

TIME-SERIES SEGMENTATION

TO ALL WHOM IT MAY CONCERN:

BE IT KNOWN THAT (1) PAUL M. BAGGENSTOSS, citizen of the United States of America, employee of the United States Government, resident of Newport, County of Newport, State of Rhode Island has invented certain new and useful improvements entitled as set forth above of which the following is a specification:

DISTRIBUTION STATEMENT A
Approved for Public Release
Distribution Unlimited

JAMES M. KASISCHKE, ESQ.
Reg. No. 36562
Naval Undersea Warfare Center
Division, Newport
Newport, RI 02841-1708
TEL: 401-832-4763
FAX: 401-832-1231

20030618 068



23523

PATENT TRADEMARK OFFICE

I hereby certify that this correspondence is being deposited with the U.S. Postal Service as U.S EXPRESS MAIL, Mailing Label No. EL578538533US in envelope addressed to: Assistant Commissioner for Patents, Washington, DC 20231 on 8/31/2001

(DATE OF DEPOSIT)

James M. Kasischke
APPLICANT'S ATTORNEY

8/31/2001
DATE OF SIGNATURE



DEPARTMENT OF THE NAVY
OFFICE OF COUNSEL
NAVAL UNDERSEA WARFARE CENTER DIVISION
1176 HOWELL STREET
NEWPORT RI 02841-1708

IN REPLY REFER TO:

Attorney Docket No. 83063
Date: 12 May 2003

The below identified patent application is available for licensing. Requests for information should be addressed to:

PATENT COUNSEL
NAVAL UNDERSEA WARFARE CENTER
1176 HOWELL ST.
CODE 00OC, BLDG. 112T
NEWPORT, RI 02841

Serial Number 09/949,409
Filing Date 8/31/01
Inventor Paul M. Baggenstoss

If you have any questions please contact James M. Kasischke, Acting Deputy Counsel, at 401-832-4736.

1 Attorney Docket No. 83063

2

3 TIME-SERIES SEGMENTATION

4

5 STATEMENT OF GOVERNMENT INTEREST

6 The invention described herein may be manufactured and used
7 by or for the Government of the United States of America for
8 governmental purposes without the payment of any royalties
9 thereon or therefor.

10

11 CROSS REFERENCE TO OTHER PATENT APPLICATIONS

12 Not applicable.

13

14 BACKGROUND OF THE INVENTION

15 (1) Field of the Invention

16 This invention generally relates to a method and system for
17 identifying data segments within a signal by using naturally
18 occurring boundaries in the signal and updating sample-by-
19 sample.

20 More particularly, the invention is directed to solving the
21 problem of dividing an input signal, such as an acoustic data
22 signal or a speech signal, consisting of multiple "events" into
23 frames where the signal within each frame is statistically
24 "consistent". Once the data has been segmented, detection and

1 classification of events is greatly facilitated. In speech
2 signals, for example, the data becomes segmented into
3 phonetically constant frames or frames in which there are an
4 integer number of pitch periods. This makes determination of
5 pitch more accurate and reliable.

6 (2) Description of the Prior Art

7 Prior to this invention, it has not been known how to
8 divide a time-series (signal) into segments with a fine enough
9 resolution corresponding to individual pitch interval
10 boundaries. The current art for optimally segmenting a time-
11 series consists of first segmenting the data into fixed-size
12 segments, then performing a second stage of segmentation to
13 group together numbers of the fixed-size segments into larger
14 blocks. This approach has a resolution no finer than the size
15 of the fixed-size segments.

16 Because speech signals contain features that are very short
17 in duration, it would be preferable to segment the data to a
18 finer resolution, such as to a resolution of one sample. The
19 current art cannot be used to segment the data to a resolution
20 of one sample because it requires first segmenting to fixed-size
21 segments large enough to extract meaningful features.
22 Furthermore, the existing dynamic-programming solution is
23 computationally impractical because the data has to be processed
24 at each delay and at each segment length.

1 Thus, a problem exists in the art whereby it is necessary
2 to develop a computationally efficient and practical method of
3 segmenting multiple events into frames to a resolution of one
4 sample necessary to identify individual pitch intervals.

5 By way of example of the state of the art, reference is
6 made to the following papers, which are incorporated herein by
7 reference. References pertaining to the prior art are contained
8 in the following references:

9

10 [1] Euler, S.A.; Juang, B.H.; Lee, C.H.; Soong, F.K.,
11 *Statistical Segmentation and Word Modeling Techniques in*
12 *Isolated Word Recognition*, 1990 International Conference on
13 Acoustics, Speech, and Signal Processing, vol.2, pp. 745 -748.

14

15 [2] Svendsen, F. Soong, *On the Automatic Segmentation of Speech*
16 *Signals*, 1987 International Conference on Acoustics, Speech, and
17 Signal Processing, pp. 77-80, Dallas, 1986.

18

19 [3] R. Bellman, S. Dreyfus, *Applied Dynamic Programming*,
20 Princeton Univ. Press, 1962

21

22 [4] R. Kenefic, *An Algorithm to Partition DFT Data into Sections*
23 *of Constant Variance*, IEEE Trans AES, July 1998

24

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24

Referring further to the current state of the art as developed in the field to date, it should be understood that detection and classification of short signals is a high priority for the Navy. Segmentation of a time series is a method that facilitates detection and classification.

In segmentation of short signals, the following is an illustration of the current state of the art. Let there be N samples $x=[x_1..x_N]$. One would like to divide these samples into a number of segments, for example:

$$x = [x_1..x_a] [x_{a+1}..x_b] [x_{b+1}..x_c] [x_{c+1}..x_N],$$

such that the total score, Q , where:

$$Q = Q(x_1..x_a) + Q(x_{a+1}..x_b) + Q(x_{b+1}..x_c) + Q(x_{c+1}..x_N)$$

is as high as possible.

To do this, the score function, $Q(n,t)$, must be known for a segment of length n ending at time t . Assuming it is known, the problem is to find the best number of segments and their start times $\{a, b, c, d, \dots\}$. The standard dynamic-programming approach disclosed in Bellman and also Soong, above, is to first compute the score for all possible segment lengths at all possible end-times. In other words, compute $Q(n,t)$ for $t = n_{min}..N$ and $n = n_{min}$ to n_{max} where n_{min} and n_{max} are the range of allowed segment lengths. The problem is solved by starting at sample n_{min} because the best solution for segmenting the data up

1 to sample n_{\min} is immediately known, it is just the value of the
2 score function Q when $n=n_{\min}$ and $t=n_{\min}$, $Q(n_{\min}, n_{\min})$. Let this be
3 called $Q_b(n_m)$. The best solutions for later samples are then
4 easily found as follows:

5
$$Q_b(t)=Q(n,t)+Q_b(t-n)$$
 maximized over n .

6 Since $Q_b(t-n)$ was already computed, all of the necessary
7 information is available. The value of n for this solution is
8 also saved and is called $n_m(t)$. This process proceeds until
9 $Q_b(N)$, $n_b(N)$. The problem is then solved. The maximum total
10 score is $Q_b(N)$ and the length of the last segment is $n_b(N)$. The
11 other segment lengths are found by working backwards. For
12 example, the length of the next-to-last segment is $n_b(N-n_b(n))$,
13 which was previously stored. This is the standard approach
14 taught in the prior art.

15 In many problems, it is needed to have the best segments
16 and also to pick the best models for each segment. In speech,
17 for example, it may be necessary to know if a segment is voiced
18 or unvoiced speech or it might be necessary to choose the best
19 model order. Let p be an index that ranges over all possible
20 models. To find the best combination of segment lengths and
21 model indexes, first the score function $Q(p,n,t)$ must be known.
22 A slight modification is then made to the above procedure by
23 carrying out the maximizations at each time over both n and p
24 jointly.

1 What has been described so far is the standard approach
2 taught by the Bellman and Soong references. The problem with
3 applying the method to speech processing and other fields is
4 that computing the score function is time-consuming and the
5 method is not practical to apply sample-by-sample as data is
6 acquired. Instead, it is necessary to apply the method to a
7 coarse resolution defined by the frame-processing interval
8 taught by Soong. Features of the data finer than the frame
9 processing interval are filtered out of the data.

10 As mentioned, sample-by-sample processing is normally
11 impractical. If the score function is computed on samples $[x_t-$
12 $n+1..x_t]$, and it is desired to move over one sample to $[x_t-$
13 $n+2..x_{t+1}]$, it is necessary to re-compute the entire score
14 function. This is because the state of the art in signal
15 processing in speech and other fields uses the Fast Fourier
16 Transform (FFT) and a "window" function such as a Hanning
17 window. Window functions are necessary to smooth transitions in
18 the data and eliminate edge effects. This is because the data
19 is processed in "chunks" which are not always aligned with the
20 naturally occurring event boundaries.

21 It should be understood that the present invention would in
22 fact enhance the functionality of the above cited art by the
23 combined effect of eliminating the window function previously
24 used, and providing sample-by-sample updates.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24

SUMMARY OF THE INVENTION

Therefore it is an object of this invention to provide an improved method of time series segmentation.

Another object of this invention is to provide a method for dividing a signal into plural segments of data.

Still another object of this invention is to provide a method for dividing a signal into plural segments of data in the absence of a window function.

Yet another object of the invention is to provide a method for dividing a signal into plural segments of data and updating segment scores thereof one sample at a time.

In accordance with one aspect of this invention, there is provided a method for segmenting a signal into segments having similar spectral characteristics is provided. Initially the method generates a table of previous values from older signal values that contains a scoring value for the best segmentation of previous values and a segment length of the last previously identified segment. The method then receives a new sample of the signal and computes a new spectral characteristic function for the signal based on the received sample. A new scoring function is computed from the spectral characteristic function. Segments of the signal are recursively identified based on the newly computed scoring function and the table of previous

1 values. The spectral characteristic function can be a selected
2 one of an autocorrelation function and a discrete Fourier
3 transform. An example is provided for segmenting a speech
4 signal.

5
6 BRIEF DESCRIPTION OF THE DRAWINGS

7 The appended claims particularly point out and distinctly
8 claim the subject matter of this invention. The various
9 objects, advantages and novel features of this invention will be
10 more fully apparent from a reading of the following detailed
11 description in conjunction with the accompanying drawings in
12 which like reference numerals refer to like parts, and in which:

13 The FIG. is an example of segmentation of speech
14 illustrating the result of the method and system of the present
15 invention.

16
17 DESCRIPTION OF THE PREFERRED EMBODIMENT

18 In general, the present invention is directed to solving
19 the problem of dividing an input signal, such as acoustic data
20 or a speech signal, consisting of multiple "events" into frames
21 where the signal within each frame is statistically
22 "consistent". Once the data has been segmented, detection and
23 classification of events is greatly facilitated. In speech
24 signals, for example, the data becomes segmented into

1 phonetically constant frames or frames in which there are an
2 integer number of pitch periods. This makes determination of
3 pitch more accurate and reliable.

4 This invention was disclosed by the inventor in the
5 following presentation, which is incorporated by reference
6 herein.

7 P. M. Baggenstoss et al., *A Theoretically Optimal*
8 *Probabilistic Classifier Using Class-Specific Features*,
9 2000 International Conference on Pattern Recognition,
10 Barcelona, Spain, September 2, 2000.

11 The invention automatically divides an arbitrary time-
12 series signal into arbitrary-length frames or segments wherein
13 the data in each frame is "consistent". This ability to
14 determine a consistent frame of data facilitates detection and
15 classification of each frame of the data as well as the data as
16 a whole. Current detectors locate events only to an FFT frame.
17 The proposed method can locate events to a resolution of one
18 sample. The results of experiments show that the segmentation
19 occurring in the present invention is as good as possible by a
20 human operator.

21 As indicated above, the problem is to divide a time-series
22 signal such as a digitized audio stream into segments
23 corresponding to the naturally occurring events in the signal.
24 The invention provides a non-windowed processing method (in

1 contrast to the state of the art which uses windowing) which
2 allows recursive update of a spectral feature function such as
3 one of a Discrete Fourier Transform (DFT) and a circular Auto
4 Correlation Function (ACF). This method has the added benefit
5 of causing the resulting segments to be perfectly aligned to
6 event boundaries.

7 The method of the present invention allows sample by sample
8 updating of spectral feature function which does not require a
9 window function. The window function is not necessary because
10 the segments will be exactly aligned to the "event" boundaries
11 in the signal. Also, because no window function is used, it is
12 possible to update the score function efficiently by accounting
13 only for the added and dropped samples.

14 When a spectral feature function such as a discrete Fourier
15 transform (DFT) is computed on samples $[x_{t-n+1}..x_t]$, denoted $X_t[k]$
16 where t is the sample and k is the transform variable, and it is
17 desired to compute it on samples $[x_{t-n+2}..x_{t+1}]$, denoted $X_{t+1}[k]$,
18 $X_t[k]$ is related to $X_{t+1}[k]$ by the following equation:

$$19 \quad X_{t+1}[k] = e^{j2k/n} [X_{t+1}[k] - (x_{t-n+1} - x_{t+1})] \quad (1)$$

20 If the spectral feature function is a circular
21 autocorrelation function (ACF) computed on samples $[x_{t-n+1}..x_t]$,
22 denoted $r_t[\tau]$ where τ is the correlation variable, and it is
23 desired to compute it on samples $[x_{t-n+2}..x_{t+1}]$, denoted $r_{t+1}[\tau]$
24 then:

1
$$r_{t+1}[\tau] = r_t[\tau] + (x_{t+1} - x_{t-n+1})(x_{t-n+1} - x_{t+1-\tau})/n \quad (2)$$

2 Score functions that are computed from the spectral feature
3 function can be computed efficiently at each sample. Other
4 types of efficiently-computed score functions are also possible.
5 Previous values of the score function for the best segmentation
6 and the length of the last segment can be stored in a table.
7 Upon segmentation, the score and length of the latest value can
8 be utilized with the table values for efficiently obtaining the
9 current best segmentation. Accordingly, at any time the
10 invention using dynamic programming can segment a stream of time
11 series data into segments having like characteristics. These
12 segments can then be classified.

13 Applying the current invention for speech processing the
14 Autocorrelation function (ACF) is used as the spectral feature
15 function. Because the ACF is sensitive to spectral features in
16 the data, the resulting segments are on boundaries where the
17 spectrum changes.

18 In speech data, a "reward" is additionally added to the
19 score function for segments matching the pitch interval exactly.
20 To determine the score of a segment, the ACF is computed, and
21 then the Levinson recursion is used to compute the linear
22 prediction error variance for every model order up to a maximum
23 (of about 16). The score for a given model order p on a segment
24 of length n is:

1 $Q(p,n) = (-n/2) (\log(\sigma^2[p,n])+1) - (p/2)*\log(n) + K$ (3)

2 where n is the segment length and $\sigma^2[p,n]$ is the prediction
3 error variance for model order p , and K is a "reward" value for
4 periodicity. The well-known Levinson-Durbin algorithm can be
5 used to compute $\sigma^2[p,n]$ from the ACF efficiently.

6 The term $(p/2)*\log(n)$ is the well-known Minimum Description
7 Length (MDL) penalty score. To "reward" the segment for
8 matching the pitch interval, a positive number K is added to
9 $Q(p,n)$ if the ACF of the segment shows "periodicity". To
10 determine periodicity in the speech application, every division
11 factor $d = 2$ up to $d = 6$ is tested. The meaning of d is the
12 number of pitch intervals in the segment. For each value of d ,
13 the smallest ACF lag in the set $\{r[0], r[n/d], r[2n/d], \dots$
14 $r[n/2]\}$ is determined. d_{max} is determined as the division factor
15 producing the largest minimum ACF value. If d_{max} is greater than
16 a fraction of $r[0]$, it can be labeled as periodic with a
17 division factor d_{max} and thus the period is n/d_{max} . The fraction
18 is established by trial and error based on the given
19 application. For speech recognition .5 has been found to be an
20 effective fraction. While this only happens rarely, it is bound
21 to happen for some segment (and all segments are tested), thus
22 the method works. The reward value used is a monotonically
23 increasing function of d_{max} .

1 The FIG. is an illustration of a segmented speech signal
2 10. Identified segments 12 are indicated by dashed lines. The
3 Autocorrelation Function is provided for three identified
4 segments (A), (B) and (C). (D) is provided as the
5 autocorrelation function of an arbitrary region of the speech
6 signal that is not indicated by the scoring function as a
7 segment. The segments (A), (B) and (C) enclose well-defined
8 events or periodic (voiced) areas of exactly 3, 2, and 5 pitch
9 intervals, respectively. Further, the illustrated non-windowed
10 ACF functions for these segments begin and end at the same
11 levels showing almost perfect periodicity. The non-windowed ACF
12 for the arbitrary region (D) which is slightly smaller than
13 segment (C) does not have this property. From this example, it
14 can be seen how the segmentation works hand-in hand with the
15 non-windowed ACF.

16 The key feature of this invention is the use of non-
17 windowed processing which permits fast computation of a spectral
18 feature function such as a DFT or ACF on a sample-by sample
19 basis. Thus, variations of the method include any method that
20 uses the DFT, ACF or other recursively computed spectral feature
21 function, as described herein.

22 In view of the above detailed description, it is
23 anticipated that the invention herein will have far reaching
24 applications other than those specifically described.

1 This invention has been disclosed in terms of certain
2 embodiments. It will be apparent that many modifications can be
3 made to the disclosed apparatus without departing from the
4 invention. Therefore, it is the intent of the appended claims
5 to cover all such variations and modifications as come within
6 the true spirit and scope of this invention.

1 Attorney Docket No. 83063

2

3

TIME-SERIES SEGMENTATION

4

5

ABSTRACT OF THE DISCLOSURE

6 A method for segmenting a signal into segments having
7 similar spectral characteristics is provided. Initially the
8 method generates a table of previous values from older signal
9 values that contains a scoring value for the best segmentation
10 of previous values and a segment length of the last previously
11 identified segment. The method then receives a new sample of
12 the signal and computes a new spectral characteristic function
13 for the signal based on the received sample. A new scoring
14 function is computed from the spectral characteristic function.
15 Segments of the signal are recursively identified based on the
16 newly computed scoring function and the table of previous
17 values. The spectral characteristic function can be a selected
18 one of an autocorrelation function and a discrete Fourier
19 transform. An example is provided for segmenting a speech
20 signal.

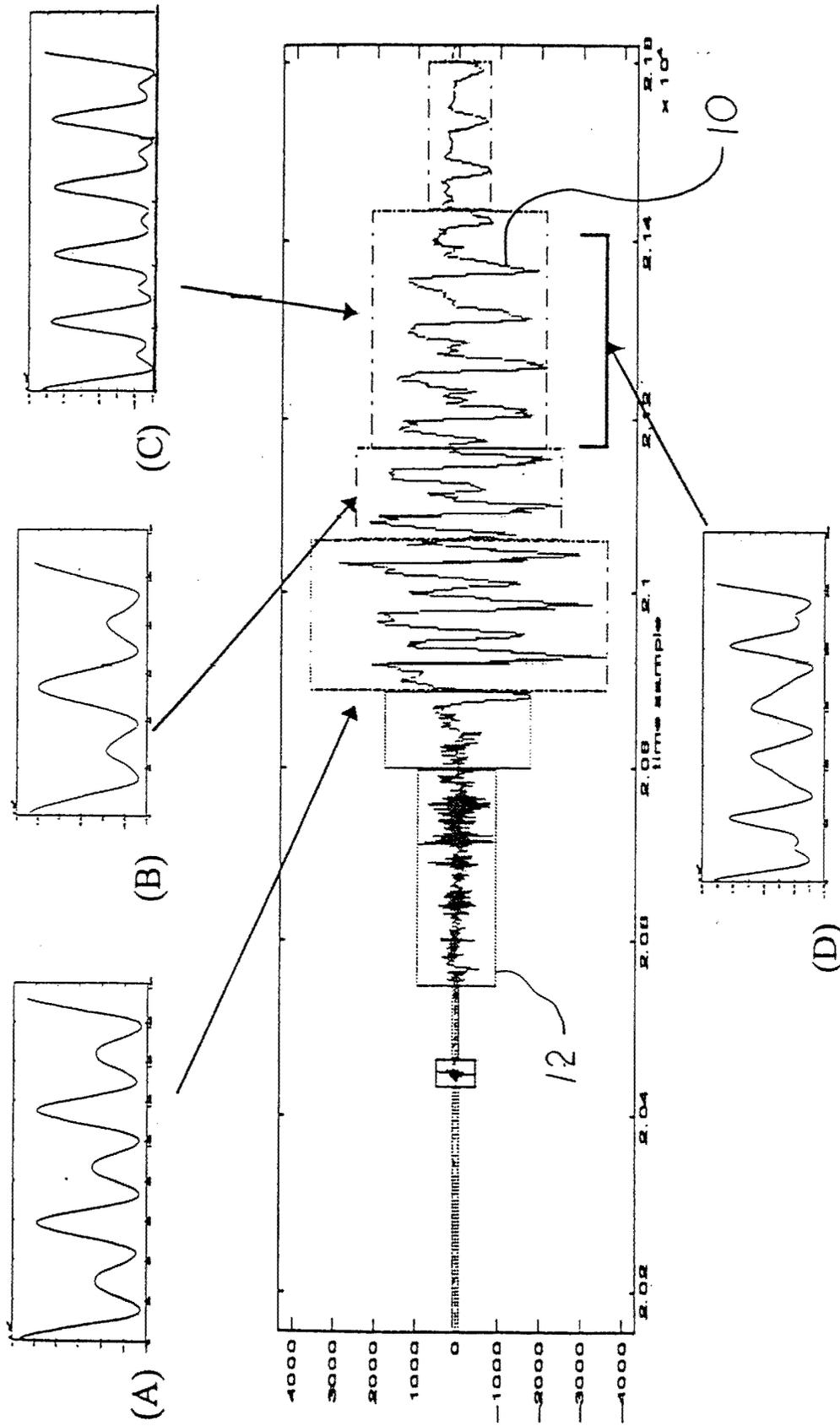


FIG.