



DEPARTMENT OF THE NAVY

NAVAL UNDERSEA WARFARE CENTER DIVISION
1176 HOWELL STREET
NEWPORT RI 02841-1708

IN REPLY REFER TO:

Attorney Docket No. 84264

Date: 14 October 2004

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 10/863,839
Filing Date 1 June 2004
Inventor Francis J. O'Brien

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

DISTRIBUTION STATEMENT A
Approved for Public Release
Distribution Unlimited

20041021 140

Attorney Docket No. 84264
Customer no. 23523

METHOD FOR SPARSE DATA TWO-STAGE
STOCHASTIC MENSURATION

TO WHOM IT MAY CONCERN:

BE IT KNOWN THAT FRANCIS J. O'BRIEN, JR, employee of the United States Government, citizen of the United States of America, 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:

MICHAEL F. OGLO, ESQ.
Reg. No. 20464
Naval Undersea Warfare Center
Division, Newport
Newport, RI 02841-1708
TEL: 401-832-4736
FAX: 401-832-1231

1 Attorney Docket No. 84264

2

3

METHOD FOR SPARSE DATA TWO-STAGE

4

STOCHASTIC MENSURATION

5

6

STATEMENT OF GOVERNMENT INTEREST

7

8

9

10

11

The invention described herein may be manufactured and used by or for the Government of the United States of America for Governmental purposes without the payment of any royalties thereon or therefore.

12

CROSS REFERENCE TO RELATED PATENT APPLICATIONS

13

14

15

16

17

18

19

20

21

22

23

24

The present application is related to the following copending applications: application of F. J. O'Brien, Jr. entitled "Detection of Randomness in Sparse Data Set of Three Dimensional Time Series Distributions," serial number 10/679,866, filed 6 October 2003 (Attorney Docket No. 83996); application of F. J. O'Brien, Jr. entitled "Enhanced System for Detection of Randomness in Sparse Time Series Distributions," filed 3 March 2004 (Attorney Docket No. 83995); application of F. J. O'Brien, Jr. and Chung T. Nguyen entitled "Method for Classifying a Random Process for Data Sets in Arbitrary Dimensions," filed on even date with the present application (Attorney Docket No. 78586); application of F. J. O'Brien, Jr. entitled "Method for Detecting

1 a Spatial Random Process Using Planar Convex Polygon Envelope,"
2 filed on even date with the present application (Attorney Docket
3 No. 83047); and application of F. J. O'Brien, Jr. entitled
4 "Multi-Stage Planar Stochastic Mensuration," filed on even date
5 with the present invention (Attorney Docket No. 83992).

6

7

BACKGROUND OF THE INVENTION

8

(1) Field of the Invention

9

10 The invention generally relates to signal processing/data
11 processing systems for processing time series distributions
12 containing a small number of data points (e.g., less than about
13 ten (10) to twenty-five (25) data points). More particularly,
14 the invention relates to a method and apparatus for classifying
15 the white noise degree (randomness) of a selected signal
16 structure comprising a time series distribution composed of a
17 highly sparse data set. As used herein, the term "random" (or
18 "randomness") is defined in terms of a "random process" as
19 measured by a selected probability distribution model. Thus,
20 pure randomness, pragmatically speaking, is herein considered to
21 be a time series distribution for which no function, mapping or
22 relation can be constituted that provides meaningful insight into
23 the underlying structure of the distribution, but which at the
same time is not chaos.

1 (2) Description of the Prior Art

2 Recent research has revealed a critical need for highly
3 sparse data set time distribution analysis methods and apparatus
4 separate and apart from those adapted for treating large sample
5 distributions. This is particularly the case in applications
6 such as naval sonar systems which require that input time series
7 signal distributions be classified according to their structure,
8 i.e., periodic, transient, random or chaotic. It is well known
9 that large sample methods often fail when applied to small sample
10 distributions, but that the same is not necessarily true for
11 small sample methods applied to large data sets. As a typical
12 example, it is extremely valuable for a sonar operator to have
13 tools with significant accuracy in distinguishing between random
14 noise and the signal produced by a contact of interest.

15 Very small data set distributions may be defined as those
16 with less than about ten (10) to twenty-five (25) measurement
17 (data) points. Such data sets can be analyzed mathematically
18 with certain nonparametric discrete probability distributions, as
19 opposed to large-sample methods which normally employ continuous
20 probability distributions (such as the Gaussian).

21 The probability theory discussed herein and utilized by the
22 present invention is well known. It may be found, for example,
23 in works such as P.J. Hoel et al., Introduction to the Theory of

1 Probability, Houghton-Mifflin, Boston, MA, 1971, which is hereby
2 incorporated herein by reference.

3 Also, as will appear more fully below, it has been found to
4 be important to treat white noise signals themselves as the time
5 series signal distribution to be analyzed, and to identify the
6 characteristics of that distribution separately. This aids in
7 the detection and appropriate processing of received signals in
8 numerous data acquisition contexts, not the least of which
9 include naval sonar applications. Accordingly, it will be
10 understood that prior analysis methods and apparatus analyze
11 received time series data distributions from the point of view of
12 attempting to find patterns or some other type of correlated data
13 therein. Once such a pattern or correlation is located, the
14 remainder of the distribution is simply discarded as being noise.
15 It is believed that the present invention will be useful in
16 enhancing the sensitivity of present analysis methods, as well as
17 being useful on its own.

18 Various aspects related to the present invention are
19 discussed in the following exemplary patents:

20 U.S. Patent No. 6,068,659, issued May 30, 2000, to Francis
21 J. O'Brien, Jr., discloses a method for measuring and recording
22 the relative degree of pical density, congestion, or crowding of
23 objects dispersed in a three-dimensional space. A Population
24 Density Index is obtained for the actual conditions of the

1 objects within the space as determined from measurements taken of
2 the objects. The Population Density Index is compared with
3 values considered as minimum and maximum bounds, respectively,
4 for the Population Density Index values. The objects within the
5 space are then repositioned to optimize the Population Density
6 Index, thus optimizing the layout of objects within the space.

7 U.S. Patent No. 5,506,817, issued April 9, 1996, to Francis
8 J. O'Brien, Jr., discloses an adaptive statistical filter system
9 for receiving a data stream comprising a series of data values
10 from a sensor associated with successive points in time. Each
11 data value includes a data component representative of the motion
12 of a target and a noise component, with the noise components of
13 data values associated with proximate points in time being
14 correlated. The adaptive statistical filter system includes a
15 prewhitener, a plurality of statistical filters of different
16 orders, stochastic decorrelator and a selector. The prewhitener
17 generates a corrected data stream comprising corrected data
18 values, each including a data component and a time-correlated
19 noise component. The plural statistical filters receive the
20 corrected data stream and generate coefficient values to fit the
21 corrected data stream to a polynomial of corresponding order and
22 fit values representative of the degree of fit of corrected data
23 stream to the polynomial. The stochastic decorrelator uses a
24 spatial Poisson process statistical significance test to

1 determine whether the fit values are correlated. If the test
2 indicates the fit values are not randomly distributed, it
3 generates decorrelated fit values using an autoregressive moving
4 average methodology which assesses the noise components of the
5 statistical filter. The selector receives the decorrelated fit
6 values and coefficient values from the plural statistical filters
7 and selects coefficient values from one of the filters in
8 response to the decorrelated fit values. The coefficient values
9 are coupled to a target motion analysis module which determines
10 position and velocity of a target.

11 U.S. Patent No. 6,466,516 B1, issued October, 15, 2002, to
12 O'Brien, Jr. et al., discloses a method and apparatus for
13 automatically characterizing the spatial arrangement among the
14 data points of a three-dimensional time series distribution in a
15 data processing system wherein the classification of said time
16 series distribution is required. The method and apparatus
17 utilize grids in Cartesian coordinates to determine (1) the
18 number of cubes in the grids containing at least one input data
19 point of the time series distribution; (2) the expected number of
20 cubes which would contain at least one data point in a random
21 distribution in said grids; and (3) an upper and lower
22 probability of false alarm above and below said expected value
23 utilizing a discrete binomial probability relationship in order
24 to analyze the randomness characteristic of the input time series

1 distribution. A labeling device also is provided to label the
2 time series distribution as either random or nonrandom, and/or
3 random or nonrandom within what probability, prior to its output
4 from the invention to the remainder of the data processing system
5 for further analysis.

6 U.S. Patent No. 6,397,234 B1, issued May 28, 2002, to
7 O'Brien, Jr. et. al., discloses a method and apparatus for
8 automatically characterizing the spatial arrangement among the
9 data points of a time series distribution in a data processing
10 system wherein the classification of said time series
11 distribution is required. The method and apparatus utilize a
12 grid in Cartesian coordinates to determine (1) the number of
13 cells in the grid containing at least-one input data point of the
14 time series distribution; (2) the expected number of cells which
15 would contain at least one data point in a random distribution in
16 said grid; and (3) an upper and lower probability of false alarm
17 above and below said expected value utilizing a discrete binomial
18 probability relationship in order to analyze the randomness
19 characteristic of the input time series distribution. A labeling
20 device also is provided to label the time series distribution as
21 either random or nonrandom, and/or random or nonrandom.

22 U.S. Patent No. 6,597,634 B1, issued July 22, 2003, to
23 O'Brien, Jr. et al., discloses a signal processing system to
24 processes a digital signal converted from to an analog signal,

1 which includes a noise component and possibly also an information
2 component comprising small samples representing four mutually
3 orthogonal items of measurement information representable as a
4 sample point in a symbolic Cartesian four-dimensional spatial
5 reference system. An information processing sub-system receives
6 said digital signal and processes it to extract the information
7 component. A noise likelihood determination sub-system receives
8 the digital signal and generates a random noise assessment of
9 whether or not the digital signal comprises solely random noise,
10 and if not, generates an assessment of degree-of-randomness. The
11 information processing system is illustrated as combat control
12 equipment for undersea warfare, which utilizes a sonar signal
13 produced by a towed linear transducer array, and whose mode
14 operation employs four mutually orthogonal items of measurement
15 information.

16 The above described references do not show a multi-stage
17 process that may be utilized to select between the most accurate
18 distribution for computing a probability for comparison with a
19 false alarm probability in order to classify sparse data as noise
20 or signal.

21

22

SUMMARY OF THE INVENTION

23

24

It is an object of the present invention to provide an improved method for analyzing sparse data.

1 It is another object of the present invention to classify
2 sparse data as signal or noise.

3 It is yet another object of the invention to provide a
4 method and apparatus including an automated measurement of the
5 spatial arrangement among a very small number of points, object,
6 measurements or the like whereby an ascertainment of the noise
7 degree (i.e., randomness) of the time series distribution may be
8 made.

9 It is yet another object of the invention to provide a
10 method and apparatus useful in naval sonar systems which require
11 acquired signal distributions to be classified according to their
12 structure (i.e., periodic, transient, random, or chaotic) in the
13 processing and use of those acquired signal distributions as
14 indications of how and from where they were originally generated.

15 Further, it is an object of the invention to provide a
16 method and apparatus capable of labeling a time series
17 distribution with (1) an indication as to whether or not it is
18 random in structure, and (2) an indication as to whether or not
19 it is random within a probability of false alarm of a specific
20 randomness calculation.

21 These and other objects, features, and advantages of the
22 present invention will become apparent from the drawings, the
23 descriptions given herein, and the appended claims. However, it
24 will be understood that above listed objects and advantages of

1 the invention are intended only as an aid in understanding
2 certain aspects of the invention, are not intended to limit the
3 invention in any way, and do not form a comprehensive or
4 exclusive list of objects, features, and advantages.

5 Accordingly, the present invention comprises a method for
6 characterizing sparse data as signal or noise which is preferably
7 utilized in conjunction with other techniques discussed in the
8 related patent applications referenced hereinbefore. The method
9 may comprise one or more method steps such as, for example only,
10 creating a virtual window having a two-dimensional area
11 containing a distribution of data points of the sparse data for a
12 selected time period and/or subdividing substantially the
13 entirety of the area of the virtual window into a plurality k of
14 cells wherein each of the plurality k of cells have the same
15 polygonal shape and define the same area value.

16 Additional steps may comprise determining a quantity Θ
17 wherein Θ represents an expected proportion of the plurality k
18 of cells which will be nonempty in a random distribution. When
19 Θ is less than a pre-selected value, then the method may
20 comprise utilizing a Poisson distribution to determine a first
21 mean of the data points. When Θ is greater than the pre-
22 selected value, then the method may comprise utilizing a binomial
23 distribution to determine a second mean of the data points.

1 The method may further comprise computing a probability p
2 from the first mean or the second mean, depending on whether Θ
3 is greater than or less than the pre-selected value. Other steps
4 may comprise determining a false alarm probability α based on a
5 total number of the plurality of k cells.

6 By comparing p with α , the method may be utilized to then
7 determine whether to characterize the sparse data as noise or
8 signal.

9 In one embodiment, the distribution of the sparse data from
10 the selected time period comprises less than about twenty-five
11 (25) data points.

12 In one example, the method the pre-selected amount discussed
13 above is equal to 0.10 such that if $\Theta \leq 0.10$, then the Poisson
14 distribution is utilized, and if $\Theta > 0.10$, then the binomial
15 distribution is utilized.

16 In one presently preferred embodiment, the step of
17 determining a probability of false alarm rate α comprises
18 setting the alarm rate α equal to 0.01 when the total number of
19 the plurality of k cells is greater than 25, and/or determining a
20 probability of false alarm rate α comprises setting the alarm
21 rate α equal to 0.05 when the total number of the plurality of k
22 of cells is greater than or equal to 5 and less than or equal to
23 25 and/or determining a probability of false alarm rate α

1 comprises setting the alarm rate α equal to 0.10 when the total
2 number of the plurality of k cells is less than 5.

3 The above and other novel features and advantages of the
4 invention, including various novel details of construction and
5 combination of parts will now be more particularly described with
6 reference to the accompanying drawings and pointed out by the
7 claims. It will be understood that the particular device and
8 method embodying the invention is shown and described herein by
9 way of illustration only and not as limitations on the invention.
10 The principles and features of the invention may be employed in
11 numerous embodiments without departing from the scope of the
12 invention in its broadest aspects.

13

14

BRIEF DESCRIPTION OF THE DRAWINGS

15 Reference is made to the accompanying drawings in which is
16 shown an illustrative embodiment of the apparatus and method of
17 the invention, from which its novel features and advantages will
18 be apparent to those skilled in the art, and wherein:

19 FIG. 1 is a hypothetical depiction in Cartesian coordinates
20 of a representative white noise (random) time series signal
21 distribution;

22 FIG. 2 is a hypothetical illustrative representation of a
23 virtual window in accordance with the invention divided into a

1 grid of square cells each having a side of length δ , and an area
2 of δ^2 ;

3 FIG. 3 is a block diagram representatively illustrating the
4 method steps of the invention;

5 FIG. 4 is a block diagram representatively illustrating an
6 apparatus in accordance with the invention; and

7 FIG. 5 is a table showing an illustrative set of discrete
8 binomial probabilities for the randomness of each possible number
9 of occupied cells of a particular time series distribution within
10 a specific probability of false alarm rate of the expected
11 randomness number.

12

13

DESCRIPTION OF THE PREFERRED EMBODIMENT

14 Referring now to the drawings, a preferred embodiment of the
15 method and apparatus of the invention will be presented first
16 from a theoretical perspective, and thereafter, in terms of a
17 specific example. In this regard, it is to be understood that
18 all data points are herein assumed to be expressed and operated
19 upon by the various apparatus components in a Cartesian
20 coordinate system. Accordingly, all measurement, signal and
21 other data input existing in terms of other coordinate systems is
22 assumed to have been re-expressed in a Cartesian coordinate
23 system prior to its input into the inventive apparatus or the
24 application of the inventive method thereto.

1 where

2 $K_I = \delta_I^2 k_I / (\Delta t \cdot \Delta Y) \leq 1$ and

3 $K_{II} = \delta_{II}^2 k_{II} / (\Delta t \cdot \Delta Y) \leq 1$

4 In cases with very small amplitudes, it may occur that $\text{int}(\Delta Y / \delta_I)$
5 ≤ 1 or $\text{int}(\Delta Y / \delta_{II}) \leq 1$. In such cases, the solution is to round
6 off either quantity to the next highest value (i.e., ≥ 2). This
7 weakens the theoretical approach, but it allows for practical
8 measurements to be made.

9 Thus, for example, if Δt (or N)=30, and $\Delta Y=20$, then $k=24$
10 and $\delta=5.0$. Accordingly, $k * \delta^2 = 24 * 25 = 600 = \cong t * \Delta Y$. In
11 essence, therefore, the above relation defining the value k
12 selects the number of squares of length δ and area δ^2 which fill
13 up the total space $\cong t * \Delta Y$ to the greatest extent possible (i.e.,
14 ideally $k * \delta^2 = \cong t * \Delta Y$).

15 From the selected partitioning parameter k , the region
16 (area) $\cong t * \Delta Y$ is carved up into k squares with the length of
17 each square being δ as defined above. In other words, the
18 horizontal (or time) axis is marked off into intervals, exactly
19 $\text{int}(\Delta t / \delta)$ of them, so that the time axis has the following
20 arithmetic sequence of cuts (assuming that the time clock starts
21 at $\Delta t = 0$):

22 $0, \delta, 2\delta, \dots, \text{int}(\Delta t / \delta) * \delta$ (3)

1 Likewise, the vertical (or measurement or amplitude) axis is cut
2 up into intervals, exactly $\text{int}(\Delta Y/\delta)$ of them, so that the
3 vertical axis has the following arithmetic sequence of cuts:

$$4 \quad \min(Y), \min(Y) + \delta, \dots, \min(Y) + \text{int}(\Delta Y/\delta) * \delta = \max(Y),$$

5 where \min is the minimum operator and δ is defined as above.

6 Based on the Poisson point process theory for a measurement
7 set of data in a time interval Δt of measurement magnitude ΔY ,
8 that data set is considered to be purely random (or "white
9 noise") if the number of partitions k are nonempty (i.e., contain
10 at least one data point of the time series distribution thereof
11 under analysis) to a specified degree. The expected number of
12 nonempty partitions in a random distribution is given by the
13 relationship:

$$14 \quad k * \Theta = k * (1 - e^{-N/k}) \quad (4)$$

15
16 where the quantity Θ is the expected proportion of nonempty
17 partitions in a random distribution and N/k is "the parameter of
18 the spatial Poisson process" corresponding to the average number
19 of points observed across all subspace partitions.

20 The boundary, above and below $k * \Theta$, attributable to random
21 variation and controlled by a false alarm rate is the so-called
22 "critical region" of the test. The quantity Θ not only
23 represents (a) the expected proportion of nonempty partitions in

1 a random distribution, but also (b) the probability that one or
2 more of the k partitions is occupied by pure chance, as is well
3 known to those in the art. The boundaries of the random process
4 are determined in the following way.

5 Let M be a random variable representing the integer number
6 of occupied cells (partitions) as illustratively shown in FIG. 2.
7 Let m be an integer (sample) representation of M . Let m_1 be the
8 quantity forming the lower random boundary of the statistic $k *
9 \Theta$ given by the binomial criterion:

10

$$11 \quad P(M \leq m) \leq \alpha_0/2, \min(\alpha/2 - \alpha_0/2) \quad (5)$$

12 where

$$13 \quad P(M \leq m) = \sum_{m=0}^{m=m_1} B(m; k, \Theta) \text{ from } m=0 \text{ to } m=m_1, \text{ and}$$

14 k and Θ are defined as above.

15

16 $B(m; k, \Theta)$ is the binomial probability function given

17 as:

18

$$19 \quad B(m; k, \Theta) = (k, m) (\Theta)^m (1-\Theta)^{k-m} \quad (6)$$

20

21 where (k, m) is the binomial coefficient, $(k, m) = k! / m! (k-$

22 $m)!$ and $\sum B(m; k, \Theta)$ from $m=0$ to $m=k$ equals 1.0.

1 The quantity α_0 is the probability of coming closest to an exact
 2 value of the pre-specified false alarm probability α , and m_1 is
 3 the largest value of m such that $P(M \leq m) \leq \alpha_0/2$. It is an
 4 objective of this method to minimize the difference between α
 5 and α_0 . The recommended values of α (the probability false alarm
 6 rate) for differing values of spatial subsets k are as follows:

7
 8 If $k > 25$, the $\alpha = 0.01$;

9 If $5 \leq k \leq 25$, then $\alpha = 0.05$; and

10 If $k < 5$, then $\alpha = 0.10$ (7)

11
 12 The upper boundary of the random process is called m_2 , and
 13 is determined in a manner similar to the determination of m_1 .

14 Thus, let m_2 be the upper random boundary of the statistic
 15 $k^* \Theta$ given by:

16
 17
$$P(M \geq m) \leq \alpha_0/2, \min(\alpha/2 - \alpha_0/2)$$
 (8)

18 where

19
$$P(M \geq m) = \sum_{m=m_2}^k B(m; k, \Theta) \leq \alpha_0/2$$

20 or

21
$$P(M \geq m) = 1 - \sum_{m=0}^{m_2} B(m; k, \Theta) \leq \alpha_0/2$$

1 α_0 is the probability of coming closest to an exact value of
2 the pre-specified false alarm probability α , and m_2 is the
3 largest value of m such that $P(M \geq m) \leq \alpha_0/2$. It is an objective
4 of the invention to minimize the difference between α and α_0 .

5 Hence, the subsystem determines if the signal structure
6 contains m points within the "critical region" warranting a
7 determination of "random".

8 The subsystem also assesses the random process hypothesis by
9 testing:

$$10 \quad H_0: \bar{P} = \Theta \text{ (Noise)} \quad (9)$$

$$11 \quad H_1: \bar{P} \neq \Theta \text{ (Signal + Noise)}$$

12 Where $\bar{P} = m/k$ is the sample proportion of signal points
13 contained in the k subregion partitions expected to be occupied
14 by a truly random (stochastic) spatial distribution. As noted
15 above, FIG. 1 shows what a hypothetical white noise (random)
16 distribution looks like in Cartesian time-space.

17 Thus, if $\Theta \approx \bar{P} = m/k$, the observed distribution conforms to
18 a random distribution corresponding to "white noise".

19 The estimate for the proportion of k cells occupied by N
20 measurements (\bar{P}) is developed in the following manner. Let each
21 of the k cells of length δ be denoted by C_{ij} and the number of
22 objects observed in each C_{ij} cell be denoted card (C_{ij}) where card

1 means "cardinality" or subset count. C_{ij} is labeled from left to
2 right starting at the lower left-hand corner $C_{11}, C_{12}, \dots, C_{46}$
3 (see FIG. 2).

4 Next to continue the example for $k = 24$ shown in FIG. 2,
5 define the following count quantity for the 6×4 partition
6 comprising whole square subsets:

7
8
$$X_{ij} = 1 \text{ if card } (C_{ij}) > 0; i = 1 \text{ to } 4, j = 1 \text{ to } 6$$

9
10
$$X_{ij} = 0 \text{ if card } (C_{ij}) = 0; i = 1 \text{ to } 4, j = 1 \text{ to } 6 \quad (10)$$

11
12 where card is the cardinality or count operator. X_{ij} is a
13 dichotomous variable taking on the individual values of 1 if a
14 cell C_{ij} has one or more objects present, and a value of 0 if the
15 box is empty.

16 Then calculate the proportion of 24 cells occupied in the
17 partition region:

18
19
$$\bar{P} = 1/24 \sum \sum X_{ij} \quad (11)$$

20
21 where the sums are taken from $j = 1$ to 6 and $i = 1$ to 4,
22 respectively.

23 The generalization of this example to any sized table is
24 obvious, and within the scope of the present invention. For the

1 general case, it will be appreciated that, for the statistics X_{ij}
2 and C_{ij} , the index j runs from 1 to $\text{int}(\Delta t/\delta)$ and the index i runs
3 from 1 to $\text{int}(\Delta Y/\delta)$.

4 In addition, another measure useful in the interpretation of
5 outcomes is the R ratio, defined as the ratio of observed to
6 expected occupancy rates:

$$7 \quad R = m/(k * \Theta) = \bar{P}/\Theta \quad (12)$$

9
10 A rigorous statistical procedure has been developed to
11 determine whether the observed R-value is indicative of "noise"
12 or "signal". The procedure renders quantitatively the
13 interpretations of the R-value whereas the prior art has relied
14 primarily on intuitive interpretation or ad hoc methods, which
15 can be erroneous.

16 In this formulation, one of two statistical assessment tests
17 is utilized depending on the value of the parameter Θ .

18 If $\Theta \leq 0.10$, then a Poisson distribution is employed. To
19 apply the Poisson test, the distribution of the N sample points
20 is observed in the partitioned space. It will be appreciated
21 that a data sweep across all cells within the space will detect
22 some of the squares being empty, some containing $k = 1$ points, k

1 = 2 points, $k = 3$ points, and so on. The number of points in
 2 each k category is tabulated in a table such as follows:

3
 4 Frequency Table of Cell Counts

k (number of cells with points)	N_k (number of points in k cells)
0	N_0
1	N_1
2	N_2
3	N_3
\vdots	\vdots
K	N_k

5
 6 From this frequency table, two statistics are of interests
 7 for the Central Limit Theorem approximation:

8 The "total", $Y = \sum_{k=0}^K kN_k$, and (13)

9 the sample mean, $\mu_0 = \frac{\sum_{k=0}^K kN_k}{\sum_{k=0}^K N_k}$.

10 Then, if $\Theta \leq 0.10$, the following binary hypothesis is of
 11 interest:

$$\begin{aligned}
 H_0 : \mu &= \mu_0(\text{NOISE}) \\
 H_1 : \mu &\neq \mu_0(\text{SIGNAL})
 \end{aligned}
 \tag{14}$$

13 The Poisson test statistic, derived from the Central Limit
 14 Theorem, Eq. (3) is as follows:

$$Z_p = \frac{Y - N\mu_0}{\sqrt{N\mu_0}}, \quad (k > 25) \quad (15)$$

where

$$Y = \sum_{k=0}^K kN_k,$$

and N is the sample size. Then

$$\mu_0 = \frac{\sum_{k=0}^K kN_k}{\sum_{k=0}^K N_k}$$

is the sample mean and sample variance. (It is

well known that $\mu = \sigma^2$ in a Poisson distributions).

The operator compares the value of Z_p against a probability of False Alarm α . α is the probability that the null hypothesis (NOISE) is rejected when the alternative (SIGNAL) is the truth.

The probability of the observed value Z_p is calculated as:

$$p = P(|z_p| \leq Z) = 1 - \frac{1}{\sqrt{2\pi}} \int_{-z_p}^{+z_p} \exp(-.5x^2) dx \quad (16)$$

where $|x|$ means "absolute value" as commonly used in mathematics.

The calculation of Eq. 6, as known to those skilled in the art, is performed in a standard finite series expansion.

On the other hand, if $\Theta > .10$, the invention dictates that the following binary hypothesis set prevail:

$$\begin{aligned}
 H_0: \mu &= k\theta(\text{NOISE}) \\
 H_1: \mu &= k\theta(\text{SIGNAL})
 \end{aligned}
 \tag{17}$$

2 The following binomial test statistic is employed to test the
 3 hypothesis:

$$Z_B = \frac{m \pm c - k\theta}{\sqrt{k\theta(1-\theta)}}
 \tag{18}$$

5 where $c = 0.5$ if $X < \mu$ (Yates Continuity correction factor used
 6 for discrete variables)

7 The quantities of Z_B have been defined previously.

8 The probability of the observed value Z_B is calculated as

$$p = P(|Z_B| \leq Z) = 1 - \frac{1}{\sqrt{2\pi}} \int_{-|z_B|}^{+|z_B|} \exp(-.5x^2) dx
 \tag{19}$$

11
 12 in a standard series expansion.

13 For either test statistic, Z_p or Z_B , the following decision
 14 rule is used to compare the false alarm rate α with the observed
 15 probability of the statistic, p :

$$\begin{aligned}
 \text{if } p &\geq \alpha \Rightarrow \text{NOISE} \\
 \text{If } p &< \alpha \Rightarrow \text{SIGNAL}
 \end{aligned}$$

17 Thus, if the calculated probability value $p > \alpha$, then the
 18 three-dimensional spatial distribution is deemed "noise";
 19 otherwise the X-Y-Z data is characterized as "signal" by the
 20 Rtest.

1 where

2 $\delta = \sqrt{(\Delta t * \Delta Y) / k} = 5.0$

3

4 Thus, the 600 square unit space of the virtual window is
5 partitioned into 24 cells of side 5.0 so that the whole space is
6 filled ($k * \delta^2 = 600$). The time-axis arithmetic sequence of cuts
7 are: 0, 5, ..., $\text{int}(\Delta t / \delta) * \delta = 30$. The amplitude axis cuts are:
8 $\text{min}(Y), \text{min}(Y) + \delta, \dots, \text{min}(Y) + \text{int}(\Delta Y / \delta) * \delta = \text{max}(Y)$.

9 Next, the probability false alarm rate is set at step 110
10 according to the value of k as discussed above. More
11 particularly, in this case $\alpha = 0.01$, and the probability of a
12 false alarm within the critical region is $\alpha/2 = 0.005$.

13 The randomness count is then calculated by first computing
14 device 16 at step 112 according to the relation $k * \Theta = k * (1 - e^{-N/k})$
15 which in this example equals 0.713. Therefore, the number of
16 cells expected to be nonempty in this example if the input time
17 series distribution is random is about 17.

18 The binomial distribution discussed above is then calculated
19 by a second computing device 18 according to the relationships
20 discussed above (step 114, FIG. 3). Representative values for
21 this distribution are shown in FIG. 5 for each number of possible
22 occupied cells m.

1 The upper and lower randomness boundaries then are
2 determined, also by second calculating device 18. Specifically,
3 the lower boundary is calculated using m_1 from FIG. 5 (step 116).
4 Then, computing the binomial probabilities results in $P(M \leq 10) =$
5 .0025. Thus, the lower bound is $m_1 = 10$. FIG. 5 also shows the
6 probabilities for $\Theta = .713$, $k = 24$.

7 The upper boundary, on the other hand, is the randomness
8 boundary m_2 from the criterion $P(M \geq m) \leq \alpha_0/2$. Computing the
9 binomial probabilities gives $P(M \geq 23) = .0032$; hence $m_2 = 23$ is
10 taken as the upper bound (step 118). The probabilities necessary
11 for this calculation also are shown in FIG. 5.

12 Therefore, the critical region is defined in this example as
13 $m_1 \leq 10$, and $m_2 \geq 23$ (step 120).

14 The actual number of cells containing one or more data
15 points of the time series distribution determined by
16 analysis/counter device 20 (step 122, FIG. 3) is then used by
17 divider 22 and a second comparator 24 in the determination of the
18 randomness of the distribution (step 124, FIG. 3). Specifically,
19 using $m = 16$ as an example, it will be seen that $\bar{P} = m/k =$
20 0.667 , and that $R = \bar{P}/\Theta = 0.667/0.713 = .93$.

21 Branching to step 124 (FIG. 3) which the sparse data
22 decision logic module performs, the R statistic value of 0.93 is
23 evaluated statistically. A more precise indicator is obtained by

1 applying the significance test in accord with the present
2 invention, as described earlier. For this calculation, we note
3 that $\theta = .713$, which invokes the Binomial probability model to
4 test the hypothesis:

5

$$\begin{aligned} H_0 : \mu &= k\theta(\text{NOISE}) \\ H_1 : \mu &= k\theta(\text{SIGNAL}) \end{aligned} \quad (21)$$

7

8 In this case, $k\theta = 17.11$. Thus, applying the Binomial test
9 gives:

10

$$\begin{aligned} Z_B &= \frac{m \pm c - k\theta}{\sqrt{k\theta(1-\theta)}} \\ &= \frac{16 - .5 - 17.11}{\sqrt{24(.713)(1-.713)}} \approx -33 \end{aligned} \quad (22)$$

13

14 The p value is computed to be:

15

$$p = P(|Z_B| \leq Z) = 1 - \frac{1}{\sqrt{2\pi}} \int_{-33}^{+33} \exp(-.5x^2) dx = .74 \quad (23)$$

17

18 Since $p = .74$ and $\alpha = 0.1$, and since $p \geq \alpha$, we conclude that the R
19 test shows the volumetric data to be random (NOISE only, with 99%

1 certainty) with the value of $R = .93$ computed for this spatial
2 distribution in 32-space.

3 It is also worth noting in this regard that the total
4 probability is $0.0023 + .0032 = .0055$, which is the probability
5 of being wrong in deciding "random". This value is less than the
6 probability of a false alarm. Thus, the actual protection
7 against an incorrect decision is much higher (by about 45%) than
8 the a priori sampling plan specified.

9 Since $m = 16$ falls inside of the critical region, i.e., $m_1 \leq$
10 $16 \leq m_2$, the decision is that the data represent an essentially
11 white noise distribution (step 126). Accordingly, the
12 distribution is labeled at step 128 by the labeling device 26 as
13 a noise distribution, and transferred back to the data processing
14 system 10 for further processing. In the naval sonar situation,
15 a signal distribution labeled as white noise would be discarded
16 by the processing system, but in some situations a further
17 analysis of the white noise nature of the distribution would be
18 possible. Similarly, the invention is contemplated to be useful
19 as an improvement on systems which look for patterns and
20 correlations among data points. For example, overlapping time
21 series distributions might be analyzed in order to determine
22 where a meaningful signal begins and ends.

23 It will be understood that many additional changes in the
24 details, materials, steps and arrangement of parts, which have

1 been herein described and illustrated in order to explain the
2 nature of the invention, may be made by those skilled in the art
3 within the principles and scope of the invention as expressed in
4 the appended claims.

1 Attorney Docket No. 84246

2

3

METHOD FOR SPARSE DATA TWO-STAGE

4

STOCHASTIC MENSURATION

5

6

ABSTRACT OF THE DISCLOSURE

7

8

9

10

11

12

13

14

15

16

17

18

19

A method is provided for characterizing sparse data as either signal or noise. In one embodiment, a two-dimensional area which contains a distribution of data points of the sparse data from a selected time period is divided into equal size cells wherein each cell has an expectation of containing at least one data sample. Based on the expected proportion of a plurality cells which will be nonempty in a random distribution, a Poisson distribution or a Binomial distribution is utilized to determine a data sample mean. A probability of a false signal is determined from the number of cells utilized. A probability is also computed from the sample mean and compared to the probability of a false signal to determine whether to characterize the sparse data as signal or noise.

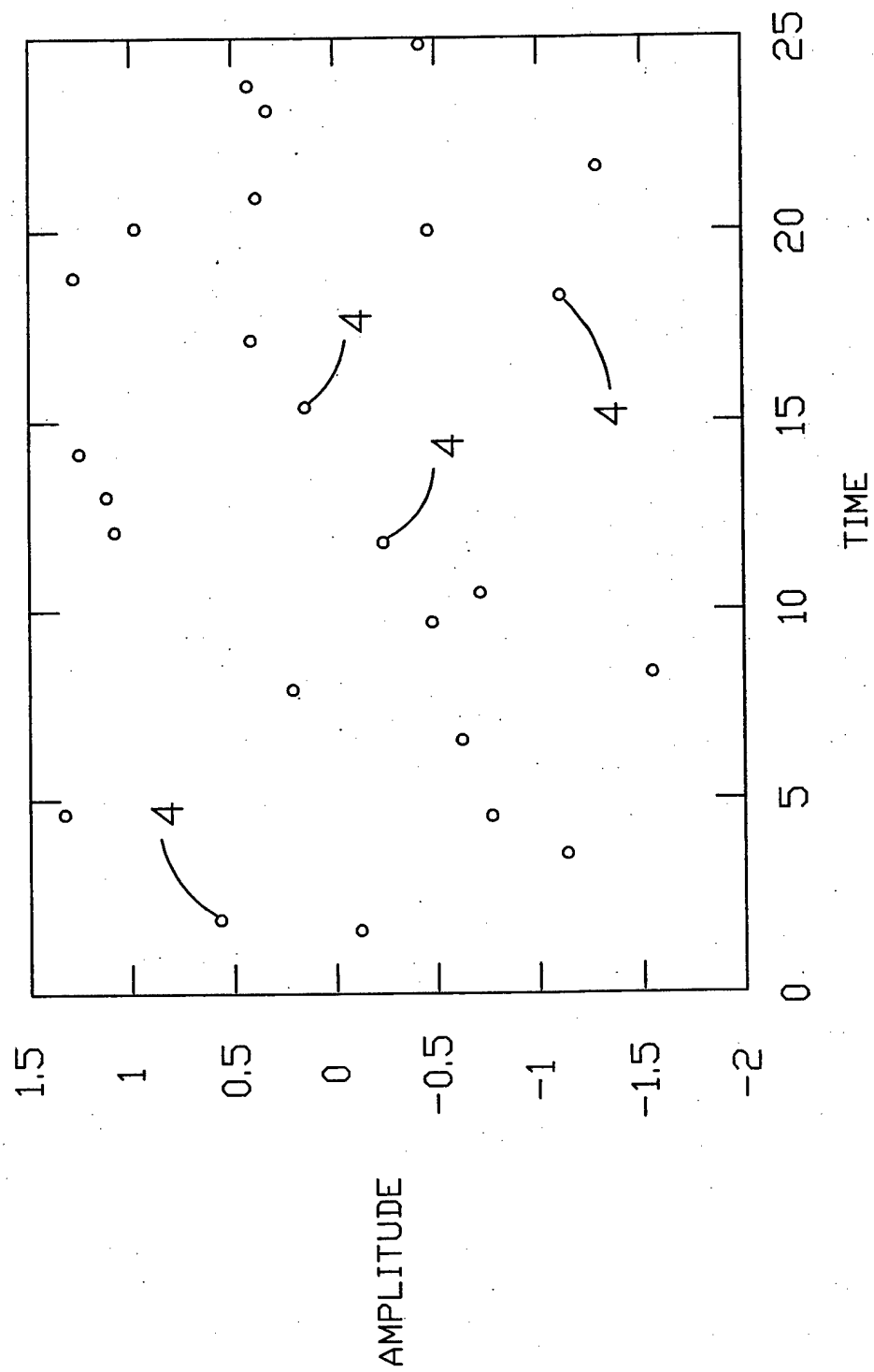


FIG. 1

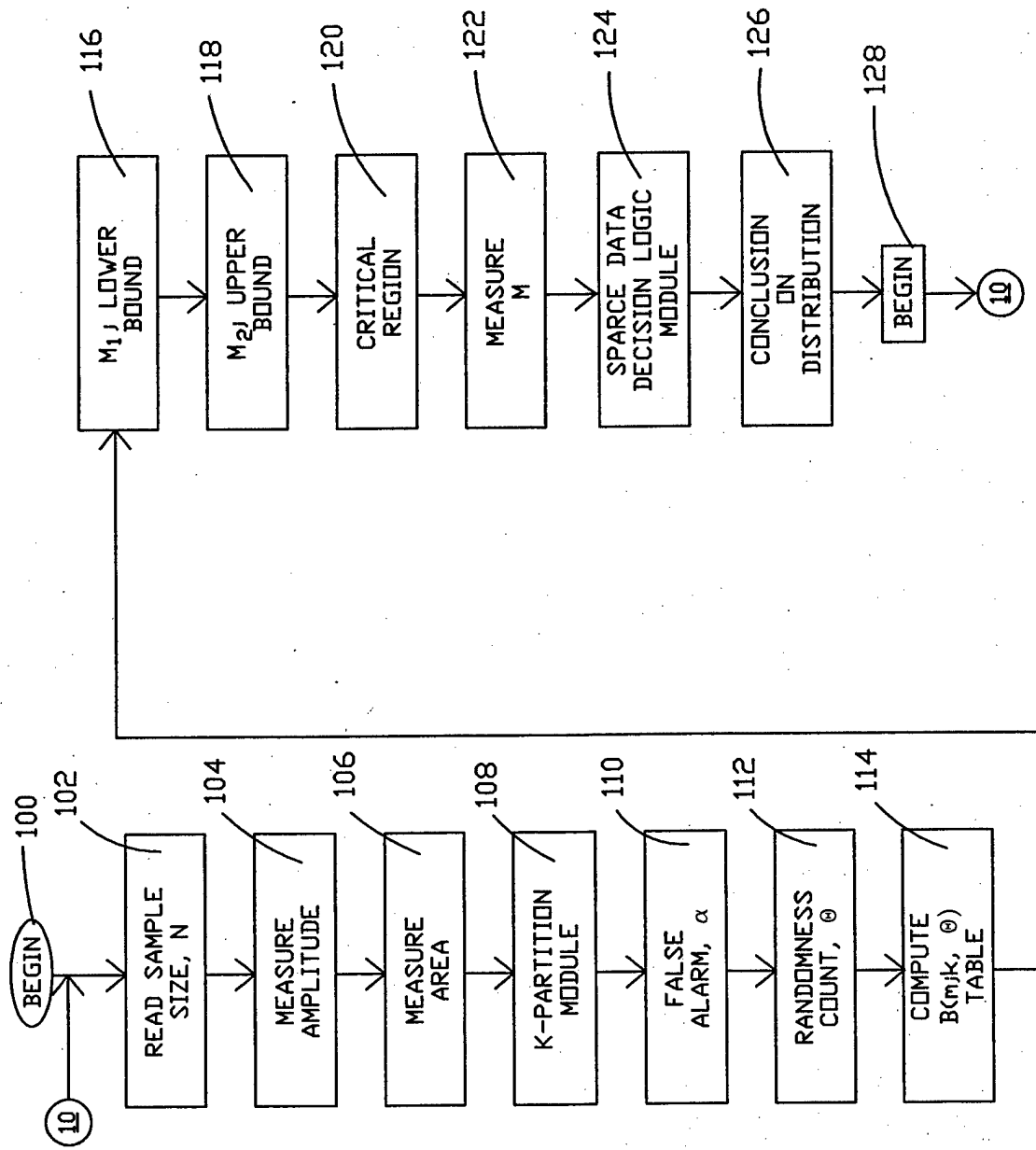


FIG. 3

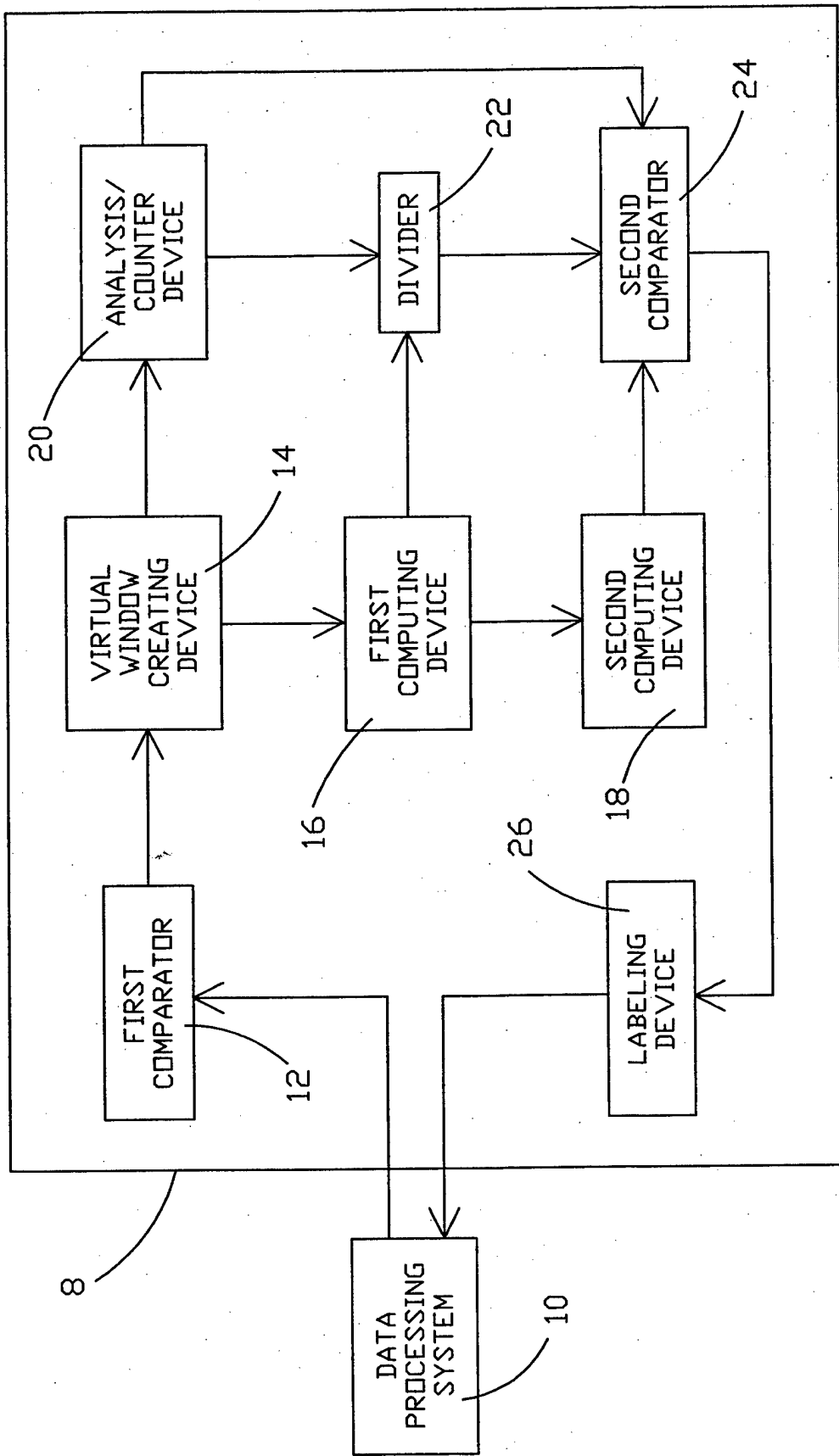


FIG. 4

BINOMIAL TABLE FOR $k=24, \theta=.713, \alpha=.01$

$$P(M=m) = \binom{k}{m} \theta^m (1-\theta)^{k-m} \quad P(M \leq m) = \sum_0^m P(M=m) \quad P(M \geq m)$$

(CUMULATIVE)

0	0	0	0
1	0	0	0
2	0	0	0
3	0	0	0
4	0	0	0
5	0	0	0
6	0	0	0
7	0	0	0
8	.0001	.0001	.0001
9	.0005	.0006	.0006
10	.0017	.0023	.0023
11	.0053	.0076	.0076
12	.0144	.0220	.0220
13	.0334	.0551	.0551
14	DATA NOT SHOWN		
15	FOR $m=14$ to 20		
16	DATA NOT SHOWN		
17	FOR $m=14$ to 20		
18	DATA NOT SHOWN		
19	FOR $m=14$ to 20		
20	DATA NOT SHOWN		
21	.0397	.9833	.0564
22	.0135	.9968	.0167
23	.0029	.9997	.0032
$m=k=24$.0003	1.0000	.0003
			$P(M \leq m) \leq \alpha_0$
			$P(M \geq m) \leq \alpha_0 / 2$

FIG. 5