# Novel search algorithms based duplicate detection

**A SASYA SREE,**

**Assistance Professor,**

**Department of Computer Science and Engineering,**

**ANNAM SRINIVASA REDDY**

**Associate Professor,**

**Department of Electronics and Communications Engineering,**

**Siddhartha Institute of Technology and Sciences,**

**Narapally, Hyderabad, Telangana – 500 088.**

## ABSTRACT

Using video fingerprinting, we offer a quick and effective approach for locating duplicate videos in a huge database. To make fingerprints, researchers used the colour layout descriptor, a frame-based descriptor that is further encoded using vector quantization (VQ). We present a novel nonmetric distance measure for identifying the similarity between a query and a database video fingerprint, and show that it beats existing distance measures in detecting duplicates with high accuracy. Existing indexing algorithms are unable to perform successful searches for high-dimensional data when employing a nonmetric distance measure. As a consequence, we create new search algorithms based on pre-computed distances, as well as new dataset pruning processes that allow for faster retrieval times. To conduct our research, we used a collection of 38 000 videos totalling 1600 hours of data. On an Intel Xeon with CPU 2.33 GHz, the replica video is delivered in 0.032 seconds, with a 97.5 percent accuracy, for independent queries with an average time of 60 seconds (about half of the typical database video length).

## 1. INTRODUCTION

The practise of answering questions that are "alike" but not "exactly" "identical" is known as similarity search. It's been used to mimic human-assisted object proximity ranking techniques including image retrieval and time series matching [10]. As a result of rapid advancements in multimedia and network technology, many applications of video databases are becoming increasingly popular, and advanced approaches for encoding, matching, and indexing movies are in high demand.

When dealing with temporal order, frame alignment, gap, and noise, all existing multidimensional sequence similarity measures, such as normalised pair wise distance (Mean distance), Dynamic Time Warping , Longest Common Subsequence (LCSS), or Edit distance

and its weighted extension Edit distance with Real Penalty (ERP), are insufficient in some aspects. Our method has the advantage of collecting visual content not just based on average distance between frame pairings, but also takes temporal order and frame alignment into consideration.

It differs from the widely used Edit distance in that it allows for element permutation and may account for the cross mapping link in the presence of 1. Under Australian copyright law, free-to-air broadcasting content can be saved and utilised for research purposes.

## 2. RELATED WORKS

All similar matching motions are affected by one or more editing actions on a single motion. Standard editing techniques such as time warping, filtering, and motion-warping are all approached in this way. The interface approach falls somewhere between automation and complete user control. Unlike prior mocap synthesis systems, which seek for plausible transitions between motion segments to generate new motion, our method quickly looks for similar motions using a query-by-example paradigm while yet allowing for significant control over the nature of the matching.

Based on their histogram pyramid representations, this vocabulary tree filters those videos that are dissimilar to the query. Second, it's a quick edit-distance-based sequence matching method that avoids redundant frame-by-frame comparisons. With respect to the lengths of the sequences under comparison, this step reduces the quadratic runtime to a linear duration. Experiments with the muscle VCD benchmark show that our method is both effective and efficient. It matches sequences 18 times faster than the original algorithms.
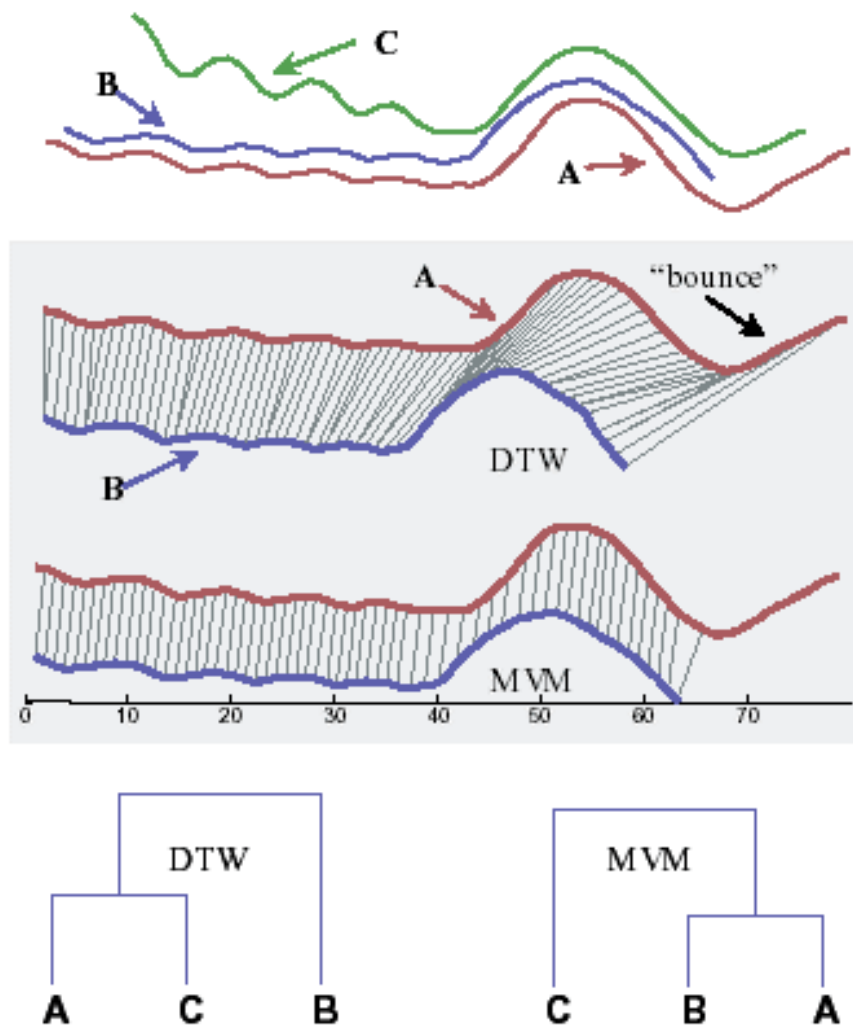
L-dimensional sequences can be seen in text and DNA strings. They do, however, consist of discrete symbols rather than continuous integers, which makes a difference in the feature extraction process. Whole matching is a technique for estimating queries on time sequences, colour images, and even 3-d MRI brain scans. The idea behind all of these methods is to employ f feature extraction procedures to map an entire sequence or image into a point in the (dimensional) feature space, then search for similar sequences or images using spatial access methods. F – iradez is the name of the resultant index, which contains points in feature space.

## 3. PROPOSED SYSTEM

*Filter-And-Refine Search Method*

The closely matched sections of the long sequence are then retrieved, and certain irrelevant subsequences are pruned using a filter-and-refine search approach. To produce a smaller collection of candidates, Maximum Size Matching is used during the filtering stage for each sub-graph constructed by the query and candidate subsequence. Sub-Maximum Similarity Matching is used during the refinement step to find the subsequence with the highest aggregate

score among all candidates, using a comprehensive video similarity model that takes into account visual content, temporal order, and frame alignment information.



**Figure 1. Temporal ordinal measure**

*Generation of Matched Sequence Clip*

Finally, it looks at sub-sampled frame-based matching, taking into account average inter-frame similarity as well as temporal order, frame alignment, gap, and noise in order to achieve reliable identification. Without enumerating all permutations, an efficient heuristic technique is designed to combine the scores of multiple elements for immediately determining the most similar subsequence according to this overall video similarity measure. The matched sequence clip will be generated from the new frames using these calculations and weights.

## 4. RESULT AND DISCUSSION

A bipartite graph (or biography) is a graph whose vertices can be separated into two distinct sets U and V, with each edge connecting a vertex in U to a vertex in V; U and V are independent sets. A bipartite graph, on the other hand, is a graph that does not contain any odd-length cycles. Matching difficulties can be modelled using bipartite graphs. A job matching problem is an

example of a bipartite graph. Assume we have a set P of persons and a set J of occupations, with some people being better suited to certain jobs than others. This can be represented as a bipartite graph (P, J, E). There is an edge in the graph between px and jy if px is acceptable for a specific task jy. The marriage theorem describes bipartite graphs in such a way that they can be perfectly matched. Modern coding theory makes heavy use of bipartite graphs, particularly to decipher code words received from the channel. This can be seen in factor graphs and Tanner graphs, for example. The system's behaviour is further constrained by additional constraints on the nodes and edges. Petri nets make use of the qualities of bipartite directed graphs and other properties to allow mathematical demonstrations of system behaviour while also making simulations of the system simple to build.

### Table 1. Analyzing frames of videos

| Manipulation Type | FMT | DCT | PCA |
|---|---|---|---|
| JPEG | 20 | 40 | 50 |
| Rotation | $10^o$ | $5^o$ | $0^o$ |
| Scaling | 10% | 10% | 0% |

Assume it has a database that stores lengthy time series, which represent continuous observations gathered over relatively long periods of time. As shown in Table 1, each of these time series can comprise observations pertaining to multiple activities that occur in sequence at varying intervals across the time range spanned by that time series. Similarity-based activity retrieval, or the ability to find the best matches for a certain activity of interest, is a useful feature in such databases.

### Table 2. 1-NN classification accuracy.

| Methods | FACE | GUN | LEAF |
|---|---|---|---|
| DTW | 32.43 | 45.00 | 67.38 |
| MVM | 68.21 | 70 | 67.29 |

Dynamic Time Warping aligns the time axis prior to calculating the distance to overcome the problem of time scaling in time series, as shown in Table 2. The total of the distances of their respective elements is the DTW distance between time series. This distance is kept to a minimum by using dynamic programming and associated elements.

## 5. CONCLUSION

Due to content alteration, visually comparable movies may have distinct orderings, resulting in certain inherent cross mappings. This video similarity model, which elegantly achieves a

compromise between discarding temporal order and resolutely holding to temporal order, is especially well-suited to dealing with this issue and may thereby facilitate trustworthy identification. Despite the fact that these studies are just concerned with colour, the proposed strategy obviously takes into account other elements. It was made to test how new features like ordinal signature effect how videos are displayed. Furthermore, user feedback might be utilised to adjust the weighting of each aspect when evaluating video similarity, resulting in a more complete and systematic depiction of the degree of likeness.

## REFERENCES

Yuan, J., Duan, L.-Y., and Xu, C. (2004) "Fast and Robust Short Video Clip Search for Copy Detection," Proc. Fifth IEEE Pacific-Rim Conf. Multimedia (PCM '04), vol. 2, pp. 479-488.

Wang, H., and Sun, H. (2003) "Survey of Compressed-Domain Features Used in Audio-Visual Indexing and Analysis," J. Visual Comm. and Image Representation, vol. 14, no. 2, pp. 150-183.

Smeulders, A.W.M., Worring, M., Santini, S. (2000) "Content-Based Image Retrieval at the End of the Early Years," IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 22, no. 12, pp. 1349-1380, Dec.

Shier, D.R.(2004) "Matchings and Assignments," Handbook of Graph Theory, J.L. Gross and Yellen, Y., eds., pp. 1103-1116, CRC Press.

Shen, H.T., Zhou, X. (2007) "Uqlips: A Real-Time Near-Duplicate Video Clip Detection System," Proc. 33rd Int'l Conf. Very Large Databases (VLDB '07), pp. 1374-1377.

Shen, H.T., Ooi, B.C., Zhou,X. 2005 "Towards Effective Indexing for Very Large Video Sequence Database," Proc. ACM SIGMOD '05, pp. 730-741.

Shao, J., and Huang, Z.(2008), "Batch Nearest Neighbor Search for Video Retrieval," IEEE Trans. Multimedia, vol. 10, no. 3, pp. 409-420.

Pua, K.M., and Miadowicz, J.Z. (2004) "Real Time Repeated Video Sequenc Identification," Computer Vision and Image Understanding, vol. 93, no. 3, pp. 310-327.

Peng, Y., and Ngo, C.W. (2006) "Clip-Based Similarity Measure for Query-Dependent Clip Retrieval and Video Summarization," IEEE Trans. Circuits and Systems for Video Technology, vol. 16, no. 5, pp. 612-627.

Naphade, M.R., Yeung, M.M., Yeo, L. (2000) "A Novel Scheme for Fast and Efficient Video Sequence Matching Using Compact Signatures," Proc. Storage and Retrieval for Image and Video Databases (SPIE '00), pp. 564-572.