# Multi-Resolution K-Means Clustering of Time Series and Application to Images

Michail Vlachos    Jessica Lin    Eamonn Keogh    Dimitrios Gunopulos

Computer Science & Engineering Department
University of California - Riverside
Riverside, CA 92521

{mvlachos, jessica, eamonn, dg}@cs.ucr.edu

## ABSTRACT

*Clustering is vital in the process of condensing and outlining information, since it can provide a synopsis of the stored data. However, the high dimensionality of multimedia data today presents an insurmountable challenge for clustering algorithms. Based on the well known fact that time series and image histograms can both be represented accurately in a lower resolution using orthonormal decompositions, we present an anytime version of the k-Means algorithm. The algorithm works by leveraging the multi-resolution property of wavelets. The dilemma of choosing the initial centers for k-Means is mitigated by assigning the final centers at each approximation level as the initial centers for the subsequent, finer approximation. This process is repeated until the clustering results stabilize or until the finest "approximation" level is reached. In addition to casting k-Means as an anytime algorithm, our approach has two other very desirable properties. We observe that even by working at coarser approximations, the achieved quality is better than the batch algorithm, and even if the algorithm is run to completion, the running time is significantly reduced. We show how this algorithm can be suitably extended to chromatic and textural features extracted from images. Finally, we demonstrate the applicability of this approach on the online image search engine scenario.*

## Keywords
Time Series, k-Means, Anytime Algorithm, Color and Texture Histograms, Image Clustering

## 1. INTRODUCTION

The vast growth of disk technology in the past decade has enabled us to store large multimedia databases, such as audio, video, images, time series, etc. While storage is no longer an impediment, it has become increasingly apparent that an interactive process is needed for humans to efficiently browse through the enormous amount of data. Clustering has proven to be an invaluable tool for distinguishing homogeneous object groups and for producing representatives of the discovered clusters. For example, image clustering is essential for many applications, such as geospatial applications (aerial and satellite photographs), medicine (distinction of tumor images), robot vision (object recognition), online searching (annotation of images), etc. In addition to the high dimensionality of such data, one has to be very careful in extracting features that are coherent with the human perceptual system. Such features include, but are not limited to, color, texture, shape, and location. Essentially, the objective is to select content descriptors that are tailored to a specific application and will lead to discovery of homochromous or homomorphic objects.

Although numerous clustering algorithms have been proposed, the majority of them work in a batch fashion, thus hindering interaction with the end users. Here we address the clustering problem by introducing a novel anytime version of the popular k-Means clustering algorithm [15, 29] based on the wavelet decomposition. Anytime algorithms are valuable for large databases, because results are produced progressively and are refined over time [20]. Their utility for data mining has been documented at length elsewhere [5, 33]. While k-Means algorithm and wavelet decompositions have both been studied extensively in the past, the major novelty of our approach is that it mitigates the problem associated with the choice of initial centers, in addition to providing the functionality of user-interaction.

The algorithm works by leveraging off the multi-resolution property of wavelets [11]. In particular, an initial clustering is performed with a very coarse representation of the data. The results obtained from this "quick and dirty" clustering are used to initialize a clustering at a finer level of approximation. This process is repeated until the clustering results stabilize or until the "approximation" *is* the original "raw" data. The clustering is said to stabilize when the objects do not change membership from the last iteration, or when the change of

membership does not improve the clustering results. Our approach allows the user to interrupt and terminate the process at any level. In addition to casting the k-Means algorithm as an anytime algorithm, our approach has two other very unintuitive properties. The quality of the clustering is often better than the batch algorithm, and even if the algorithm is run to completion, the time taken is typically much less than the time taken by the batch algorithm. We first formulate it as a generic time series problem, since the histograms we extract from images can be well treated as time series. In particular, the high dimensionality and high feature correlation suggest that a time-series algorithm can be suitably applied to image histograms. In this setting, we illustrate the speed-up and scalability of the algorithm, in addition to its improved accuracy. Then we show how this novel algorithm can be applied on histograms extracted from image colors and texture. In addition, we demonstrate how a clustering post-filtering step can enhance the interpretability of the results from on-line image search engines.

The rest of this paper is organized as follows. In Section 2 we review related work, and introduce the necessary background on the wavelet transform and k-Means clustering. In Section 3, we introduce our algorithm. Section 4 contains a comprehensive comparison of our algorithm to classic k-Means on real datasets. In Section 5 we apply the algorithm on image data. Section 6 summarizes our findings and offers suggestions for future work.

## 2. BACKGROUND AND RELATED WORK

Since our work draws on the confluence of clustering, wavelets and anytime algorithms, we provide the necessary background on these areas in this section.

### 2.1 Background on Clustering

One of the most widely used clustering approaches is hierarchical clustering, due to the great visualization power it offers [24]. Hierarchical clustering produces a nested hierarchy of similar groups of objects, according to a pairwise distance matrix of the objects. One of the advantages of this method is its generality, since the user does not need to provide any parameters such as the number of clusters. However, its application is limited to only small datasets, due to its quadratic (or higher order) computational complexity.

A faster method to perform clustering is k-Means [5, 29]. The basic intuition behind k-Means (and a more general class of clustering algorithms known as iterative refinement algorithms) is shown in Table 1:

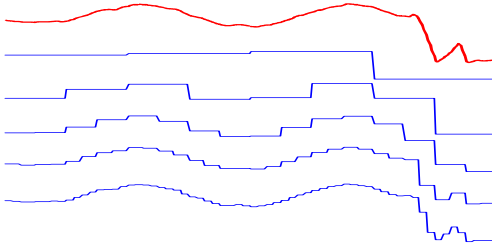| **Algorithm** *k-Means* |
|---|
| 1.  Decide on a value for *k*. |
| 2.  Initialize the *k* cluster centers (randomly, if necessary). |
| 3.  Decide the class memberships of the *N* objects by assigning them to the nearest cluster center. |
| 4.  Re-estimate the *k* cluster centers, by assuming the memberships found above are correct. |
| 5.  If none of the *N* objects changed membership in the last iteration, exit. Otherwise goto 3. |

**Table 1:** An outline of the k-Means algorithm

The k-Means algorithm for N objects has time complexity of O(kNrD) [29], with k the number of clusters specified by the user, r the number of iterations until convergence, and D the dimensionality of the points. The shortcomings of the algorithm are its tendency to favor spherical clusters, and the fact that the knowledge on the number of clusters, k, is required in advance. The latter limitation can be mitigated by placing the algorithm in a loop, and attempting all values of k within a large range. Various statistical tests can then be used to determine which value of k is most parsimonious. Since k-Means is essentiality a hill-climbing algorithm, it is guaranteed to converge on a local but not necessarily global optimum. In other words, the choices of the initial centers are critical to the quality of results. Nevertheless, in spite of these undesirable properties, for clustering large datasets of time-series, k-Means is preferable due to its faster running time.

In order to scale the various clustering methods to massive datasets, one can either reduce the number of objects, N, by sampling [5], or reduce the dimensionality of the objects [1, 6, 14, 27, 31, 39, 40, 24, 25]. In the case of time-series, the objective is to find a representation at a lower dimensionality that preserves the original information and describes the original shape of the time-series data as closely as possible. Many approaches have been suggested in the literature, including the Discrete Fourier Transform (DFT) [1, 14], Singular Value Decomposition [27], Adaptive Piecewise Constant Approximation [25], Piecewise Aggregate Approximation (PAA) [7, 40], Piecewise Linear Approximation [24] and the Discrete Wavelet Transform (DWT) [6, 31]. While all these approaches have shared the ability to produce a high quality reduced-dimensionality approximation of time series, wavelets are unique in that their representation of data is intrinsically multi-resolution. This property is critical to our proposed algorithm and will be discussed in detail in the next section.

Although we choose the Haar wavelet for this work, the algorithm can generally utilize any wavelet basis. The preference for the Haar wavelet is mainly based on its simplicity and its wide usage in the data mining community.

## 2.2 Background on Wavelets

Wavelets are mathematical functions that represent data or other functions in terms of the averages and differences of a prototype function, called the analyzing or mother wavelet [11]. In this sense, they are similar to the Fourier transform. One fundamental difference is that wavelets are localized in time. In other words, some of the wavelet coefficients represent small, local subsections of the data being studied, as opposed to Fourier coefficients, which always represent global contributions to the data. This property is very useful for multi-resolution analysis of data. The first few coefficients contain an overall, coarse approximation of the data; additional coefficients can be perceived as "zooming-in" to areas of high detail. Figure 1 illustrates this idea.



**Figure 1.** The Haar Wavelet can represent data at different levels of resolution. Above we see a raw time series, with increasing faithful wavelet approximations below.

The Haar Wavelet decomposition works by averaging two adjacent values on the time series function at a given resolution to form a smoothed, lower-dimensional signal, and the resulting coefficients are simply the differences between the values and their averages [6]. The coefficients can also be computed by averaging the differences between each pair of adjacent values. The coefficients are crucial for reconstructing the original sequence, as they store the detailed information lost in the smoothed signal. For example, suppose we have a time series data T = (2 8 1 5 9 7 2 6). Table 2 shows the decomposition at different resolutions. As a result, the Haar wavelet decomposition is the collection of the coefficients at all resolutions, with the overall average being its first component: (5 –1 1 2 –3 –2 1 –2). It is clear to see that the decomposition is completely reversible and the original sequence can be reconstructed from the coefficients. For example, to get the signal of the second level, simply compute 5 ± (-1) = (4, 6).

| Resolution | Averages | Differences (coefficients) |
|---|---|---|
| 8 | (2 8 1 5 9 7 2 6) | |
| 4 | (5 3 8 4) | (-3 –2 1 –2) |
| 2 | (4 6) | (1 2) |
| 1 | 5 | -1 |

**Table 2**. Haar Wavelet Decomposition on time series (2 8 1 5 9 7 2 6)

Recently there has been an explosion of interest in using wavelets for time series data mining. Researchers have introduced several non-Euclidean, wavelet-based distance measures [22, 36]. Chan and Fu [6] have demonstrated that Euclidean distance indexing with wavelets is competitive to Fourier-based techniques [14].

## 2.3 Background on Anytime Algorithms

Anytime algorithms are algorithms that trade execution time for quality of results [20]. In particular, an anytime algorithm always has a best-so-far answer available, and the quality of the answer improves with execution time. The user may examine this answer at any time, and choose to terminate the algorithm, temporarily suspend the algorithm, or allow the algorithm to run to completion.

The usefulness of anytime algorithms for data mining has been extensively documented [5, 33]. Suppose a batch version of an algorithm takes a week to run (not an implausible scenario in data mining massive data sets). It would be highly desirable to implement the algorithm as an anytime algorithm. This would allow a user to examine the best current answer after an hour or so as a "sanity check" of all assumptions and parameters. As a simple example, suppose the user had accidentally set the value of $k$ to 50 instead of the desired value of 5. Using a batch algorithm the mistake would not be noted for a week, whereas using an anytime algorithm the mistake could be noted early on and the algorithm restarted with little cost.

The motivating example above could have been eliminated by user diligence! More generally, however, data mining algorithms do require the user to make choices of several parameters, and an anytime implementation of *k-Means* would allow the user to interact with the entire data mining process in a more efficient way.

## 2.4 Related Work

Bradley et. al. [5] suggest a generic technique for scaling the *k-Means* clustering algorithms to large databases by attempting to identify regions of the data that are compressible, that must be retained in main memory, and regions that may be discarded. However, the generality of the method contrasts with our algorithm's explicit exploitation of the structure of the data type of interest.

Our work is more similar in spirit to the dynamic time warping similarity search technique introduced by Chu et. al. [7]. The authors speed up linear search by examining the time series at increasingly finer levels of approximation.

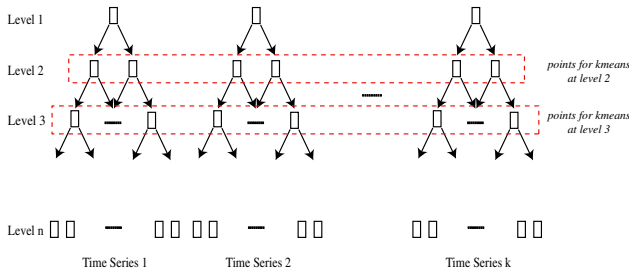## 3. OUR APPROACH – THE I-kMEANS ALGORITHM

As noted in Section 2.1, the complexity of the k-Means algorithm is *O(kNrD)*, where *D* is the dimensionality of data points (or the length of a sequence, as in the case of

time-series). For a dataset consisting of long time-series, the *D* factor can burden the clustering task significantly. This overhead can be alleviated by reducing the data dimensionality.

Another major drawback of the k-Means algorithm derives from the fact that the clustering quality is greatly dependant on the choice of initial centers (i.e., line 2 of Table 1). As mentioned earlier, the k-Means algorithm guarantees local, but not necessarily global optimization. Poor choices of the initial centers can degrade the quality of clustering solution and result in longer execution time (See [15] for an excellent discussion of this issue). Our algorithm addresses these two problems associated with k-Means, in addition to offering the capability of an anytime algorithm, which allows the user to interrupt and terminate the program at any stage.

We propose using the wavelet decomposition to perform clustering at increasingly finer levels of the decomposition, while displaying the gradually refined clustering results periodically to the user. We compute the Haar Wavelet decomposition for all time-series data in the database. The complexity of this transformation is linear to the dimensionality of each object; therefore, the running time is reasonable even for large databases. The process of decomposition can be performed off-line, and the time series data can be stored in the Haar decomposition format, which takes the same amount of space as the original sequence. One important property of the decomposition is that it is a lossless transformation, since the original sequence can always be reconstructed from the decomposition.

Once we compute the Haar decomposition, we perform the k-Means clustering algorithm, starting at the second level (each object at level *i* has $2^{(i-1)}$ dimensions) and gradually progress to finer levels. Since the Haar decomposition is completely reversible, we can reconstruct the approximate data from the coefficients at any level and perform clustering on these data. We call the new clustering algorithm ***I-kMeans***, where *I* stands for "interactive." Figure 2 illustrates this idea.



**Figure 2**. k-Means is performed on each level on the reconstructed data from the Haar wavelet decomposition, starting with the second level.

The intuition behind this algorithm originates from the observation that the general shape of a time series sequence can often be approximately captured at a lower resolution. As shown in Figure 1, the shape of the time series is well preserved, even at very coarse approximations. Because of this desirable property of wavelets, clustering results typically stabilize at a low resolution, thus saving time by eliminating the need to run at full resolution (the raw data).

On every level except for the starting level (i.e. level 2), which uses random initial centers, the initial centers are selected based on the final centers from the previous level. More specifically, the final centers computed at the end of level *i* will be used as the initial centers on level *i+1*. Since the length of the data reconstructed from the Haar decomposition doubles as we progress to the next level, we project the centers computed at the end of level *i* onto level *i+1* by doubling each coordinate of the centers. This way, they match the dimensionality of the points on level *i+1*. For example, if one of the final centers at the end of level 2 is (0.5, 1.2), then the initial center used for this cluster on level 3 is (0.5, 0.5, 1.2, 1.2). This approach mitigates the dilemma associated with the choice of initial centers, which is crucial to the quality of clustering results [15]. It also contributes to the fact that our algorithm often produces better clustering results than the k-Means algorithm. The pseudo-code of the algorithm is provided in Table 3.

| **Algorithm** I-*k*Means |
|---|
| 1.  Decide on a value for *k*. |
| 2.  Initialize the *k* cluster centers (randomly, if necessary). |
| 3.  Run the *k*-Means algorithm on the level$_i$ representation of the data |
| 4.  Use final centers from level$_i$ as initial centers for level$_{i+1}$. This is achieved by projecting the k centers returned by *k*-Means algorithm for the $2^i$ space in the $2^{i+1}$ space. |
| 5.  If none of the *N* objects changed membership in the last iteration, exit. Otherwise goto 3. |

**Table 3:** An outline of the I-kMeans algorithm

The algorithm achieves the speed-up by doing the vast majority of reassignments (Line 3 in Table 1), at the lower resolutions, where the costs of distance calculations are considerably lower. As we gradually progress to finer resolutions, we already start with good initial centers. Therefore, the number of iterations *r* until convergence will typically be much lower.

The I-kMeans algorithm allows the user to monitor the quality of clustering results as the program executes. The user can interrupt the program at any level, or wait until the execution terminates once the clustering results stabilize. One surprising and highly desirable finding from the experimental results (as shown in the next section), is that even if the program is run to completion (until the last

level, with full resolution), the total execution time is generally less than that of clustering on the raw data.

# 4. EXPERIMENTAL EVALUATION on TIME SERIES

To show that our approach is superior to the k-Means algorithm for clustering time series, we performed a series of experiments on publicly available real datasets. For completeness, we ran the I-kMeans algorithm for all levels of approximation, and recorded the cumulative execution time and clustering accuracy at each level. In reality, however, the algorithm stabilizes in early stages and can automatically terminate much sooner. We compare the results with that of k-Means on the original data. Since both algorithms start with random initial centers, we execute each algorithm 100 times with different centers. However, for consistency we ensure that for each execution, both algorithms are seeded with the same set of initial centers. After each execution, we compute the error (more details will be provided in Section 4.2) and the execution time on the clustering results. We compute and report the averages at the end of each experiment. By taking the average, we achieve better objectiveness than taking the best (minimum), since in reality, we would not have the knowledge of the correct clustering results, or the "oracle," to compare against our results (as in the case with one of our test datasets).
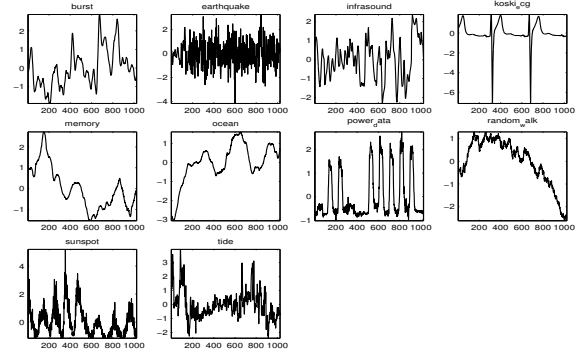
## 4.1 Datasets and Methodology

We tested on two publicly available, real datasets. The dataset cardinalities range from 1,000 to 8,000. The length of each time series has been set to 512 on one dataset, and 1024 on the other. Each time series is z-normalized to have mean value of 0 and standard deviation of 1.

• **JPL:** This dataset consists of readings from various inertial sensors from Space Shuttle mission STS-57. The data is particularly appropriate for our experiments since the use of redundant backup sensors means that some of the data is very highly correlated. In addition, even sensors that measure orthogonal features (i.e. the X and Y axis) may become temporarily correlated during a particular maneuver; for example, a "roll reversal" [13]. Thus, the data has an interesting mixture of dense and sparse clusters. To generate datasets of increasingly larger cardinalities, we extracted time series of length 512, at random starting points of each sequence from the original data pool.
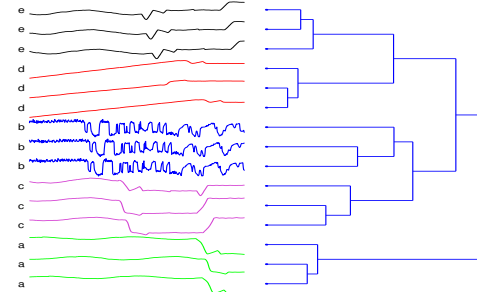
• **Heterogeneous:** This dataset is generated from a mixture of 10 real time series data from the UCR Time Series Data Mining Archive [26]. Figure 3 shows the 10 time-series we use as seeds. We produced variations of the original patterns by adding small time warping (2-3% of the series length), and interpolated Gaussian noise. Gaussian

noisy peaks are interpolated using splines to create smooth random variations.

In the Heterogeneous dataset, we know that the number of clusters (*k*) is 10. However, for the JPL dataset, we lack this information. Finding k is an open problem for the k-Means algorithm and is out of scope of this paper. To determine the optimal *k* for k-Means, we attempt different values of *k*, ranging from 2 to 8. Nonetheless, our algorithm out-performs the k-Means algorithm regardless of *k*. In this paper we only show the results with *k* equals to 5. Figure 4 shows that our algorithm produces the same results as does the hierarchical clustering algorithm, which is generally more costly.



**Figure 3**. Real time series data from UCR Time Series Data Mining Archive. We use these time series as seeds to create our Heterogeneous dataset.



**Figure 4**. On the left-hand side, we show three instances from each cluster discovered by the I-kMeans algorithm. We can visually verify that our algorithm produces intuitive results. On the right-hand side, we show that hierarchical clustering (using average linkage) discovers the exact same clusters. However, hierarchical clustering is more costly than our algorithm.

## 4.2 Error of Clustering Results

In this section we compare the clustering quality for the I-kMeans and the classic k-Means algorithm.

Since we generated the heterogeneous datasets from a set of given time series data, we have the knowledge of correct clustering results in advance. In this case, we can simply compute the clustering accuracy by summing up the number of correctly classified objects for each cluster *c* and then dividing by the dataset cardinality. This is done by the use of a confusion matrix. Note the accuracy computed
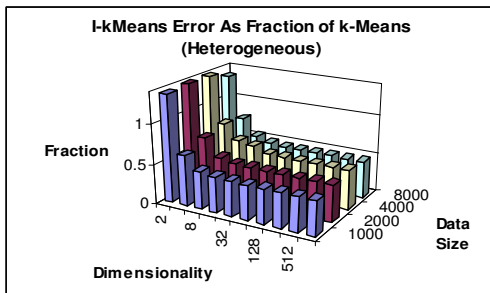
here is equivalent to "recall," and the error rate is simply (1-accuracy).

The error is computed at the end of each level. However, it's worth mentioning that in reality, the correct clustering information would not be available in advance. The incorporation of such known results in our error calculation merely serves the purpose of demonstrating the quality of both algorithms.
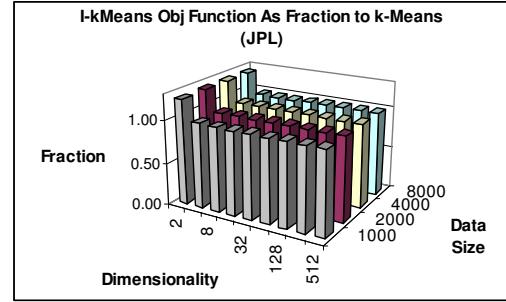
For the JPL dataset, we do not have prior knowledge of correct clustering information (which conforms more closely to real-life cases). Lacking this information, we cannot use the same evaluation to determine the error.

Since the k-Means algorithm seeks to optimize the objective function by minimizing the sum of squared intra-cluster errors, we can evaluate the quality of clustering by using the objective functions. However, since the I-kMeans algorithm involves data with smaller dimensionality except for the last level, we have to compute the objective functions on the raw data, in order to compare with the k-Means algorithm. We show that the objective functions obtained from the I-kMeans algorithm are better than those from the k-Means algorithm. The results are consistent with the work of [8], in which the authors show that dimensionality reduction reduces the chances of the algorithm being trapped in a local minimum. Furthermore, even with the additional step of computing the objective functions from the original data, the I-kMeans algorithm still takes less time to execute than the k-Means algorithm.

In figures 5 and 6, we show the errors/objective functions from the I-kMeans algorithm as a fraction of those obtained from the k-Means algorithm. As we can see from the plots, our algorithm stabilizes at early stages and consistently results in smaller error than the classic k-Means algorithm.
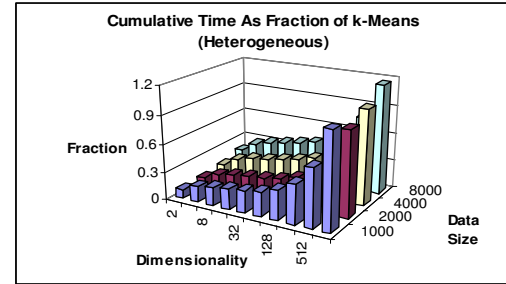
**Figure 5.** Error of I-kMeans algorithm on the Heterogeneous dataset, presented as fraction of the error from the k-Means algorithm. Our algorithm results in smaller error than the k-Means after the second stage (i.e. 4 dimensions), and stabilizes typically after the third stage (i.e. 8 dimensions).
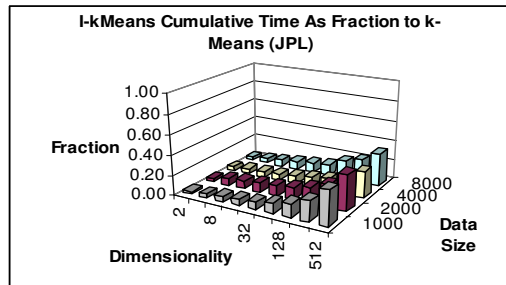
**Figure 6.** Objective functions of I-kMeans algorithm on the JPL dataset, presented as fraction of error from the k-Means algorithm. Again, our algorithm results in smaller objective functions (i.e. better clustering results) than the k-Means, and stabilizes typically after the second stage (i.e. 4 dimensions).

### 4.3 Running Time

In Figures 7 and 8, we present the cumulative running time for each level on the I-kMeans algorithm as a fraction to the k-Means algorithm. The cumulative running time for any level i is the total running time from the starting level (level 2) to level i. In most cases, even if the I-kMeans algorithm is run to completion, the total running time is still less than that of the k-Means algorithm. We attribute this improvement to the good choices of initial centers for successive levels after the starting level, since they result in very few iterations until convergence. Nevertheless, we have already shown in the previous section that the I-kMeans algorithm finds the best result in relatively early stage and does not need to run through all levels.
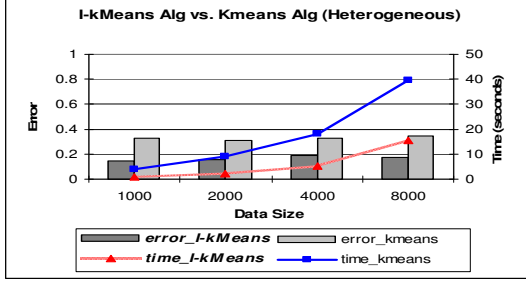
**Figure 7**. Cumulative running time for the Heterogeneous dataset. Our algorithm typically cuts the running time by half as it does not need to run through all levels to retrieve the best results.
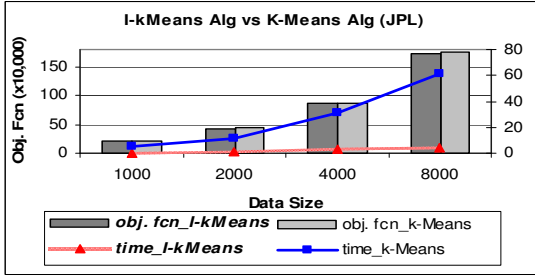
**Figure 8**. Cumulative running time for the JPL dataset. Our algorithm typically takes only 30% of time.

## 4.4  I-kMeans Algorithm vs. k-Means Algorithm

In this section (Fig. 9 and 10), rather than showing the error/objective function on each level, as in Section 4.2, we present only the error/objective function returned by the I-kMeans algorithm when it stabilizes or, in the case of JPL dataset, outperforms the k-Means algorithm in terms of the objective function. We also present the time taken for the I-kMeans algorithm to stabilize. We compare the results to those of the k-Means algorithm. From the figures we can observe that our algorithm achieves better clustering accuracy at significantly faster response time.
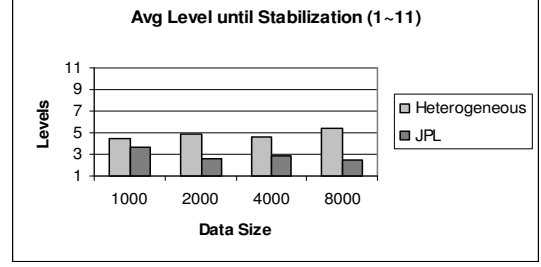


**I-kMeans Alg vs. Kmeans Alg (Heterogeneous)**

**Figure 9**. The I-kMeans algorithm is highly competitive with the k-Means algorithm. The errors and execution time are significantly smaller.



**I-kMeans Alg vs K-Means Alg (JPL)**

**Figure 10**. I-kMeans vs. k-Means algorithms in terms of objective function and running time for JPL dataset. Our algorithm outperforms the k-Means algorithm. The running time remains small for all data sizes because the algorithm terminates at very early stages.

Figure 11 shows the average level where the I-kMeans algorithm stabilizes or, in the case of JPL, outperforms the k-Means algorithm in terms of objective function. Since the length of the time series data is 1024 in the Heterogeneous dataset, there are 11 levels. Note that the JPL dataset has only 10 levels since the length of the time series data is only 512. We skip level 1, in which the data has only one dimension (the average of the time series) and is the same for all sequences, since the data has been normalized (zero mean). Each level $i$ has $2^{(i-1)}$ dimensions. From the plot we can see that our algorithm generally stabilizes at levels 3-6 for the Heterogeneous dataset and at levels 2-4 for the JPL dataset. In other words, the I-kMeans algorithm operates on data with a maximum dimensionality of 32 and 8, respectively, rather than 1024 and 512.



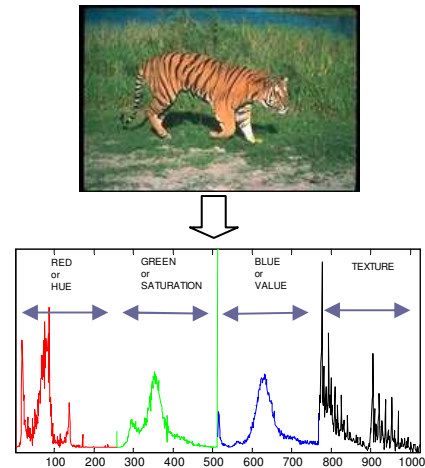**Avg Level until Stabilization (1~11)**

**Figure 11.** Average level until stabilization. The algorithm generally stabilizes between level 3 and level 6 for Heterogeneous dataset, and between level 2 and 4 for the JPL dataset.

## 5.  APPLICATION TO IMAGES

In this section we provide preliminary results regarding the applicability of our approach for images.

In our case study, we utilize two descriptors: color and texture, and we compare two color representations, the RGB and HSV chromatic spaces. In general, the RGB components are highly correlated and do not necessarily coincide with the human visual system. We validate this observation empirically, since our experiments show that the HSV model results in clusters of better quality. From each image, color histograms of size 786 (256 for each component) are extracted for each color model.

For the texture we apply gabor-wavelet filters [38] at various scales and orientations, which results in an additional vector of size 256. The color and texture information are concatenated and the final vector of size 1024 is treated as a time-series. Therefore, the method that we proposed in the previous sections can be applied unaltered to images. Notice that this representation can very easily facilitate a weighted scheme, where according to the desired matching criteria we may favor the color or the texture component.



**Figure 12.** An example of the image vector that is extracted for the purposes of our experiments.

A well known observation in the image retrieval community is that the use of large histograms suffers from the "curse of dimensionality" [19]. Our method is therefore particularly applicable in this scenario, since the wavelet (or any other) decomposition can help reduce the dimensionality effectively. It should also be noted that the manipulation of the image features as time-series is perfectly valid. The characteristics of 'smoothness' and auto-correlation of the time-series are evident here, since adjacent bins in the histograms have very similar values.
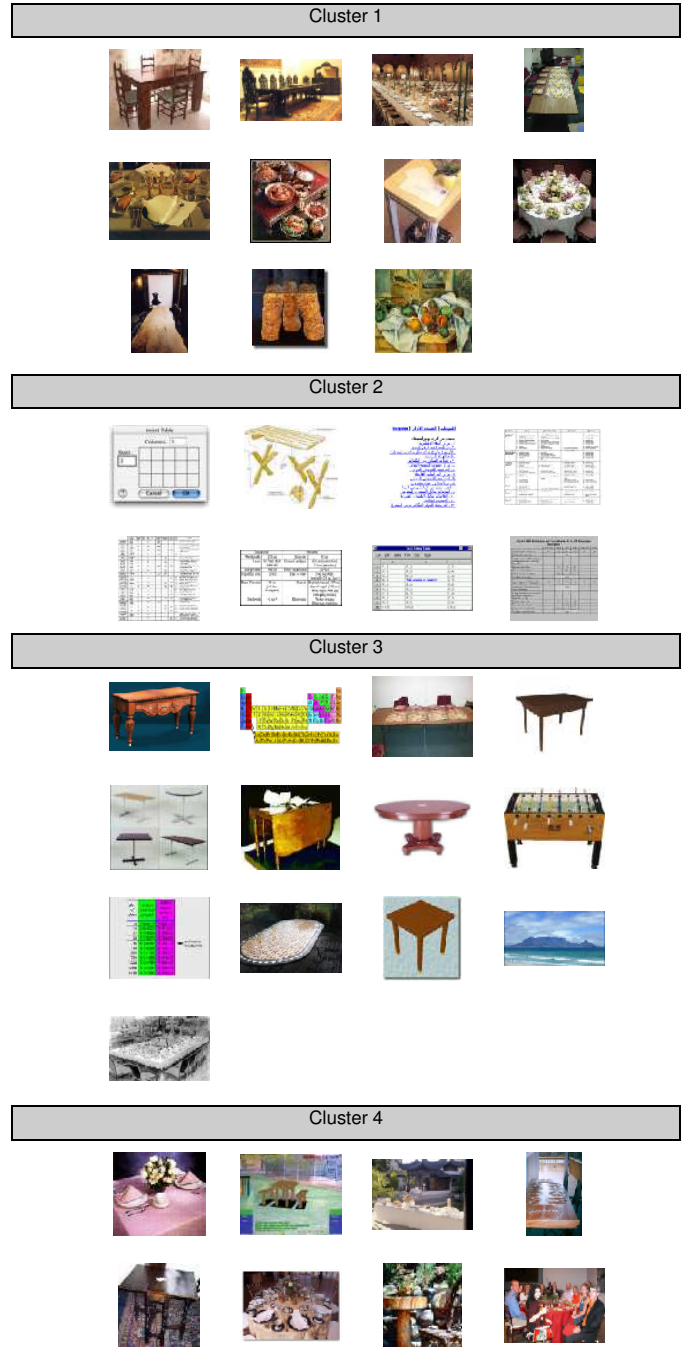
We demonstrate the usefulness of our approach on two scenarios. First, we perform clustering on a subset of images taken from the Corel database. The dataset consists of 6 clusters, each with 100 images. The clusters represent diverse topics: *diving, fireworks, tigers, architecture, snow and stained glass.* This is by no means a large dataset; however, it can very well serve as an indicator of the speed and the accuracy of the I-kMeans algorithm. We compare the algorithm with k-Means, and with hierarchical clustering using Ward's linkage method. Our approach achieves the highest recall rate and is also the fastest method. Since the performance of k-Means greatly depends on the initial choice of centers, we also provided 'good' initial centers by picking one image randomly from each cluster. Even in this scenario our algorithm was significantly faster and more accurate. The results, shown in Table 4, are generated using the L1 distance metric. Although we also tested on the Euclidean distance, we discover that using L1 metric generally results in better accuracy for our test cases.

|  | Error | | | Time (sec) | | |
|---|---|---|---|---|---|---|
|  | Hier. | k-Means | I-kMeans | Hier. | k-Means | I-kMeans |
| **RGB** | **0.31** | 0.37 | 0.37 | 50.7 | 11.18 | **4.94** |
| **HSV** | 0.36 | **0.26** | **0.26** | 49.8 | 11.12 | **6.82** |
| **HSV+ TEXTR** | 0.24 | 0.25 | **0.19** | 55.5 | 11.33 | **3.14** |
| **HSV+ TEXTR+ good centers** | N/A | 0.22 | **0.17** | N/A | 8.26 | **2.89** |

**Table 4.** Clustering Error and Running Time averaged over 100 runs. Our method is even more accurate than hierarchical clustering, which requires almost an order of magnitude more time.
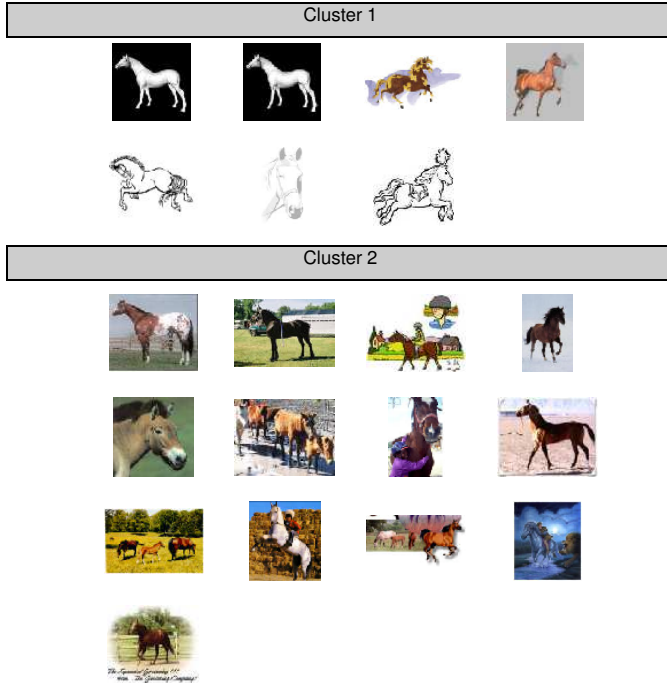
For the second paradigm we used a real-world example from the image search feature of Google. All online image search engines gather information based only on keywords. Therefore, the image query results can be from very diverse disciplines and have very different content. For example, if one searches for '*bass*', the result would be a mixture of images about 'fish' and 'musical instruments.' In this specific case, a post-filtering of the matches using *texture features* can create 2 separate clusters, and as a consequence, lead to a more intuitive presentation of the results. In order for such a post-filtering step to become a

reality, it is obvious that one must utilize an extremely lightweight algorithm. We posed several queries on Google and we grouped the results into clusters. We present representative results for two queries, the words "table" and "horse". We retrieved the first 2 image pages returned (40 images in total) for the query *table* and the first one (20 images) for the word *horse*. The images are passed on to our algorithm and are clustered into four groups and two groups, respectively. Figure 13 and 14 present the results.



**Figure 13.** Clustering results on the image query '*table*' posed at Google.

For the query "table" we observe that we can distinguish almost perfectly between the 'array' and the 'object' tables returned by Google, while for the word "horse" there is an obvious separation between the hand-drawn pictures and the photographic images. These experiments suggest that online image search could be augmented by a clustering step, in the spirit of the well-known "Scatter/Gather" framework [21].



**Figure 14.** Clustering results on the image query '*horse*' posed at Google

The results could be improved by using relevance feedback. In addition, compact histogram representations, as well as the use of more robust distance functions such as the *dLog* distance proposed in [35], could further boost the performance for the proposed algorithm.

# 6. CONCLUSIONS AND FUTURE WORK

We have presented an approach to perform incremental clustering of time series at various resolutions using the Haar wavelet transform. Using *k-Means* as our clustering algorithm, we reuse the final centers at the end of each resolution as the initial centers for the next level of approximation. This approach mitigates the dilemma associated with the choices of initial centers for k-Means and significantly improves the execution time and clustering quality. Our experimental results indicate that this approach yields faster execution time than the traditional *k-Means* approach, in addition to improving the clustering quality of the algorithm. The anytime algorithm stabilizes at very early stages, eliminating the needs to operate on high dimensionality. In addition, the anytime algorithm allows the user to terminate the program at any stage.

Since image histograms extracted from colors and textures can be suitably treated as time series, we further demonstrate the efficacy of our algorithm on the image data. In future work we plan to investigate the following:

- Extending our algorithm to other iterative refinement clustering techniques, especially the EM-algorithm.

- Examining the possibility of re-using the results (i.e. objective functions that determine the quality of clustering results) from the previous stages to eliminate the need to re-compute all the distances.

- Generalizing our approach to a broader class of algorithms and decompositions. For example, even though we have used the wavelet decomposition in this work, our algorithm can be easily generalized to Fourier coefficients as well.

## 7. REFERENCES

[1] Agrawal, R., Faloutsos, C. & Swami, A. (1993). Efficient similarity search in sequence databases. In proceedings of the 4[th] Int'l Conference on Foundations of Data Organization and Algorithms. Chicago, IL, Oct 13-15. pp 69-84.

[2] André-Jönsson, H. & Badal. D. (1997). Using signature files for querying time-series data. In proceedings of Principles of Data Mining and Knowledge Discovery, 1[st] European Symposium. Trondheim, Norway, Jun 24-27. pp 211-220.

[3] Bay, S. (1999). UCI Repository of Kdd databases [http://kdd.ics.uci.edu/]. Irvine, CA: University of California, Department of Information and Computer Science

[4] Bozkaya, T., Yazdani, N. & Ozsoyoglu, Z. M. (1997). Matching and indexing sequences of different lengths. In proceedings of the 6[th] Int'l Conference on Information and Knowledge Management. Las Vegas, NV, Nov 10-14. pp 128-135.

[5] Bradley, P., Fayyad, U., & Reina, C. (1998). Scaling clustering algorithms to large databases. In proceedings of the 4[th] Int'l Conference on Knowledge Discovery and Data Mining. New York, NY, Aug 27-31. pp 9-15.

[6] Chan, K. & Fu, A. W. (1999). Efficient time series matching by wavelets. In proceedings of the 15[th] IEEE Int'l Conference on Data Engineering. Sydney, Australia, Mar 23-26. pp 126-133.

[7] Chu, S., Keogh, E., Hart, D., Pazzani, M. (2002). Iterative Deepening Dynamic Time Warping for Time Series. In proceedings of the 2[nd] SIAM International Conference on Data Mining.

[8] Ding, C., He, X., Zha, H. & Simon, H. (2002). Adaptive Dimension Reduction for Clustering High Dimensional Data. In proceedings of the 2[nd] IEEE International Conference on Data Mining. Dec 9-12. Maebashi, Japan. pp 147-154.

[9] Das, G., Gunopulos, D. & Mannila, H. (1997). Finding similar time series. In proceedings of Principles of Data

Mining and Knowledge Discovery, 1[st] European Symposium. Trondheim, Norway, Jun 24-27. pp 88-100.

[10] Das, G., Lin, K., Mannila, H., Renganathan, G. & Smyth, P. (1998). Rule discovery from time series. In proceedings of the 4[th] Int'l Conference on Knowledge Discovery and Data Mining. New York, NY, Aug 27-31. pp 16-22.

[11] Daubechies, I. (1992). Ten Lectures on Wavelets. Number 61 in CBMS-NSF regional conference series in applied mathematics, Society for Industrial and Applied Mathematics, Philadelphia.

[12] Debregeas, A. & Hebrail, G. (1998). Interactive interpretation of kohonen maps applied to curves. In proceedings of the 4[th] Int'l Conference of Knowledge Discovery and Data Mining. New York, NY, Aug 27-31. pp 179-183.

[13] Dumoulin, J. (1998). NSTS 1988 News Reference Manual. http://www.fas.org/spp/civil/sts/

[14] Faloutsos, C., Ranganathan, M. & Manolopoulos, Y. (1994). Fast subsequence matching in time-series databases. In proceedings of the ACM SIGMOD Int'l Conference on Management of Data. Minneapolis, MN, May 25-27. pp 419-429.

[15] Fayyad, U., Reina, C. & Bradley. P (1998). Initialization of iterative refinement clustering algorithms. In Proceedings of the 4[th] International Conference on Knowledge Discovery and Data Mining. New York, NY, Aug 27-31. pp 194-198.

[16] Gavrilov, M., Anguelov, D., Indyk, P. & Motwani, R. (2000). Mining the stock market: which measure is best? In proceedings of the 6[th] ACM Int'l Conference on Knowledge Discovery and Data Mining. Boston, MA, Aug 20-23. pp 487-496.

[17] Ge, X. & Smyth, P. (2000). Deformable markov model templates for time-series pattern matching. In proceedings of the 6[th] ACM SIGKDD Int'l Conference on Knowledge Discovery and Data Mining. Boston, MA, Aug 20-23. pp 81-90.

[18] Geurts, P. (2001). Pattern extraction for time series classification. In proceedings of Principles of Data Mining and Knowledge Discovery, 5[th] European Conference. Freiburg, Germany, Sept 3-5. pp 115-127.

[19] Gibson, S. & Harvey, R. (2001). Analyzing and Simplifying Histograms Using Scale-Trees. In proceedings of the 11[th] Int'l Conference on Image Analysis and Processing. Palermo, Italy. Sept 26-28.

[20] Grass, J. & Zilberstein, S. (1996). Anytime algorithm development tools. Sigart Artificial Intelligence. Vol 7, No. 2, April. ACM Press.

[21] Hearst, M. A. & Pedersen, J. O. (1996). Reexamining the Cluster Hypothesis: Scatter/Gatter on Retrieval Results. In proceedings of the 19[th] Annual International ACM SIGIR Conference on Research and Development in Information Retrieval. Zurich, Switzerland. Aug 18-22.

[22] Huhtala, Y., Kärkkäinen, J. & Toivonen, H. (1999). Mining for similarities in aligned time series using wavelets. Data Mining and Knowledge Discovery: Theory, Tools, and Technology, SPIE Proceedings Series, Vol. 3695. Orlando, FL, Apr. pp 150-160.

[23] Kalpakis, K., Gada, D. & Puttagunta, V. (2001). Distance measures for effective clustering of ARIMA time-series. In proceedings of the IEEE Int'l Conference on Data Mining. San Jose, CA, Nov 29-Dec 2. pp 273-280.

[24] Keogh, E. & Pazzani, M. (1998). An enhanced representation of time series which allows fast and accurate classification, clustering and relevance feedback. In proceedings of the 4[th] Int'l Conference on Knowledge Discovery and Data Mining. New York, NY, Aug 27-31. pp 239-241.

[25] Keogh, E., Chakrabarti, K., Pazzani, M. & Mehrotra, S. (2001). Locally adaptive dimensionality reduction for indexing large time series databases. In proceedings of ACM SIGMOD Conference on Management of Data. Santa Barbara, CA. pp 151-162.

[26] Keogh, E. & Folias, T. (2002). The UCR Time Series Data Mining Archive [http://www.cs.ucr.edu/~eamonn/TSDMA/ index.html].

[27] Korn, F., Jagadish, H. & Faloutsos, C. (1997). Efficiently supporting ad hoc queries in large datasets of time sequences. In proceedings of the ACM SIGMOD Int'l Conference on Management of Data. Tucson, AZ, May 13-15. pp 289-300.

[28] Li, C., Yu, P. S. & Castelli, V. (1998). MALM: a framework for mining sequence database at multiple abstraction levels. In proceedings of the 7[th] ACM CIKM Int'l Conference on Information and Knowledge Management. Bethesda, MD, Nov 3-7. pp 267-272.

[29] McQueen, J. (1967). Some methods for classification and analysis of multivariate observation. L. Le Cam and J. Neyman (Eds.), 5[th] Berkeley Symp. Math. Stat. Prob., 1,pp 281-297.

[30] Ng, Y., Jordan, M. & Weiss, Y.(2001). On Spectral clustering: analysis and an algorithm.. Advances in Neural Information Processing Systems 14.

[31] Popivanov, I. & Miller, R. J. (2002). Similarity search over time series data using wavelets. In proceedings of the 18[th] Int'l Conference on Data Engineering. San Jose, CA, Feb 26-Mar 1. pp 212-221.

[32] Qu, Y., Wang, C. & Wang, X. S. (1998). Supporting fast search in time series for movement patterns in multiples scales. In proceedings of the 7[th] ACM CIKM Int'l Conference on Information and Knowledge Management. Bethesda, MD, Nov 3-7. pp 251-258.

[33] Smyth, P,. Wolpert, D. (1997). Anytime exploratory data analysis for massive data sets. In proceedings of the 3[rd] Int'l Conference on Knowledge Discovery and Data Mining. Newport Beach, CA. pp 54-60

[34] Shahabi, C., Tian, X. & Zhao, W. (2000). TSA-tree: a wavelet based approach to improve the efficiency of multi-level surprise and trend queries. In proceedings of the 12[th] Int'l Conference on Scientific and Statistica Database Management. Berlin, Germany, Jul 26-28. pp 55-68.

[35] Stehling, R. O., Nascimento, M. A. & Falcão, A. X. (2002). A Compact and Efficient Image Retrieval Approach Based on Border/Interior Pixel Classification. In proceedings of the

ACM Intl. Conf. on Information and Knowledge Management. McLean, VA. Nov 4-9.

[36] Struzik, Z. & Siebes, A. (1999). The Haar wavelet transform in the time series similarity paradigm. In proceedings of Principles of Data Mining and Knowledge Discovery, 3rd European Conference. Prague, Czech Republic, Sept 15-18. pp 12-22.

[37] Wijk, J.J. van, E. van Selow.(1999). Cluster and calendar-based visualization of time series data. In Proceedings 1999 IEEE Symposium on Information Visualization (InfoVis'99), IEEE Computer Society, pp 4-9.

[38] Wu, P., Manjunath, B. S., Newsam, S. & Shin, H. D. (2000). A Texture Descriptor for Browsing and Similarity Retrieval. *Journal of Signal Processing: Image Communication.* Volume 16, Issue 1-2, page 33-43.

[39] Wu, Y., Agrawal, D. & El Abbadi, A. (2000). A comparison of DFT and DWT based similarity search in time-series databases. In proceedings of the 9th ACM CIKM Int'l Conference on Information and Knowledge Management. McLean, VA, Nov 6-11. pp 488-495.

[40] Yi, B. & Faloutsos, C. (2000). Fast time sequence indexing for arbitrary lp norms. In proceedings of the 26th Int'l Conference on Very Large Databases. Cairo, Egypt, Sept 10-14. pp 385-394.l Database Management. Berlin, Germany, Jul 26-28. pp 55-68.