



## Audio Fingerprinting: Nearest Neighbor Search in High Dimensional Binary Spaces\*

MATTHEW L. MILLER

*NEC Laboratories, 4 Independence Way, Princeton, NJ 08540*

MANUEL ACEVEDO RODRIGUEZ

*EPFL, 1015 Lausanne, Switzerland and Eurecom Institute, BP 193-06904, Sophia-Antipolis, France*

INGEMAR J. COX

*Department of Computer Science, University College London, Gower Street, London WC1E 6BT*

**Abstract.** Audio fingerprinting is an emerging research field in which a song must be recognized by matching an extracted “fingerprint” to a database of known fingerprints. Audio fingerprinting must solve the two key problems of representation and search. In this paper, we are given an 8192-bit binary representation of each five second interval of a song and therefore focus our attention on the problem of high-dimensional nearest neighbor search. High dimensional nearest neighbor search is known to suffer from the curse of dimensionality, i.e. as the dimension increases, the computational or memory costs increase exponentially. However, recently, there has been significant work on efficient, approximate, search algorithms. We build on this work and describe preliminary results of a probabilistic search algorithm. We describe the data structures and search algorithm used and then present experimental results for a database of 1,000 songs containing 12,217,111 fingerprints.

### 1. Introduction

Audio fingerprinting seeks to recognize a song by extracting a compact representation from an arbitrary, say 5 s, interval and comparing this fingerprint to a database of known fingerprints. There are two primary applications of such technology; usage monitoring to determine broadcast usage fees and media linking in which, for example, a consumer is able to transmit a fingerprint using his or her cell phone to a database

and receive back identification information as well as a web site where the song can be purchased.

Audio fingerprinting falls within the domain of classical pattern recognition and must therefore solve both the representation and matching problems. The problem is more tractable than many pattern recognition problems (e.g. three dimensional face recognition) in that the range of distortions is relatively small. Nevertheless, the problem is still difficult because, for the applications we envision, the database must be very large, perhaps one million songs. If we assume approximately 10,000 unique fingerprints per song, this means we may wish to search through about 10 billion fingerprints.

Ideally, the fingerprint should be a compact representation that is invariant to a variety of common

\*This work was performed while the author was at NEC Research Institute, Princeton. © [2002] IEEE, Reprinted, with permission, from Miller, M.L.; Rodriguez, M.A.; Cox, I.J.; Multimedia Signal Processing, 2002 IEEE Workshop on, 2002 Page(s): 182–185.

distortions, including additive noise, low pass filtering, subsampling, lossy compression and small offsets in the origin of the fingerprint. Distortions to a song will introduce distortions in the corresponding fingerprint. Obviously, the less the fingerprint is distorted, the easier the subsequent search will be.

In this paper, we assume a representation developed by Haitsma et al. [6] in which each 5 s interval is represented by an 8192-bit fingerprint. Other representations are possible. See for example, [1, 3, 4].

Haitsma et al.'s fingerprint is the concatenation of 256 32-bit words. Each word is derived by taking the Fourier transform of  $5/256$  second overlapping intervals, thresholding these values and subsequently computing a 32-bit hash value. Although the words are non-overlapping in time, the concatenated fingerprints have substantial temporal overlap. Assuming that the duration of the average song is 3 min, the number of fingerprints per song is  $3 \times 60 \times 256/5 = 9216 \approx 10,000$ .

Given this 8192-bit representation, the focus of our work has been to develop an efficient search algorithm. Previous work in the context of audio search is either linear in the size of the database [3] or has assumed that there exists a subset of the fingerprint that can be matched perfectly. There are efficient methods for exact matching, and these techniques can be used to identify candidate songs which are then tested against the entire fingerprint [5, 9, 10]. For large databases of songs, sublinear search complexity is required. Empirically, it appears that the assumption of an exact match between a subset of the query and database fingerprints is good. However, the search method described next does not require this condition to be met.

Our objective is to find the bit-sequence in the database with the smallest hamming distance, or bit error rate (hamming distance/8192), from a given query. This search can be characterized as a nearest-neighbor search in a very high (8192) dimensional space. Of course, this assumes that the nearest fingerprint in the database to the query is the correct match. However, under some distortions, this assumption may not be valid. Interestingly, it is less important to correctly match the query to its corresponding fingerprint as it is to match the fingerprint to its corresponding song. Since a song is composed of many fingerprints, incorrectly matching a query to a fingerprint does not necessarily lead to a song recognition error. This issue is discussed in more detail in the experimental evaluation of Section 3.

Multi-dimensional nearest-neighbor search is a well-studied problem. Proposed solutions generally

create a tree structure, the leaf nodes representing the data (fingerprints), and searching becomes a traversal of the tree. Specific algorithms differ in how this tree is constructed and traversed. Two related data structures, *kd*-trees and vantage point or *vp*-trees, have been extensively studied. However, both data structures succumb to a problem known as “the curse of dimensionality”: as dimensionality increases, an increasing percentage of the tree must be searched in order to locate the nearest neighbor to a query. If the data has much more than 30 or so dimensions, most algorithms end up searching essentially all of the tree.

Recent work [2, 7, 8, 11] appears to acknowledge the fact that a perfect search that guarantees to find the nearest neighbor in a high dimensional space is not feasible. However, the curse of dimensionality can be removed if the search is approximate. For example, Yianilos [11] describes an algorithm that, with probability,  $p$ , will find a neighbor within a Euclidean distance  $r$  of the query when the datum are uniformly distributed within an  $n$ -dimensional hypercube. Unfortunately, this work has not been extended to the binary case and Hamming rather than Euclidean distance.

In this paper, we develop an approximate search algorithm for high dimensional binary vectors. Section 2 first describes the algorithm. Section 3 then presents experimental results on a database of 1000 songs and 12,217,111 fingerprints. Finally, Section 4 summarizes our results and discusses possible avenues of future work.

## 2. Algorithm

In the following subsections, we describe an approximate search algorithm for binary vectors in a high dimension space.

Given the set of known fingerprints, we first construct a 256-ary tree. Each 8192-bit fingerprint is represented as 1024 8-bit bytes. The value of each consecutive byte in the fingerprint determines which of the 256 possible children to descend. A path from the root node to a leaf defines a fingerprint.

As the depth of the tree increases, it is common to find nodes with only a single child. This is because the number of actual fingerprints is very much less than the total possible number. For efficiency purposes, we compress such sequences of nodes with only one child into a single node that represents multiple bytes. We

consider the level,  $l$ , of a node to be the number of bytes represented by the path from the root to that node.

Our search algorithm is guided by a table,  $T$ , indexed by a node level,  $l$ , and a number of errors,  $e$ . During a search, when we visit any node  $n_{l,i}$ , at level  $l$  in the tree, we will have examined  $b = 8l$  bits of the vector and seen  $e$  errors between the first  $b$ -bits of the query and the first  $b$ -bits represented by the path to node  $n_{l,i}$ . The probability,  $p$ , of observing  $x$  errors in  $b$ -bits with a bit error rate (BER) of  $r$  is a binomial distribution, i.e.

$$p(x|b, r) = \binom{b}{x} r^x (1-r)^{b-x} \quad (1)$$

assuming that the bit errors are uniformly distributed over the entire fingerprint.

The probability that we will observe *at least*  $E$  errors in  $b$ -bits is simple one minus the cumulative probability of  $p(x|b, r)$ , i.e.

$$p(e_o \geq E|b, r) = 1 - \sum_0^E p(x|b, r) dx. \quad (2)$$

Having observed  $e$  errors in  $b$ -bits, we can use Eq (2) to calculate the bit error rate  $r_t$ , between the query and the closest fingerprint under the node, such that the probability of observing *at least*  $e$  errors is greater than a threshold  $p_t$ . That is,  $r_t$  is such that  $p(e_o \leq e|b, r_t) = p_t$ .

If we increase the probability,  $p_t$ , of observing at least  $e$  errors, then the calculated bit error rate  $r_t$  will increase. That is, the higher the bit error rate, the more likely it is to observe *at least*  $e$  errors. Thus, in order to provide a conservative estimate of the bit error rate, we set the probability threshold,  $p_t = 0.4$ . This means we can be reasonably certain that, if we were to find the closest fingerprint under this node, its overall bit error rate would be greater than  $r_t$ .

In order to prune the search, we estimate the bit error rate at each node we traverse based on the number of bits,  $b$ , and errors,  $e$ , we have observed in descending from the root to the node. If this estimated bit error rate,  $r_t$  is greater than the bit error rate of the best candidate fingerprint we have found so far, then we can be confident that we will not find a closer fingerprint below this node.

In order to determine the estimated bit error rate,  $r_t$ , for a particular node in the tree, we pre-compute a table of  $r_t$  values for a given probability threshold,  $p_t$ . The table is indexed by  $b$  (rows) and  $e$  (columns, where the

Table 1. Table of  $r_t$  values.

0.00	0.06	0.17	0.28	0.40	0.52	0.61	0.61	0.61
0.00	0.03	0.08	0.14	0.20	0.26	0.32	0.38	0.44
0.00	0.02	0.06	0.09	0.13	0.17	0.21	0.25	0.29
0.00	0.02	0.04	0.07	0.10	0.13	0.16	0.19	0.22
0.00	0.01	0.03	0.06	0.08	0.10	0.13	0.15	0.17
0.00	0.01	0.03	0.05	0.07	0.09	0.11	0.13	0.15
0.00	0.01	0.02	0.04	0.06	0.07	0.09	0.11	0.12
0.00	0.01	0.02	0.04	0.05	0.06	0.08	0.09	0.11

first column represents 0 errors encountered). Part of such a table is given in Table 1 for  $p_t = 0.4$  and a range of  $b$  and  $e$ .

Given the tree and associated table, we now describe the search algorithm.

### 2.1. Search Algorithm

Prior to finding an initial candidate fingerprint for our query, we decide on a maximum tolerable error rate and assign this to a variable, `best_err_rate`, which represents the best error rate seen so far. In essence, we lie to the algorithm, telling it we have already found a fingerprint that is this close to the query. This limits the search, so that if the algorithm finds it is unlikely to find a closer fingerprint, it will fail, rather than searching more of the tree. For our experiments, this value was 0.55, which is probably larger than what we would use in a real application.

The search begins by comparing the first byte of the query with the children of the root node. For each child node we calculate the cumulative number of bit errors seen so far. This is simply the sum of the parent errors and the Hamming distance between the 8-bit value represented by the child and the corresponding 8-bit value in the query. Then a test is applied to each child, in order of increasing error, to determine whether to search that child. If the `best_error_rate` is greater than the  $r_t$  threshold determined by our table, then the child is searched. This is because the table indicates that we expect the bit error rate under the child to be less than the best rate seen so far. Conversely, if the corresponding  $r_t$  threshold is greater than the best error rate seen so far, then we will not search this or subsequent children of the parent node.

The search continues recursively and when a leaf node is reached, the error rate associated with the retrieved fingerprint is compared to the best error rate

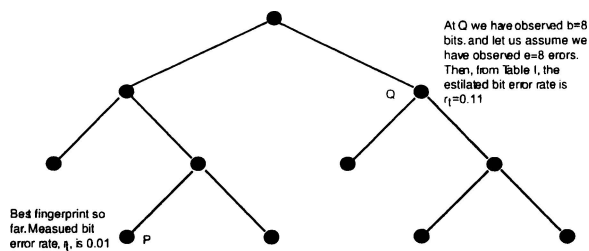


Figure 1. A simple tree illustrating how the tree is pruned. The child node *P* represents the best fingerprint found so far. The bit error rate at this node is 0.01. When we examine node *Q*, we have observed 8 bits and 8 errors. Using Table 1 we determine that the estimated bit error rate is  $r_t = 0.11$ . Since this is larger than for *P*, we can prune this node.

seen so far. If it is less, then we update the best error rate to this new value and assign this fingerprint as our best candidate nearest neighbor so far.

This is illustrated in Fig. 1 where we have shown a tree of depth 4.

Pseudo-code for this algorithm is given below.

```

query = signature obtained from
        input
best_fingerprint = none found
best_err_rate = max err rate
                tolerable

root.errs = 0
search( root )

search( node ) {
  if node is leaf {
    err_rate = node.errs /
              (num bits in query)
    if err_rate < best_err_rate {
      best_fingerprint = fingerprint
                        for this node
      best_err_rate = err_rate }}
  else {
    for each child of node
      child.errs = node.errs +
                  hamming distance
                    between child.bits
                    and
                    corresponding bits
                    of
                    query
    for each child of node, in
    increasing
    order of child.errs

```

```

if best_err_rate >
  T[child.level][child.errs]
  search(child)
else
  return }}

```

Note that the above algorithm can be used to implement a variety of different searching behaviors by defining the threshold table in a variety of different ways. In particular, we can obtain an exact search, which is guaranteed to find the closest neighbor, by letting  $T[l][e]=e/8192$  for each level,  $l$ . This is the error rate we will obtain if there are no further errors in the bytes below this node, and hence the minimum possible error rate obtainable after observing  $e$  errors. The algorithm would then search every node that has even the slightest possibility of yielding a fingerprint closer than the closest found so far, and is thus guaranteed to find the closest.

### 3. Experimental Results

Prior to investigating our search algorithm, we first tested whether fingerprints from one song exhibited any correlation with fingerprints from other songs. Figure 2 illustrates a typical histogram of bit error rates, and shows that inter-song correlation appears fairly random, the average bit error rate being 0.5.

We next examined the correlation between fingerprints from the same song. Figure 3 shows that there is a strong temporal correlation in the vicinity of a

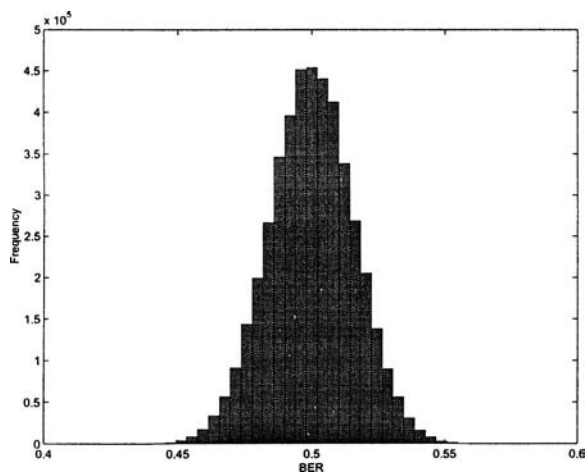


Figure 2. Histogram of bit error rates between a fingerprint from one song and the fingerprints of all other songs.

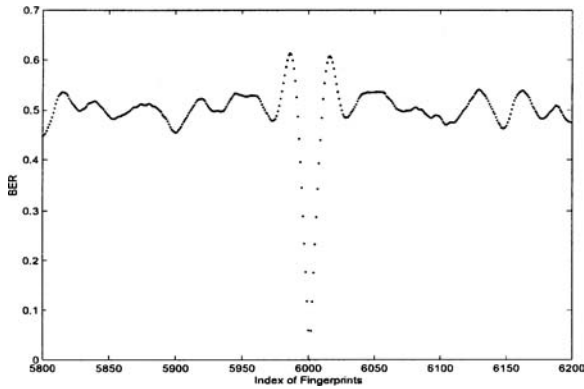


Figure 3. Computed Hamming distance between a fingerprint from one song and all the fingerprints from the same song.

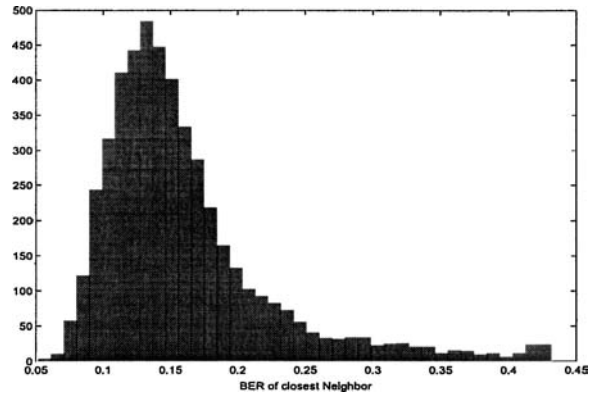


Figure 4. Distribution of distances between queries and their closest fingerprints in the database.

fingerprint. The duration of this correlation is approximately 0.33 s. Fortunately however, even if we match a fingerprint to one of its temporal neighbors, we will still correctly identify the song. In fact, the more intra-song correlation that exists, the easier the identification becomes.

In order to evaluate our search algorithm, we digitized 1,000 songs and computed 12,217,111 corresponding fingerprints. 4913 queries were generated by randomly selecting 5 s portions from the song database and playing them through inexpensive speakers attached to a PC. These song snippets were then digitized using an inexpensive microphone.

For each query, an exhaustive search of the database was performed in order to determine the closest fingerprint in the database. This provided ground truth data with which to compare with our approximate search results. The distribution of distances is shown in Fig. 4. We see that the Hamming distance between the majority of queries and their nearest neighbors is approximately 1065, which is equivalent to a bit error rate of 13%. However, the range of Hamming distances is between 423 and 3532 with corresponding bit error rates of 5 and 43% respectively.

Since we know which song each query is derived from, we were also able to determine the validity of our assumption that the nearest neighbor identifies the correct song. In 53 cases, the nearest neighbor did not correctly identify the song. This is approximately a 1% error rate. However, as we are here concerned only with the performance of our search, and not with the performance of the signature, we regard the search as finding the “correct” song if it finds the song that

contains the closest signature, regardless of whether this is the song that generated the query.

### 3.1. Search Results

Different thresholds for  $p_t$  and corresponding tables for  $r_t$  will produce different search statistics. Typically, the smaller the threshold  $p_t$ , the closer our approximation will be to an exact search, but this comes at the expense of searching more nodes in the tree. Clearly, we seek to find a compromise threshold in which acceptable retrieval accuracy is obtained while searching as little of the tree as possible.

Here, we describe results in which  $p_t = 0.40$ . The tree contains 12,217,111 fingerprints. This is rather large, given that the tree only contains 1,000 songs. However, there is considerable redundancy in the tree since every overlapping fingerprint is represented.

With the queries drawn from the distribution of Fig. 4, we search an average of 419,380 nodes, which is 2.53% of the nodes in the tree.

On average, 85% of the correct nearest neighbors are found. However, the erroneous song rate is only 1%. We attribute this discrepancy to the significant temporal correlation between fingerprints in a song.

Clearly, the distribution of queries can have a very significant effect on the overall search performance. To investigate this effect, we tabulated the number of nodes examined as a function of the distance between the query and its nearest neighbor. Figure 5 shows that for queries with nearest neighbors less than a Hamming distance of 2000, the percentage of nodes visited is much less than 10%. However, as the Hamming

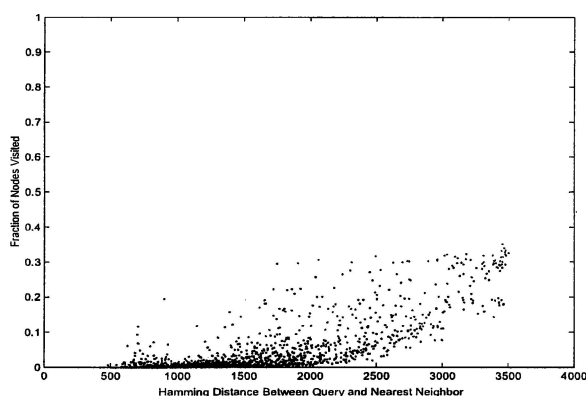


Figure 5. Number of nodes examined as a function of the known distance between a query and its nearest neighbor in the database.

distance increases from 2000 to 3500 a much greater fraction of the database must be searched.

If only queries with a Hamming distance of 2900 or less are considered, then the erroneous song detection drops to 0.3% and the average number of nodes visited is 346,770 or 2.09%. The number of correctly identified fingerprints does not improve significantly and is 86%. However, as previously noted, it is the song detection rate that is important.

#### 4. Discussion

We described preliminary results for an approximate search algorithm which can be used to identify songs based on 5 s samples. These samples are represented as a 8192-bit vector and we proceeded to develop an approximate search algorithm for high-dimensional nearest neighbor search.

This search algorithm constructs a 256-ary tree and an associated table that is constructed assuming that the errors between a query and its nearest neighbor are uniformly distributed. Branches of the tree are pruned by comparing the bit error rate of the current best candidate with the likely bit error rate we expect from the nearest neighbor below a node in the tree.

For a database of 1000 songs and 12,217,111 fingerprints, we demonstrated a song recognition rate of 99% while on average only searching 2.53% of the nodes in the tree. Queries can be more or less noisy and the distribution of queries can significantly affect the search. If only queries with a Hamming distance of less than 2900 to their nearest neighbor are considered, then we are able to achieve a song recognition rate of

99.7% while only visiting an average of 2.09% of the tree.

Significant future work is possible. This includes studying the effect of the threshold probability,  $p_t$  on the song recognition rate and investigating how the performance scales with the size of the song database. We also believe that a more accurate Bayesian model of the search process might be formulated, leading to a better table,  $T$ , and further reducing the number of nodes searched for the same probability of finding the closest fingerprint.

#### Acknowledgment

The authors thank Ton Kalker and Philips Corporation for use of their fingerprint extraction software.

#### References

1. E. Allamanche, J. Herre, O. Hellmuth, B. Bernhard Frobach, and M. Cremer, "AudioID: Towards Content-Based Identification of Audio Material," in *Proc. Int. Conf. on Web Delivering of Music*, 2001.
2. S. Arya, D.M. Mount, N.S. Nctanyahu, R. Silverman, and A.Y. Wu, "An Optimal Algorithm for Approximate Nearest Neighbor Searching in Fixed Dimensions," *J. of the ACM*, vol. 45, no. 6, 1998, pp. 891–923.
3. Y. Cheng, "Music Database Retrieval Based on Spectral Similarity," in *Int. Symp. on Music Information retrieval*, 2001.
4. D. Fragoulis, G. Rousopoulos, T. Panagopoulos, C. Alexiou, and C. Papaodysseus, "On the Automated Recognition of Seriously Distorted Musical Recordings," *IEEE Trans. on Signal Processing*, vol. 49, no. 4, 2001, pp. 898–908.
5. J. Haitsma and T. Kalker, "A Highly Robust Audio Fingerprinting System," in *3rd Int. Conf. on Music Information Retrieval ISMIR*, 2002.
6. J. Haitsma, T. Kalker, and J. Oostveen, "Robust Audio Hashing for Content Identification," in *Content Based Multimedia Indexing*, 2001.
7. P. Indyk and R. Motwani, "Approximate Nearest Neighbors: Towards Removing the Curse of Dimensionality," in *Proc. of the 30th annual ACM Symp. on Theory of Computing*, 1998, pp. 604–613.
8. J.M. Kleinberg, "Two Algorithms for Nearest-Neighbor Search in High Dimensions," in *Proc. of the 29th annual ACM Symp. on Theory of Computing*, 1997, pp. 599–608.
9. H. Neuschmied, H. Mayer, and E. Battle, "Content-Based Identification of Audio Titles on the Internet," in *Proc. Int. Conf. on Web Delivering of Music*, 2001.
10. J. Oostveen, T. Kalker, and J. Haitsma, "Feature Extraction and a Database Strategy for Video Fingerprinting," in *5th Int. Conf. on Visual Information Systems*, vol. LNCS 2314, pp. 117–128. Springer Verlag, 2002.
11. P.N. Yianilos, "Locally Lifting the Curse of Dimensionality for nearest Neighbor Search," in *Proc. of the 11th annual ACM-SIAM Symp. on Discrete Algorithms*, 2000, pp. 361–370.

**Ingemar J. Cox** is currently Professor and Chair of Telecommunications in the Departments of Electronic Engineering and Computer Science at University College London and Director of UCL's Adastral Park Postgraduate Campus. He is currently a holder of a Royal Society Wolfson Fellowship. He received his B.Sc. from University College London and Ph.D. from Oxford University. He was a member of the Technical Staff at AT&T Bell Labs at Murray Hill from 1984 until 1989 where his research interests were focused on mobile robots. In 1989 he joined NEC Research Institute in Princeton, NJ as a senior research scientist in the computer science division. At NEC, his research shifted to problems in computer vision and he was responsible for creating the computer vision group at NECL. He has worked on problems to do with stereo and motion correspondence and multimedia issues of image database retrieval and watermarking. In 1999, he was awarded the IEEE Signal Processing Society Best Paper Award (Image and Multidimensional Signal

Processing Area) for a paper he co-authored on watermarking. From 1997–1999, he served as Chief Technical Officer of Signafy, Inc, a subsidiary of NEC responsible for the commercialization of watermarking. Between 1996 and 1999, he led the design of NEC's watermarking proposal for DVD video disks and later collaborated with IBM in developing the technology behind the joint "Galaxy" proposal supported by Hitachi, IBM, NEC, Pioneer and Sony. In 1999, he returned to NEC Research Institute as a Research Fellow. He is a senior member of the IEEE, a Fellow of the IEE and a Fellow of the Royal Society for Arts and Manufactures. He is co-editor in chief of the IEE Proc. on Information Security and an associate editor of the IEEE Trans. on Information Forensics and Security. He is co-author of a book entitled "Digital Watermarking" and the co-editor of two books, 'Autonomous Robots Vehicles' and 'Partitioning Data Sets: With Applications to Psychology, Computer Vision and Target Tracking'.