$ , where  $c$  is a constant.

359 as follows:

$$\begin{aligned} T(n) &= 1 + f(n)T(n_l) + (1 - f(n))T(n_r) \\ &\geq 1 + f(n)(1 + \log(kf(n))) + (1 - f(n))(1 + \log(k(1 - f(n)))) \\ &= 2 + f(n)\log k + f(n)\log f(n) + \\ &\quad (1 - f(n))\log k + (1 - f(n))\log(1 - f(n)) \\ &= 2 + \log k + \\ &\quad f(n)\log f(n) + (1 - f(n))\log(1 - f(n)). \end{aligned}$$

360 According to Lemma 5.1, the term of  $f(n)\log f(n) + (1 - f(n))\log(1 - f(n))$   
361 has the minimum value of  $-1$ , when  $f(n) = 1 - f(n) = \frac{1}{2}$ . As a result, even  
362 when the node  $n$  is intermediate, we can show that  $T(n) \geq 1 + \log k$ . And  
363 this proves the theorem.  $\square$

364 As identified while we prove Theorem 5.1, the lower bound of our cost  
365 model is minimized when for every node, the conditional probability for  
366 accessing the left child node after accessing this node is exactly half. This  
367 does not directly indicate that our cost model is minimized at such case.  
368 Nonetheless, such case can be one of promising candidates for constructing  
369 optimal kd-trees given our cost model.

## 370 5.2. Median-based Partitioning

371 The derived lower bound of our probabilistic cost model for a kd-tree is  
372 minimized when for every node, the conditional probability of accessing its

373 left node is exactly equal to that of the right node. In practice it may be  
 374 impossible to partition image features of a node into left and right nodes  
 375 with the equal conditional probabilities, since SIFT image features in each  
 376 dimension have discrete values and there may be multiple features that have  
 377 the same value along a chosen partitioning normal. The median point among  
 378 image features of a node, however, can be a fairly good candidate to make  
 379 nearly equal conditional probabilities for the left and right child nodes.

380 In the same vein, partitioning normals that are close to be the principal  
 381 eigenvectors computed with Principal Component Analysis (PCA) can serve  
 382 as an excellent choice, especially when used together with median-based par-  
 383 titioning, since they can lead to nearly equal conditional probabilities for  
 384 child nodes. These partitioning normals have been demonstrated to work  
 385 quite well in practice [8, 9].

386 In this section we show that median-based partitioning produces near-  
 387 optimal kd-trees in terms of minimizing our cost model.

Consider a kd-tree constructed using median-based partitioning. Accord-  
 ing to our cost model, the cost of a node  $n$  of the kd-tree containing  $k$  image  
 features can be written as follows:

$$T(n) = 1 + \frac{\lceil k/2 \rceil}{k} T(n_l) + \frac{\lfloor k/2 \rfloor}{k} T(n_r),$$

388 where  $\lceil k/2 \rceil$  and  $\lfloor k/2 \rfloor$  image features are assigned to the left and right nodes  
 389 respectively.

390 Given a kd-tree constructed by median based partitioning, one can ob-  
 391 serve that the cost of the tree only depends on the number of image features.  
 392 Let  $S(k)$  denote the cost of a kd-tree that is computed by median-based par-  
 393 titioning and has set of  $k$  image features.  $S(k)$  can be recursively defined as  
 394 follows:

$$395 \quad S(k) = \begin{cases} 1 & \text{if } k = 1, \\ 1 + \frac{l}{2l}S(l) + \frac{l}{2l}S(l) & \text{if } k = 2l \text{ for some } l \geq 1, \\ 1 + \frac{l+1}{2l+1}S(l+1) + \frac{l}{2l+1}S(l) & \text{if } k = 2l + 1 \text{ for some } l \geq 1. \end{cases}$$

396

397 When  $k$  is a power of two, we can easily show that  $S(k)$  realizes the lower  
 398 bound, which is  $1 + \log k$ , in the next lemma (Lemma 5.2). As a result, in this  
 399 case kd-trees computed by median-based partitioning are optimal in terms  
 400 of our cost model.

401 **Lemma 5.2.** *If  $k = 2^t$  for an integer  $t \geq 0$ , then  $S(k) = 1 + \log k = 1 + t$ .*

402 **Proof:** We use induction to prove this lemma. In the basic step,  
 403  $t = 0$  and  $k = 1$ . Then  $S(k) = 1$  by its definition. Now suppose that  
 404  $k = 2^t$  for some  $t > 0$ . Assuming that  $S(2^{t-1}) = 1 + t - 1$ , we get that  
 405  $S(k) = S(2^t) = 1 + \frac{1}{2}S(2^{t-1}) + \frac{1}{2}S(2^{t-1}) = 1 + S(2^{t-1}) = 1 + 1 + t - 1 = 1 + t$ .  
 406 □

407 In practice, however,  $k$  may not be a power of two. We now introduce a  
 408 series of lemmas to prove that in general cases, costs,  $S(k)$ , of kd-trees with  
 409  $k$  image features constructed by median-based partitioning is near-optimal.

410 Suppose that two kd-trees are constructed by using median-based parti-

411 tioning, but one kd-tree is constructed with  $k$  image features, while another  
 412 one with  $k + 1$  features. We then have the following lemma.

413 **Lemma 5.3.**  $S(k + 1) > S(k)$  for  $k \geq 1$ .

414 **Proof:** We prove this lemma by induction. It is easy to see that  
 415  $S(1) = 1$ ,  $S(2) = 1 + \frac{1}{2}S(1) + \frac{1}{2}S(1) = 2$ , and  $S(3) = 1 + \frac{2}{3}S(2) + \frac{1}{3}S(1) = 8/3$ .  
 416 Therefore, when  $k = 1$  or  $k = 2$ , we have  $S(k + 1) > S(k)$ . We now show  
 417 that  $S(k + 1) > S(k)$ , when  $k = 2l + 1$ , an odd number greater than 1. We  
 418 assume that  $S(t + 1) > S(t)$  for  $1 \leq t \leq 2l$ . Then,

$$\begin{aligned}
 S(k + 1) &= S((2l + 1) + 1) = S(2l + 2) \\
 &= 1 + \frac{l + 1}{2l + 2}S(l + 1) + \frac{l + 1}{2l + 2}S(l + 1) = 1 + S(l + 1) \\
 &= 1 + \frac{l + 1}{2l + 1}S(l + 1) + \frac{l}{2l + 1}S(l + 1) \\
 &> 1 + \frac{l + 1}{2l + 1}S(l + 1) + \frac{l}{2l + 1}S(l) \\
 &= S(2l + 1) = S(k)
 \end{aligned}$$

419 Let us consider the case of  $k = 2l$ , an even number greater than 2. We  
 420 can assume that  $S(t + 1) > S(t)$  for  $1 \leq t \leq 2l - 1$ .



$$\begin{aligned}
S(k+1) &= S(2l+1) \\
&= 1 + \frac{l+1}{2l+1}S(l+1) + \frac{l}{2l+1}S(l) \\
&> 1 + \frac{l+1}{2l+1}S(l) + \frac{l}{2l+1}S(l) \\
&= 1 + S(l) = 1 + \frac{l}{2l}S(l) + \frac{l}{2l}S(l) = S(2l) = S(k)
\end{aligned}$$

421 This proves the lemma. □

422 **Corollary 5.1.**  $S(k+t) \geq S(k)$  for any integer  $t \geq 0$ .

Let a function  $\varphi : \mathbb{Z}^+ \rightarrow \mathbb{Z}^+$  to denote  $\varphi(k) = 2^{\lceil \log k \rceil}$ ;  $\varphi(k)$  is the smallest positive integer greater than or equal to  $k$  such that  $\varphi(k)$  is a power of two. If  $k$  itself is a power of two then  $\varphi(k) = k$ . This function's useful properties are given below:

$$\varphi(k) \geq k \tag{5}$$

$$\log(\varphi(k)) \leq 1 + \log k \tag{6}$$

$$S(\varphi(k)) = 1 + \log(\varphi(k)) \tag{7}$$

423 Equation 5 is obvious from its definition. Equation 6 comes from the  
424 fact that  $\log(\varphi(k)) = \log(2^{\lceil \log k \rceil}) = \lceil \log k \rceil \leq 1 + \log k$ . Finally, Equation 7  
425 follows from Lemma 5.2 and the fact that  $\varphi(k)$  is a power of two. We now  
426 show a tight upper bound of  $S(K)$  based on the function  $\varphi$ .

427 **Lemma 5.4.**  $S(k) \leq 2 + \log k$

428 **Proof:**

$$\begin{aligned} 429 \quad S(k) &\leq S(\varphi(k)) && \text{[ by Corollary 5.1 and (5) ]} \\ &= 1 + \log(\varphi(k)) && \text{[by Equation 7]} \\ &\leq 2 + \log k && \text{[by Equation 6]} \end{aligned}$$

430 □

431

432 For a kd-tree constructed by median-based partitioning, Lemma 5.4 shows  
433 that its cost is at most  $2 + \log k$ , where  $k$  is the number of image features.  
434 Theorem 5.1 shows that the cost of any kd-tree with  $k$  image features is lower  
435 bounded by  $1 + \log k$ . As a result, the following theorem naturally holds.

436 **Theorem 5.2 (Near-Optimality).** *Given our cost model (Eq. 2), let  $C$*   
437 *to be a cost of a kd-tree with  $k$  image features constructed by median-based*  
438 *partition. Then  $C \leq 1 + OPT$ , where  $OPT$  is the minimum cost of any*  
439 *kd-tree with the same set of those image features.*

440 This theorem indicates that median-based partitioning produces near-  
441 optimal kd-trees given our cost model.

## 442 6. Experimental Validations

443 We randomly choose 500 images from the Caltech 101 image benchmark  
444 and extract around 100 K SIFT features from those images for constructing  
445 the kd-trees. 500 different kd-trees are constructed by randomly choosing  
446 the partitioning values given the partitioning normals. For each kd-tree, we

447 evaluate our cost model with the tree by measuring the cost value associ-  
448 ated with the root node of the kd-tree. We have also measured the number  
449 of traversed nodes to find the nearest neighbor features given 10 K query  
450 SIFT features. These 10 K SIFT features are extracted from different 500  
451 query images for our SIFT-based image retrieval. These query images are  
452 randomly chosen from the same image benchmark, but are not used for the  
453 tree construction.

454 In this configuration, we found that our cost model shows a positive cor-  
455 relation, 0.526, against the average number of nodes traversed with the DFT.  
456 This correlation score is lower than the other correlation score presented in  
457 this paper. This is mainly because when we perform the exact search queries,  
458 the traversal methods tend to traverse most of nodes of kd-trees irrespective  
459 of the qualities of kd-trees. Because of this behavior of the exact nearest  
460 neighbor query, approximate search queries are typically used in practice.  
461 We have also measured the correlation by applying the same setting, but  
462 this time we use BBFS and culling with the UKBench benchmark. The  
463 correlation score goes up to 0.79.

464 To further analyze the correlation of our cost model against the quality  
465 of kd-tree, we measure the time spent to traversing kd-trees. Our cost model  
466 shows a high correlations, 0.74, with the average time spent on performing  
467 exact search queries that run by using the BBFS in the UK Bench image  
468 benchmark and the Caltech 101 image benchmark.

469 We check the correlation of our cost model with approximate nearest

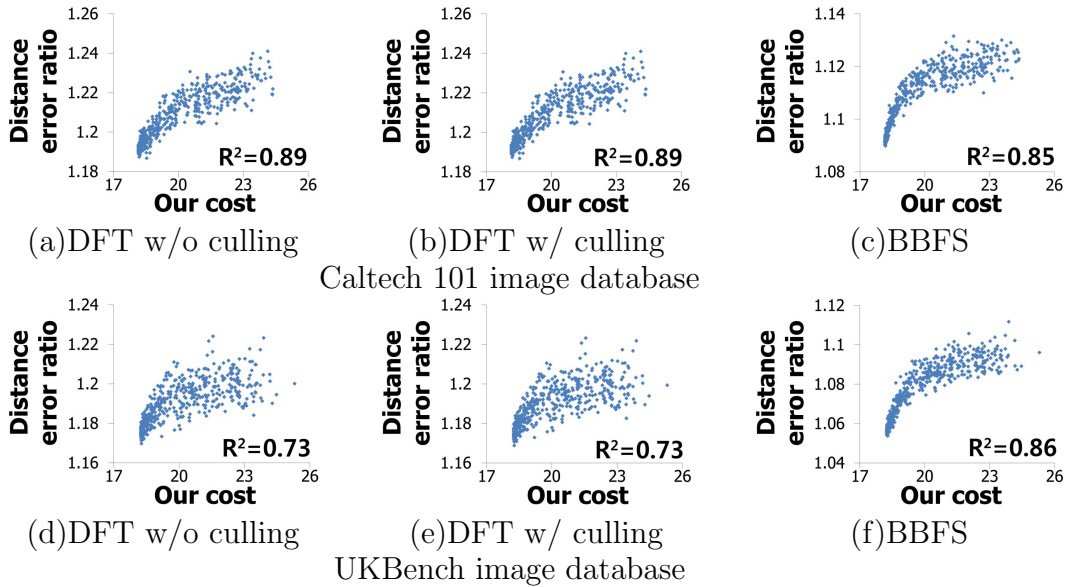


Figure 4: This figure shows correlations of our cost model against the distance error ratios achieved by (a)(d) the DFT w/o culling, (b)(e) DFT with culling, and (c)(f) BBFS in the Caltech 101 and UKBench image databases.

470 neighbor search queries. In this case, we measure the correlation between  
 471 values of our cost model and distance error ratios achieved when we limit the  
 472 number of traversed nodes to 1 K nodes. Fig. 4 shows the correlations of our  
 473 cost model against the distance error ratios achieved by different traversal  
 474 methods for the Caltech 101 and UKBench image benchmarks. Our cost  
 475 model shows high correlations, 0.89, 0.89, and 0.85, over DFT w/o culling,  
 476 DFT w/ culling, and BBFS respectively in the Caltech 101 benchmark. We  
 477 observe that our cost model shows similar, high correlations, 0.73, 0.73, and  
 478 0.86, against DFT w/o culling, DFT w/ culling, and BBFS respectively in  
 479 the UKBench image benchmark.

480 We have also found that median partitioning consistently gives the lowest

481 cost in term of both our cost model and the distance error ratios, as con-  
482 jectured in Sec. 5.2. More specifically speaking, a kd-tree constructed by a  
483 partitioning value among the median,  $m$ , and  $m \pm 1$  given a node shows the  
484 best result, compared to those 500 different kd-trees that are constructed by  
485 randomly choosing the partitioning values.

## 486 7. Conclusion

487 We have presented a probabilistic cost model that measures the expected  
488 number of traversed nodes during the kd-tree traversal. Our cost model has  
489 demonstrated to have high correlations with the observed numbers of tra-  
490 versed nodes as well as the time spent on image search under different culling,  
491 traversal methods including exact and approximate queries. Furthermore,  
492 our cost model can explain why the commonly adopted partitioning meth-  
493 ods such as median-based and PCA-based techniques work well and showed  
494 that they can achieve a near-optimal quality for the kd-tree construction.

495 **Limitations and future work:** There are many avenues for future  
496 research directions. First, we would like to see how well our theoretical  
497 and experimental results extend to web-scale image databases that consist  
498 of billions of images. In Sec. 4, we showed that culling and kd-tree traversal  
499 methods do not change the relative qualities of different kd-trees much and  
500 thus we do not consider them in our probabilistic cost model. However, it may  
501 be possible to consider a particular traversal algorithm (e.g., using multiple  
502 randomized kd-trees [10]) with culling and back-tracking properties in a cost

503 model, and use it to construct a kd-tree that has a higher quality than the  
504 one constructed by the proposed probabilistic model. Moreover, we would  
505 like to extend our current cost model to consider the dimensionality of data  
506 points for more accurate estimation of costs. Also, we would like to apply  
507 our cost model and its optimal partitioning theorem to incrementally update  
508 the kd-trees for dynamic data sets in image retrieval. Prior approaches for  
509 dynamic image datasets are based on heuristic techniques that decide when  
510 and where to reconstruct kd-trees. However, through our cost model and  
511 optimal partitioning value theorem, we can approach this problem in a more  
512 rigorous manner. We wish that our cost model can serve as a theoretical  
513 basis that leads to more rigorous techniques for various problems related to  
514 nearest neighbor search queries including recent hashing techniques [29].

## 515 **Acknowledgements**

516 This work was supported in part by MEST/NRF (2011-0030822), MCST/  
517 KOCCA/CT/R&D 2011, MKE/KEIT [KI001810035261], MKE/MCST/IITA  
518 [2008-F-033-02], BK, DAPA/ADD (UD110006MD), MEST/NRF/WCU (R31-  
519 2010-000-30007-0), KMCC, and MSRA.

## 520 **References**

- 521 [1] H. Samet, Foundations of MultiDimensional and Metric Data Struc-  
522 tures, Morgan Kaufmann, 2006.

- 523 [2] R. Datta, D. Joshi, J. Li, J. Z. Wang, Image retrieval: Ideas, influences,  
524 and trends of the new age, *ACM Computing Survey* 40 (2) (2008) 1–60.
- 525 [3] K. Q. Weinberger, J. Blitzer, L. K. Saul, Distance metric learning for  
526 large margin nearest neighbor classification, in: *NIPS*, 2006.
- 527 [4] D. Lowe, Distinctive image features from scale-invariant keypoints, *IJCV*  
528 60 (2) (2004) 91–110.
- 529 [5] J. H. Friedman, J. L. Bentley, R. A. Finkel, An algorithm for finding  
530 best matches in logarithmic expected time, *ACM TOMS* 3 (3) (1977)  
531 209–226.
- 532 [6] J. Philbin, O. Chum, M. Isard, J. Sivic, A. Zisserman, Object retrieval  
533 with large vocabularies and fast spatial matching, in: *IEEE Conference*  
534 *on Computer Vision and Pattern Recognition*, 2007, pp. 1–8.
- 535 [7] S. Arya, D. M. Mount, N. S. Netanyahu, R. Silverman, A. Wu, An op-  
536 timal algorithm for approximate nearest neighbor searching, in: *Symp.*  
537 *on Discrete Alg.*, 1994, pp. 573–582.
- 538 [8] C. Silpa-Anan, R. Hartley, S. Machines, A. Canberra, Optimised kd-  
539 trees for fast image descriptor matching, in: *IEEE Conference on Com-*  
540 *puter Vision and Pattern Recognition*, 2008, pp. 1–8.
- 541 [9] Y. Jia, J. Wang, G. Zeng, H. Zha, X.-S. Hua, Optimizing kd-trees for  
542 scalable visual descriptor indexing, in: *CVPR*, 2010.

- 543 [10] M. Muja, D. G. Lowe, Fast approximate nearest neighbors with auto-  
544 matic algorithm configuration, in: VISSAPP, 2009, pp. 331–340.
- 545 [11] J. Sivic, A. Zisserman, Video google: A text retrieval approach to object  
546 matching in videos, in: ICCV, 2003.
- 547 [12] D. Nistér, H. Stewénus, Scalable recognition with a vocabulary tree, in:  
548 CVPR, 2006.
- 549 [13] M. Datar, N. Immorlica, P. Indyk, V. S. Mirrokni, Locality-sensitive  
550 hashing scheme based on p-stable distributions, in: SoCG, 2004.
- 551 [14] D. MacDonald, B. K.S., Heuristics for ray tracing using space subdivi-  
552 sion, *Visual Computing* 6 (3) (1990) 153–166.
- 553 [15] I. Wald, V. Havran, On building fast kd-trees for ray tracing, and on  
554 doing that in  $O(N \log N)$ , in: Proceedings of the 2006 IEEE Symposium  
555 on Interactive Ray Tracing, 2006, pp. 61–69.
- 556 [16] W. Hunt, W. Mark, G. Stoll, Fast kd-tree construction with an adap-  
557 tive error-bounded heuristic, in: IEEE Symposium on Interactive Ray  
558 Tracing 2006, 2006, pp. 81–88.
- 559 [17] S. Yoon, S. Curtis, D. Manocha, Ray tracing dynamic scenes using se-  
560 lective restructuring, *Eurographics Symp. on Rendering (2007)* 73–84.
- 561 [18] J. L. Bentley, Multidimensional binary search trees used for associative  
562 searching, *Communications of the ACM* 18 (1975) 509–517.



- 563 [19] D. E. Knuth, The Art of Computer Programming, Vol. 1, Addison-  
564 Wesley, 1969.
- 565 [20] C. Faloutsos, T. Sellis, N. Roussopoulos, Analysis of object oriented  
566 spatial access methods, in: ACM SIGMOD, 1987.
- 567 [21] J. G. Cleary, Analysis of an algorithm for finding nearest neighbors in  
568 euclidean space, ACM Transactions on Mathematical Software (TOMS)  
569 5 (1979) 183–192.
- 570 [22] C. Böhm, A cost model for query processing in high dimensional data  
571 spaces, ACM Transaction on Database Systems 25 (2000) 129–178.
- 572 [23] D. T. Lee, C. K. Wong, Worst-case analysis for region and partial re-  
573 gion searches in multidimensional binary search trees and balanced quad  
574 trees, Acta Informatica 9 (1) (1977) 23–29.
- 575 [24] J. F. Sproull, Refinements to nearest-neighbor searching in  $k$ -  
576 dimensional trees, Algorithmica 6 (1991) 579–589.
- 577 [25] S. Arya, D. M. Mount, O. Narayan, Accounting for boundary effects in  
578 nearest neighbor searching, in: Symposium on Computational Geometry  
579 (SoCG), 1995.
- 580 [26] S. Berchtold, C. Böhm, D. A. Keim, H.-P. Kriegel, A cost model for  
581 nearest neighbor search in high-dimensional data space, in: Proceed-  
582 ings of the sixteenth ACM SIGACT-SIGMOD-SIGART symposium on  
583 Principles of database systems, 1997.

- 584 [27] D. G. Lowe, Object recognition from local scale-invariant features, in:  
585 the International Conference on Computer Vision, 1999.
- 586 [28] A. Torralba, A. Oliva, Statistics of natural image categories, *Network:  
587 Computation in Neural Systems* 14 (3) (2003) 391–412.
- 588 [29] J.-P. Heo, Y. Lee, J. He, S.-F. Chang, S.-E. Yoon, Spherical hashing, in:  
589 CVPR, 2012, to appear.