# New Results From a Mathematical Study of an Adaptive Quantizer

By DEBASIS MITRA

*We consider a general class of multibit adaptive quantizers in which the quantizer function is modified at every sampling instant according to a recursive law with the transitions depending on the value of the quantizer output. We obtain a rather comprehensive set of basic properties of the device which explain the interrelationship of different aspects of the device behavior and their dependence on the parameters of the adaptation algorithm. For the quantitative analysis of the device, we give formulas and bounds for the mean time required for the quantizer function to adapt from an arbitrary initial state to the optimal. A feature new with this work is a unified treatment and a common body of results for quantizers with both bounded and unbounded range. This paper extends all the analytical results reported in an earlier paper, which dealt with a restricted class of quantizers having only four levels.*

*We also present new results from a computational investigation on quantizers up to four bits (sixteen levels). These results indicate, for well-designed examples of the respective classes, the kinds of improvement in performance that can be expected in going from three-bit (eight-level) to four-bit quantizers and from uniform to nonuniform quantizers.*

## I. INTRODUCTION

In a recent paper[1] we obtained a number of fundamental properties of a class of two-bit (four-level) adaptive quantizers useful for coding speech and other continuous signals with a large dynamic range. We also developed formulas for the quantitative analysis of the device. In the present paper, we consider a general, multibit adaptive quantizer and obtain extensions to all the results previously reported. A feature new with this work is a unified treatment and a common body of results for quantizers with both bounded and unbounded range, the former being the case of practical interest.

In the final section of the paper, Section IV, we present results from a computational investigation on adaptive quantizers up to four bits.

Readers familiar with quantizers and whose primary interest is in the performance of the device may skip the earlier sections that contain the development of the mathematical results. Section IV includes a comparison of the performances of uniform and nonuniform quantizers for normally distributed input sequences.

A quantizer with $2N$ levels is shown in Fig. 1. In the figure, *input* refers to the $n$th sample of the continuous signal, $x(n)$, where $n = 0$, $1, \cdots$; *output* refers to the level that is coded before transmission at that time. We let $\xi_1 = 1$ and call $\Delta$ the step size.* In uniform quantizers, $\xi_i = i$ and the vertical axis is also subdivided into equal intervals in the range $(\eta_1 \Delta, \eta_N \Delta)$. In adaptive quantizers which are of interest here, the step size, and hence the entire quantizer function, is time-variable, and the step size at the $n$th sampling instant is denoted by $\Delta(n)$. The parameters $\{\xi_i\}$ and $\{\eta_i\}$ are predetermined and do not change with time.

In this paper, the main algorithm for step-size adaptation is

$$\Delta(n + 1) = M_i \Delta(n) \quad \text{if} \quad \xi_{i-1} \Delta(n) \leqq |x(n)| < \xi_i \Delta(n), \qquad (1)$$

where $M_1, M_2, \cdots, M_N$, called multipliers, are fixed constants. The following natural restrictions are imposed on the multipliers:

$$M_1 < 1 < M_N \quad \text{and} \quad M_1 \leqq M_2 \leqq \cdots \leqq M_N. \qquad (2)$$

Even so, a great deal of the flexibility of the quantizer is incorporated in the multipliers and, to some extent, in the parameters $\{\xi_i\}$ and $\{\eta_i\}$. Observe that the algorithm in (1) utilizes only unit memory and that it is not necessary to transmit to the receiver separate information on the step size.

We shall also be considering the following important variation of (1) in which the step sizes $\{\Delta(n)\}$ are constrained to be within a specific bounded interval $[\bar{K}, \bar{L}]$; suppose $\xi_{i-1} \Delta(n) \leqq |x(n)| < \xi_i \Delta(n)$, then

$$\begin{aligned} \Delta(n + 1) &= M_i \Delta(n) & \text{if} \quad & \bar{K} \leqq M_i \Delta(n) \leqq \bar{L} \\ &= \bar{K} & \text{if} \quad & M_i \Delta(n) \leqq \bar{K} \\ &= \bar{L} & \text{if} \quad & \bar{L} \leqq M_i \Delta(n). \end{aligned} \qquad (3)$$

We call the associated device the *saturating adaptive quantizer*. There are situations where it is attractive to have the interval $[\bar{K}, \bar{L}]$ relatively small.

The most restrictive assumption that is made about the input sequence $\{x(n)\}$ is that it is a sequence of independent random variables (see Sections 1.1 and 1.2 for a discussion). However, in differential PCM schemes in which the quantizer is used together with a

---

* For notational convenience, we also let $\xi_0 = 0$ and $\xi_N = \infty$.
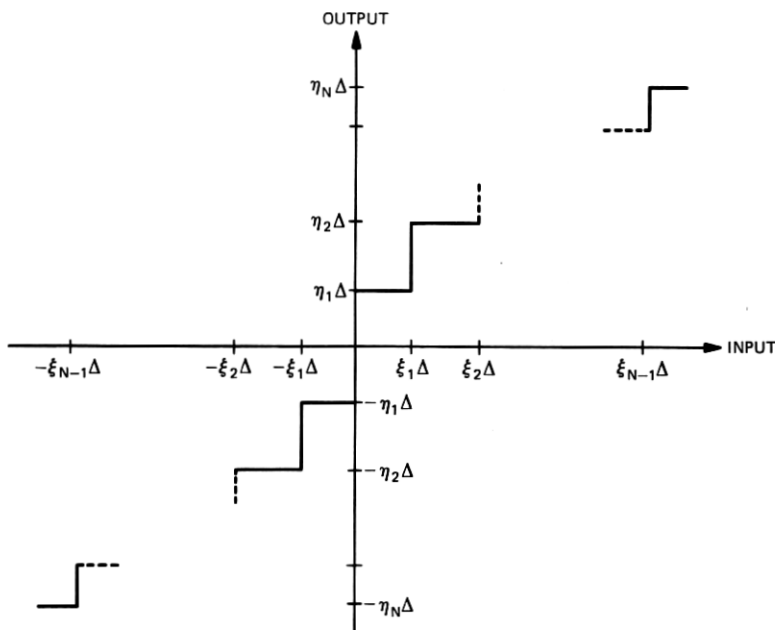
Fig. 1—The quantizer function.

predicting filter in the feedback loop, the effect of the restriction is diminished.

With $\xi_i = i$, the adaptation algorithm in (1) is due to Cummiskey, Flanagan, and Jayant,[2,3] who have also implemented speech coding by a four-bit quantizer. References 1, 2, and 4 may be consulted for a fuller account of the antecedents of the quantizer and related work that has been done in this area. Goodman and Gersho[4] have also examined the general multibit quantizer from a theoretical point of view, and their work complements rather well the work described here.

We briefly summarize here the main features of this paper.

(*i*) The theory that we give here applies to quantizers having bounded range and finite alphabet, with the important properties and relations holding also for quantizers with unbounded range. However, as may be expected, differences do exist between the two types of quantizers. For instance, a key relation in the work of Goodman and Gersho,[4] who do not consider finite range quantizers, called the *design equation*, holds exclusively for the class they consider.

(*ii*) The single most important property of either type of quantizer— ordinary or saturating—that we find is a localization property which states that, for independent identically distributed inputs, there exists a strong localization of the mass of the stationary step-size distribution

about an easily identifiable central value. See Theorem 1, Section 2.2, for a statement of this property. The localization property, together with certain scaling properties of the central state, provides the key to the synthesis of the adaptive quantizers.

(*iii*) A property of the quantizers having important implications is that, under certain conditions, as the range of the multipliers is decreased to approach unity, then the stationary step-size distribution becomes increasingly concentrated about the central step size. A result of this type is given in Ref. 4, where it is shown that a "spread function" has the appropriate behavior. However, the definition of the spread function is novel, and connections, if any, with the dispersion of mass in the distribution are not established. In Section 2.4 we establish the property directly in terms of the mass of the distribution.

(*iv*) In Section III we develop, as design aids, formulas and bounds on the mean adaptation time, i.e., mean time required for the step size to adapt from arbitrary initial values to the central step size.

The mathematical analysis is of a random walk on the integers, in which the state transition probabilities depend on the states. Random walks of the type considered here are encountered in other areas; for instance, in various schemes (up-and-down method, transformed up-and-down method[5-7]) for estimating a quantile of an unknown distribution by using only response, nonresponse data, as is required in bioassay, sensitivity data analysis, and psychological testing. The central properties of the random walk that we obtain here are new and of general interest.

### 1.1 Assumptions and background

Let $\sigma > 0$ denote a scale parameter and let $\mathcal{G}$ denote an equivalence class of distributions $F_\sigma(z)$, $z \geq 0$, in which the distributions are identical to within a scaling operation, i.e.,

$$F_\sigma(\sigma z) = F_1(z). \tag{4}$$

For instance, $\mathcal{G}$ may be the class of half normal distributions, in which case $\sigma^2$ is the variance and $F_1(z) = \Pr[|x| \leq z]$, where $x$ is normal with zero mean and unit variance. In what follows we let $\{x_\sigma(n)\}$ denote a sequence of independent random variables, each with the distribution function $\Pr[|x_\sigma(n)| \leq z] = F_\sigma(z)$.

We recall certain known facts about optimal nonadaptive quantization where $\{x_\sigma(n)\}$ forms the input sequence, $F_\sigma(z)$ is known, and, for some suitable choice of a fidelity criterion such as $E[\{y(n) - x_\sigma(n)\}^2]$ where $\{y(n)\}$ is the output of the quantizer, the optimal step size $\Delta_\sigma$ is computed. With the rms criterion and the inputs normally distrib-

uted, Max[8] has computed $\hat{\Delta}_\sigma$ and, for the nonuniform case, the corresponding optimal parameters $\{\hat{\xi}_i\}$, $\{\hat{\eta}_i\}$ for quantizers with various levels, $N$. A convenient way of presenting such results for any $\mathcal{G}$ is as

$$F_1(\hat{\Delta}_1) = \alpha, \tag{5}$$

where $\alpha$ is some constant, since optimal (nonadaptive) step sizes $\hat{\Delta}_\sigma$ corresponding to the scale parameter $\sigma$ are obtained from

$$\hat{\Delta}_\sigma = \sigma\hat{\Delta}_1. \tag{6}$$

In this paper we show that, when $\sigma$ is fixed and $\{x_\sigma(n)\}$ forms the input to the quantizer, then the step size, a random variable evolving according to either (1) or (3), has a natural center $C_\sigma$. We show, for instance, that the stationary step-size distribution is localized about $C_\sigma$ and that the degree of localization may be arbitrarily increased, although at the cost of other aspects of performance. There are two important facts to note about $C_\sigma$. First, by virtue of its explicit definition, $C_1$ can be made to take almost any desired value by suitable choice of the multipliers. Second, as we show in the following section, the central step size has a scaling property similar to (6). We are therefore in a position to incorporate the results of optimal nonadaptive quantization by identifying $\hat{\Delta}_1$ with $C_1$.

### 1.2 Central state

We consider only quantizers with multipliers having the following form:

$$M_i = \gamma^{m_i} \quad i = 1, 2, \cdots, N, \tag{7}$$

where $\gamma$ is some real number greater than 1 and the $m_i$'s take integral values. With (2), this implies

$$m_1 < 0 < m_N \quad \text{and} \quad m_1 \leqq m_2 \leqq \cdots \leqq m_N. \tag{2'}$$

We shall further take the set of $m_i$'s to be relatively prime, i.e., their greatest common divisor is 1. If, as we shall assume, the initial step size is of the form $\gamma^i$, $i$ integral, then the step size is always of that form and the space of possible step sizes forms a lattice.

Consider an independent identically distributed input sequence $\{x_1(n)\}$, where $\Pr\{|x_1(n)| \leqq z\} = F_1(z)$ and $F_1(\cdot)$ is an element of $\mathcal{G}$. We drop the subscript that identifies the scaling. For $z \geqq 0$, let*

$$B(z) \triangleq \sum_{r=1}^{N} m_r\{F(\xi_r z) - F(\xi_{r-1} z)\}. \tag{8}$$

---

* $F(0) = 0$, $F(z) \to 1$ as $z \to \infty$ and $F(z)$ is monotonic, strictly increasing with $z$.

Since it is also true that

$$B(z) = m_N - \sum_{r=1}^{N-1} (m_{r+1} - m_r)F(\xi_r z), \tag{8'}$$

it is clear that $B(z)$ is a monotonic, strictly decreasing function of $z$; also, $B(0) = m_N > 0$ and $B(z) \to m_1 < 0$ as $z \to \infty$. Hence, there exists a unique integer $i$ with the property that

$$B(\gamma^{i-1}) > 0 \geq B(\gamma^i). \tag{9}$$

We denote $\gamma^i$ by $C$ and refer to it as the *central step size*. All step sizes are considered to be of the form $C\gamma^i$, $i = 0, \pm 1, \pm 2, \cdots$.

*Remarks*:

(*i*) The parameters $\{m_i\}$ and $\gamma$ may be selected to make the resulting central step size $C$ approximate as closely as desired any given real positive number, $\hat{\Delta}$. First, by making $\gamma$ close to unity the grid of possible step sizes can be made sufficiently fine. Second, the integral parameters $\{m_i\}$ can be chosen to make $\sum m_r\{F(\xi_r\hat{\Delta}) - F(\xi_{r-1}\hat{\Delta})\}$ sufficiently small.

(*ii*) So far, we have been concerned with the central step size for the probability distribution $F(z)$, corresponding to the particular scale parameter $\sigma = 1$. To demonstrate the behavior of the central step size with various scale parameters, let $C_\sigma$ denote the central step size corresponding to the input probability distribution, $F_\sigma(z)$, and let $B_\sigma(z)$ be defined like $B(z)$ in (8) with $F(\cdot)$ replaced by $F_\sigma(\cdot)$. Let $\underline{C}_\sigma$ be the unique solution of

$$B_\sigma(\underline{C}_\sigma) = 0, \tag{10}$$

where, of course, $\underline{C}_\sigma$ may not be of the form $\gamma^i$, $i$ integral. However,

$$C_\sigma/\gamma < \underline{C}_\sigma \leq C_\sigma. \tag{11}$$

We observe that $\underline{C}_\sigma$ scales, i.e.,

$$\underline{C}_\sigma = \sigma \underline{C}_1. \tag{12}$$

The above follows from the following property of the functions $\{B_\sigma(\cdot)\}$:

$$B_\sigma(\sigma z) = B_1(z).$$

From (11) and (12),

$$\boxed{C_\sigma/\gamma < \sigma\underline{C}_1 \leq C_\sigma}, \tag{13}$$

and it is in this sense that we say that the central step size scales.

## 1.3 Basic equations

We define a Markov chain and obtain the transition equations for the ordinary quantizer with the inputs being $\{x(n)\}$, which are independent identically distributed, and $\Pr\{|x(n)| \leqq z\} = F(z)$. Let

$$\omega(n) \triangleq \log_\gamma \Delta(n) - \log_\gamma C,$$

so that

$$\omega(n+1) = \omega(n) + m_r \quad \text{if} \quad \xi_{r-1}C\gamma^{\omega(n)} \leqq |x(n)| < \xi_r C\gamma^{\omega(n)}, \quad (14)$$

where $1 \leqq r \leqq N$. We have in (14) a Markov chain on $0, \pm 1, \pm 2, \cdots$, with the central step size $C$ corresponding to the 0 state. Let

$$p(i; n) \triangleq \Pr[\omega(n) = i].$$

The state transition equations are

$$p(i; n+1) = \sum_{r=1}^{N} b^{(r)}(i - m_r)p(i - m_r; n), \quad (15)$$

where the transition probabilities are

$$b^{(r)}(i) \triangleq F(\xi_r C\gamma^i) - F(\xi_{r-1}C\gamma^i), \quad 1 \leqq r \leqq N. \quad (16)$$

The qualitative results that we obtain are based on the following two relations that do not depend on the particular distribution $F(z)$.

$(i)$ $\quad 0 \leqq F(\xi_r\gamma^i) < F(\xi_r\gamma^{i+1}) \leqq 1$

$$\text{for all } i \text{ and } 1 \leqq r \leqq (N-1). \quad (17)$$

$(ii)$ $\quad \displaystyle\sum_{r=1}^{N} m_r b^{(r)}(-1) > 0 \geqq \sum_{r=1}^{N} m_r b^{(r)}(0). \quad (18)$

The latter condition follows from the definition of the central step size.

The 0 state of the random walk has the following important property: There is a net drift to the left (right) from states to the right (left) of the 0 state.

$$E[\omega(n+1)|\omega(n) = i] - i = \sum_{r=1}^{N} m_r b^{(r)}(i) < 0 \quad \text{if} \quad i > 0$$
$$> 0 \quad \text{if} \quad i < 0 \Bigg|. \quad (19)$$

The above super- and submartingale properties are the basis for the existence of a stochastic Liapunov function (Appendix A) and the bound given in Section 3.2.

### 1.4 Saturating adaptive quantizer

Any hardware implementation of the quantizers will incorporate some scheme for restricting the range of step sizes. In addition, there are reasons for desiring the step size to be bounded. For instance, by limiting the step sizes at both ends, it is possible to devise automatic schemes for "forgetting" the effects of past channel errors.[9] In such algorithms, the step size may be bounded to fairly small intervals.

For the saturating adaptive quantizer, eq. (3), suppose that

$$\xi_{r-1} C \gamma^{\omega(n)} \leq |x(n)| < \xi_r C \gamma^{\omega(n)}$$

for some $r$, $1 \leq r \leq N$. We obtain the following equation analogous to (14):

$$
\begin{aligned}
\omega(n+1) &= \omega(n) + m_r & \text{if} \quad -K \leq \omega(n) + m_r \leq L \\
&= -K & \text{if} \quad \omega(n) + m_r \leq -K \\
&= L & \text{if} \quad L \leq \omega(n) + m_r,
\end{aligned}
\quad (20)
$$

where $K$ and $L$ are fixed positive integers. The ordinary quantizer is obtained if $K, L \to \infty$.

We observe the following: The central state for the saturating adaptive quantizer may be defined exactly as in the ordinary type of quantizer; the important martingale properties, expressed in eq. (19) for the ordinary quantizer, carry over to the saturating type. The time-dependent transition equations of the saturating quantizer are characterized by numerous involved boundary equations. However, the bulk of the equations are of the form given in (15):

$$
p(i; n+1) = \sum_{r=1}^{N} b^{(r)}(i - m_r) p(i - m_r; n)
$$

$$
-K + m_N \leq i \leq L + m_1. \quad (15')
$$

We do not give the remaining equations since we have no direct need for the time-dependent equations. In Appendix B we give, following the method and notation of Section 2.1, a complete set of reduced equations satisfied by the stationary probabilities.

## II. STATIONARY DISTRIBUTIONS

Appendix A establishes the existence and uniqueness of a finite stationary distribution for the step size in the quantizers. The following sections establish the main qualitative properties of the stationary distributions for both the ordinary and saturating adaptive quantizers.

If we set $p(i; n+1) = p(i; n) = p(i)$ in the time-dependent equations, then the stationary probabilities are given by $\{p(i)\}$. Thus, the stationary probabilities of the ordinary adaptive quantizer are

obtained from

$$p(i) = \sum_{r=1}^{N} b^{(r)}(i - m_r)p(i - m_r), \quad i = 0, \pm 1, \pm 2, \cdots \quad (21)$$

and the normalization equation,

$$\sum_{-\infty}^{\infty} p(i) = 1.$$

### 2.1 A useful reduction of the equations for stationary probabilities

In each equation in (21), the maximum difference in the indices of the state probabilities is $(m_N - m_1)$. By exploiting a property of the stationary distribution, we now obtain a set of new equations where the maximum difference in the indices is $(m_N - m_1 - 1)$. The reduced set of equations together with the normalization equation is complete. A simple interpretation and the motivation of the reduced equation is given in Ref. 1; remark (ii) below gives an additional probabilistic interpretation. The reduced equations are important to us, as they allow us to consider only a smaller set of solutions.

For any integral $j$,

$$\sum_{-\infty}^{j} p(i) = \sum_{i=-\infty}^{j} \sum_{r=1}^{N} b^{(r)}(i - m_r)p(i - m_r)$$

$$= \sum_{i=-\infty}^{j-m_N} \left\{ \sum_{r=1}^{N} b^{(r)}(i) \right\} p(i) + \sum_{i=j-m_N+1}^{j-m_{N-1}} \left\{ \sum_{r=1}^{N-1} b^{(r)}(i) \right\} p(i)$$

$$+ \cdots + \sum_{i=j-m_2+1}^{j-m_1} b^{(1)}(i)p(i).$$

Since $\sum_{r=1}^{N} b^{(r)}(i) = 1$, the above reduces to

$$\sum_{i=j-m_N+1}^{j} p(i) = \sum_{r=1}^{N-1} \sum_{i=j-m_{r+1}+1}^{j-m_r} \left\{ \sum_{s=1}^{r} b^{(s)}(i) \right\} p(i). \quad (22)$$

Define for $1 \leq r \leq N$ and all integral $i$,

$$\psi^{(r)}(i) \triangleq \sum_{s=1}^{r} b^{(s)}(i). \quad (23)$$

The quantities $\{\psi^{(r)}(i)\}$ may be directly obtained from the input distribution, since $\psi^{(r)}(i) = F(\xi_r C\gamma^i)$. From (22) we obtain the reduced equations

$$\boxed{\sum_{i=j-m_N+1}^{j} p(i) = \sum_{r=1}^{N-1} \sum_{i=j-m_{r+1}+1}^{j-m_r} \psi^{(r)}(i)p(i) \quad j = 0, \pm 1, \pm 2, \cdots}.$$

$$(24)$$

In these equations, the set $[j - m_{r+1} + 1, j - m_r]$ is to be treated as empty if $m_r = m_{r+1}$.

*Remarks*:

(*i*) The manipulations leading to (24) are justified since they involve bounded quantities, as is implied by the existence of a unique finite stationary distribution.

(*ii*) Equation (24) is equivalent to the following identity, which is intuitively plausible and may be proven independently:

$$\Pr_s [\omega(n) \le j \text{ and } \omega(n + 1) \ge j + 1]$$
$$= \Pr_s [\omega(n + 1) \le j \text{ and } \omega(n) \ge j + 1],$$

where the subscript $s$ is being used to identify stationary probabilities.

(*iii*) Equation (24) may be used to give a simple proof of an identity (called simply an identity in Ref. 1 and "the design equation" in Ref. 4) involving the stationary state probabilities of the ordinary quantizer. Sum both sides of (24) for all integral $j$:

$$\sum_{j=-\infty}^{\infty} \sum_{i=j-m_N+1}^{j} p(i) = \sum_{j=-\infty}^{\infty} \sum_{r=1}^{N-1} \sum_{i=j-m_{r+1}+1}^{j-m_r} \psi^{(r)}(i) p(i).$$

The left-hand side is simply $m_N$ and the right-hand side is

$$m_N - \sum_{r=1}^{N} m_r q_r,$$

where

$$q_r \triangleq \sum_{i=-\infty}^{\infty} \{ \psi^{(r)}(i) - \psi^{(r-1)}(i) \} p(i).$$

Hence,

$$\sum_{r=1}^{N} m_r q_r = 0. \tag{25}$$

Equation (25) has a natural interpretation if we recognize that $q_r$ is the stationary $r$th step occupancy probability, i.e.,

$$q_r = \Pr_s [\xi_{r-1} \Delta(n) \le |x(n)| < \xi_r \Delta(n)]. \tag{26}$$

The steps leading to eq. (24) may be repeated for the saturating adaptive quantizer, and a similar reduction may be achieved. These equations are given in Appendix B. The main recursion is identical to that of the ordinary quantizer, namely, eq. (24), and holds for all integral $j$, $-K + m_N \le j \le L + m_1 + 1$. Observe that the range over which (24) is valid, for the saturating quantizer, is such that

every state probability is included in at least one component of the recursion.

It may be verified by the reader that the identity in (25), the design equation of Ref. 4, does not hold for the saturating quantizer.

## 2.2 Localization property of the stationary distribution

We prove a fundamental distribution-free property of the stationary distribution of the step size. For both the ordinary and the saturating adaptive quantizers, we obtain sharp geometric bounds on almost all the stationary state probabilities as a function of the distance of the state from the 0 state. The actual bounds obtained are somewhat stronger than the above statement implies, since the rate parameter in the geometric bound itself decreases monotonically with increasing distance from the 0 state. These bounds show that a strong localization of the mass of the stationary distribution about the 0 state (central step size) is inherent in the random walk. Also, we found that it was necessary to prove a result like Theorem 1 before the effects of the multipliers on the dispersion of the stationary distribution could be quantified.

It is necessary to define certain vectors and matrices of dimensions $(m_N - m_1 - 1)$ and $(m_N - m_1 - 1) \times (m_N - m_1 - 1)$, respectively. Let $\mathbf{P}_i$ denote the column vector with the following components:*

$$\mathbf{P}_i \triangleq [p(i), p(i + 1), \cdots, p(i + m_N - m_1 - 2)]^t. \qquad (27)$$

Equation (24) may be used to construct matrices $\{\mathbf{A}_i\}$, which govern the transitions of the above vectors in the following manner:

$$\mathbf{P}_{i+1} = \mathbf{A}_i \mathbf{P}_i. \qquad (28)$$

By examining (24) we observe that the elements of $\mathbf{A}_i$ depend on the quantities $\psi^{(r)}(i), \cdots, \psi^{(r)}(i + m_N - m_1 - 1)$, $1 \leqq r \leqq N$, and the subscript $i$ indicates this dependence.

*Theorem 1 (Localization Property): Let $i > 0$. For both the ordinary and saturating adaptive quantizers, there exists a constant weight vector with positive elements, $\boldsymbol{\lambda}$, and a constant, $r > 1$, depending only on $\mathbf{A}_i$ such that, for all $j \geqq i$,*

$$(\boldsymbol{\lambda}^t \mathbf{P}_j) \leqq \left( \frac{1}{r} \right)^{j-i} (\boldsymbol{\lambda}^t \mathbf{P}_i). \qquad (29)$$

*There exists the $L_1$-norm, $|\mathbf{x}| \triangleq \sum \lambda_k |x_k|$, of the vectors $\{\mathbf{P}_j\}$ which decreases geometrically as $|j - i|$ increases.*

---

* The superscript $t$ denotes the transpose.

An identical statement with $|j - i|$ replacing the index $j - i$ in (29) is also true for $i < 0$ and all $j \leqq i$.

*Remarks*[*]:

(*i*) When $r$ and $\lambda$ in (29) are as constructed by us in the proof of the theorem, then the inequality in (29) becomes an equality if $\mathbf{A}_k = \mathbf{A}_i$ for $k = i, i + 1, \cdots, j$. This indicates that it is not possible to obtain tighter geometric bounds without making further assumptions on the distribution $F(z)$.

Using Theorem 1, we can give the following point-wise bound on the stationary state probabilities for both the ordinary and saturating adaptive quantizers:[†] let $i > 0$; then, for $j \geqq i$

$$
p(j + m_N - m_1 - 2) \leqq \left(\frac{1}{r}\right)^{j-i}(\mathbf{1}^t\mathbf{P}_i) \leqq \left(\frac{1}{r}\right)^{j-i}. \qquad (30)
$$

Similarly, for $i < 0$ and all $j \leqq i$,

$$
p(j - m_N + m_1 + 2) \leqq \left(\frac{1}{r}\right)^{i-j}(\mathbf{1}^t\mathbf{P}_i) \leqq \left(\frac{1}{r}\right)^{i-j}. \qquad (30')
$$

The proof of (30) is as follows. Let $\lambda_m$ denote the largest element of the vector $\lambda$ occuring in Theorem 1 so that $1 \leqq m \leqq m_N - m_1 - 1$. From Theorem 1,

$$
\lambda_m p(j + m - 1) \leqq \lambda^t\mathbf{P}_j \leqq \left(\frac{1}{r}\right)^{j-1}(\lambda^t\mathbf{P}_i) \leqq \left(\frac{1}{r}\right)^{j-i}\lambda_m(\mathbf{1}^t\mathbf{P}_i),
$$

and the inequalities in (30) follow.

*Remarks*:

(*ii*) Observe that for the bounds in (29) and (30) we may use any $i$, $0 < i \leqq j$, as the reference state. The choice of the best reference state depends on the behavior of $r$ with $i$ which, in turn, depends on the distribution $F(z)$. The main distribution-free property of $r(i)$, namely, statement (*iii*) of Lemma 1, indicates an advantage of choosing a large $i$ for the reference state. In Section 2.4, we prove an assertion by implicitly using more than one reference state $i$.

The proof of Theorem 1 relies on two lemmas that we state here and prove in Appendix C.[‡]

---

[*] This remark implies the tightness of the bound in (29), which is lacking for the bound obtained in Ref. 1 for the two-bit quantizer.

[†] The vector $\mathbf{1}$ has every element equal to unity.

[‡] Observe that neither $\mathbf{A}_i$ nor $\mathbf{A}_i^{-1}$ is a nonnegative matrix so that the usual Frobenius theory does not apply.

*Lemma 1*: *For every* $i > 0$,

(*i*) $\mathbf{A}_i$ *is nonsingular and* $\mathbf{A}_i^{-1}$ *has a unique positive real eigenvalue, say,* $r$. *Furthermore,* $r > 1$.

(*ii*) *Every element of the corresponding left eigenvector of* $\mathbf{A}^{-1}$, $\boldsymbol{\lambda}$, *is of the same sign and nonzero, hence we may take* $\boldsymbol{\lambda}$ *to be a positive vector.*

(*iii*) $r$ *which depends on* $i$ *is monotonic, strictly increasing with* $i$.

*Lemma 2*: *For* $j \geqq i > 0$,

$$\boldsymbol{\lambda}^t[\mathbf{A}_j^{-1} - \mathbf{A}_i^{-1}]\mathbf{P}_{j+1} \geqq 0. \tag{31}$$

*Remarks*:

(*iii*) It is not the case that $\boldsymbol{\lambda}^t[\mathbf{A}_j^{-1} - \mathbf{A}_i^{-1}] \geqq 0$, so that (31) is not true if $\mathbf{P}_{j+1}$ is taken to be an arbitrary nonnegative vector.[*] In proving Lemma 2 it is necessary to take into account the fact that the vector $\mathbf{P}_j$, from which $\mathbf{P}_{j+1}$ evolves according to eq. (28), is itself nonnegative, and this implies that $\mathbf{P}_{j+1}$ is restricted to a cone that is a proper subset of the nonnegative quadrant.

*Proof of Theorem 1*: For $j \geqq i > 0$,

$$\begin{aligned}
\boldsymbol{\lambda}^t\mathbf{P}_j = \boldsymbol{\lambda}^t\mathbf{A}_j^{-1}\mathbf{P}_{j+1} &= \boldsymbol{\lambda}^t[\mathbf{A}_j^{-1} - \mathbf{A}_i^{-1}]\mathbf{P}_{j+1} + \boldsymbol{\lambda}^t\mathbf{A}_i^{-1}\mathbf{P}_{j+1} \\
&= \boldsymbol{\lambda}^t[\mathbf{A}_j^{-1} - \mathbf{A}_i^{-1}]\mathbf{P}_{j+1} + r\boldsymbol{\lambda}^t\mathbf{P}_{j+1} \text{ from Lemma 1} \\
&\geqq r\boldsymbol{\lambda}^t\mathbf{P}_{j+1} \qquad\qquad\quad \text{from Lemma 2.} \tag{32}
\end{aligned}$$

Hence, $(\boldsymbol{\lambda}^t\mathbf{P}_j) \leqq (1/r)^{j-i}(\boldsymbol{\lambda}^t\mathbf{P}_i)$ for all $j \geqq i$, as was to be proved.

As every element of $\mathbf{P}_j$ is nonnegative, the $L_1$-norm $|\mathbf{P}_j|$ is equal to $\boldsymbol{\lambda}^t\mathbf{P}_j$. Finally, we may transfer the result that holds for $i > 0$ to the case of $i < 0$ by a simple renumbering of states in the manner that has been indicated in Ref. 1.

The notation common with Ref. 1 conceals some rather significant differences in both the main result (29) and its proof. In Ref. 1, the corresponding result involved $\boldsymbol{\lambda}$ and $r$, which were elements of the eigensystem of an additional matrix $\tilde{\mathbf{A}}_i$ obtained in an involved way from $\mathbf{A}_i$. The result in Lemma 2 has no counterpart in Ref. 1. The geometric bound obtained in Ref. 1 is peculiar to two-bit ($N = 2$) quantizers, and does not directly generalize. Also, the bound obtained here is stronger even for the case $N = 2$.

### 2.3 Lower bounds on the steepness factors, r(i)

Theorem 1 and the subsequent bound in (30) indicates that $r(i)$ is a local measure of the rate with which the stationary probabilities

---

[*] A vector is nonnegative if every element is nonnegative. The nonnegative quadrant in $R^n$ is the set of all nonnegative vectors of dimension $n$.

change, and for this reason we find it natural to call $r(i)$ the local steepness factor. Here we go back to the definition of $r(i)$ as being the unique positive real root of the polynomial $C(\mu)$, eq. (60), to obtain the following bound on $r(i)$, which has the advantages of being explicit and being dependent only on the transition probabilities at state $i$. We make free use of this bound in the following section.

$$
r(i) \geqq \rho(i) \triangleq \left[ \frac{\sum\limits_{r=1}^{\mu} (-m_r)\{\psi^{(r)}(i) - \psi^{(r-1)}(i)\}}{\sum\limits_{r=\mu+1}^{N} m_r\{\psi^{(r)}(i) - \psi^{(r-1)}(i)\}} \right]^{1/m_N - m_1 - 1}, \quad (33)
$$

where, of the $N$ multipliers, only $\mu$ multipliers have values not exceeding unity, i.e.,

$$
m_1, m_2, \cdots, m_\mu \leqq 0
$$

and

$$
m_{\mu+1}, m_{\mu+2}, \cdots, m_N > 0.
$$

The bound $\rho(i)$ has certain interesting properties. First, observe that, by virtue of the definition of the central state [eqs. (8) and (9)], $\rho(i) > 1$ for all $i > 0$. Also, the sequence $\rho(i)$, is, like $\{r(i)\}$, monotonic, increasing with $i$. The numerator and denominator of the bracketed expression have interesting probabilistic interpretations: The numerator (denominator) is the expected change in state conditional on the transition being from state $i$ to all states $i' \leqq i (i' > i)$.

The proof of eq. (33) is involved, and for the sake of brevity we omit giving it.

### 2.4 Effect of $\gamma$ on the stationary distribution

We show in this section that the mass of the stationary distribution of the step size can be concentrated about the central step size to an arbitrary extent by making $\gamma$ sufficiently close to unity. To show this, we first put together, from the results of the preceding two sections, a rather explicit bound on the stationary probability of the step size exceeding a particular value for a given $\gamma$, i.e., $\Pr_s[\Delta > C\gamma^i]$. This bound is in a form that allows direct comparison with the corresponding probability arising from the choice of $\gamma' = \sqrt{\gamma}$. By successively taking $\gamma$ to be the square root of the preceding value, the bound on the probability can be made as small as desired. This procedure for proving the assertion is similar to the one we developed in Ref. 1. We restrict our attention to step sizes that exceed the central step size, i.e., $i > 0$, since a parallel argument holds for $i < 0$.

In the following discussion the quantity $(m_N - m_1 - 2)$ arises frequently, and it is convenient to denote this quantity by the symbol

$\nu$. Clearly, $\nu$ is a measure of the spread in the log of the multipliers. For $i > 0$ and $r = r(i)$, we have from eq. (29) that

$$(\Sigma\lambda_i) \sum_{j=i+\nu}^{\infty} p(j) \le \sum_{j=i}^{\infty} \lambda^i \mathbf{P}_j \le \lambda^i \mathbf{P}_i \sum_{j=0}^{\infty} \left(\frac{1}{r}\right)^j = \lambda^i \mathbf{P}_i \frac{r}{r-1}. \quad (34)$$

Now

$$r \ge \rho(i), \quad (35)$$

where $\rho(i)$ is defined in eq. (33), and

$$\frac{\lambda^i \mathbf{P}_i}{\Sigma\lambda_i} \le \max\,[p(i), \cdots, p(i+\nu)].$$

Since

$$\Pr{}_s\,[\Delta \ge C\gamma^{i+\nu}] = \sum_{j=i+\nu}^{\infty} p(j),$$

we have, from eq. (34),

$$\boxed{\Pr{}_s\,[\Delta \ge C\gamma^{i+\nu}] \le \frac{\rho(i)}{\rho(i)-1} \max\,[p(i), \cdots, p(i+\nu)]}. \quad (36)$$

Finally, from Eq. (30), for $i \ge \nu + 1$,

$$\boxed{\max\,[p(i), \cdots, p(i+\nu)] \le \left[\frac{1}{\rho(1)}\right]^{i-\nu-1}}. \quad (37)$$

Equations (36) and (37) together give us the desired bound on the stationary probability of the step size exceeding a given value, which we now compare with a similar bound that holds for $\gamma' = \sqrt{\gamma}$. The prime superscript is used on symbols to denote the functional dependence of the associated quantities on $\gamma'$. In establishing the reference, i.e., central, step size corresponding to $\gamma'$, minor differences exist depending on whether [see eqs. (8) and (9)]

(i) $$B(\gamma^{i-1}) > 0 \ge B(\gamma^{i-\frac{1}{2}})$$

or

(ii) $$B(\gamma^{i-\frac{1}{2}}) > 0 \ge B(\gamma^i). \quad (38)$$

We consider only (ii), in which case: $\omega'(n) = 2i \iff \omega(n) = i$, and all the transition probabilities are simply related: $\psi^{(r)}(2i)' = \psi^{(r)}(i)$. As a consequence of the latter property, we have

$$\rho'(2i) = \rho(i). \quad (39)$$

Repeating the arguments leading to eqs. (36) and (37), we have

$$\mathrm{Pr}'_s \left[ \Delta \geq C\sqrt{\gamma}^{2i+\nu} \right] \leq \frac{\rho'(2i)}{\rho'(2i) - 1} \max \left[ p'(2i), \cdots, p'(2i + \nu) \right] \quad (40)$$

and

$$\max \left[ p'(2i), \cdots, p'(2i + \nu) \right] \leq \left[ \frac{1}{\rho'(2)} \right]^{2i-\nu-2}. \quad (41)$$

By the fact that $\rho'(2i) = \rho(i)$, we have

$$\mathrm{Pr}'_s \left[ \Delta \geq C\sqrt{\gamma}^{2i+\nu} \right] \leq \frac{\rho(i)}{\rho(i) - 1} \left[ \frac{1}{\rho(1)} \right]^{i-\nu-1} \left[ \frac{1}{\rho(1)} \right]^{i-1}. \quad (42)$$

Comparison with eqs. (36) and (37) completes the demonstration.

## III. TRANSIENT RESPONSE

In this section, we are interested in the random time, called the adaptation time, taken for the step size of the quantizer to adapt from some arbitrary initial value to the central step size. It is necessary to have the adaptation time relatively small if the quantizer is to adequately track the scale variations of the input process. Also, it is reasonable to expect that, as $\gamma$ is made large, the increased range of the multipliers [eq. (7)] will give the desired tracking. However, as a counterbalance, we already know from the preceding section that, with the correct choice of the log of the multipliers, $\{m_i\}$, the quality of steady-state performance is increasingly impaired as the value of $\gamma$ is raised. From this brief discussion (see Ref. 1 for a more detailed discussion), it is clear that it is useful to have formulas for the efficient computation of the mean adaptation time and bounds that provide insight on the dependence of the time on the multipliers.

### 3.1 Mean time for first passage to the central state

We consider only the saturating adaptive quantizer since, as $K$ and $L$ are made large, the quantities obtained for this model approximate corresponding quantities for the ordinary adaptive quantizer. Also, for the usual reason only the case of positive initial states, $\omega(0) > 0$, is considered.

Let the initial step $\omega(0) = i > 0$ and let $T(i)$ denote the mean value of the random time $\tau$ where $\omega(\tau) \leq 0$ and $\omega(n) > 0$ for all $n < \tau$. It can be shown that, as a consequence of the recurrence and irreducibility of the Markov chain (see Appendix A), the mean first passage time, $T(i)$, is finite with probability 1. If the first transition results in a transition to the state $i + m_r$, the process continues as if the initial state had been $i + m_r$. The conditional expectation of the first passage time is therefore $T(i + m_r) + 1$. From this argument, we

deduce that the following recursion is satisfied by the mean first passage time,

$$T(i) = \sum_{r=1}^{N} b^{(r)}(i)\{T(i + m_r) + 1\} \quad -m_1 + 1 \leq i \leq L - m_N \quad , \quad (43)$$

where, as in eq. (16), $b^{(r)}(i) = F(\xi_r C \gamma^i) - F(\xi_{r-1} C \gamma^i)$. Of course, $\sum_{r=1}^{N} b^{(r)}(i) = 1$. The recursive relation in (43) may be used to generate the entire sequence $\{T(i)\}$, provided $(m_N - m_1)$ initial conditions can be found. Now, by the same argument that led to eq. (43), we have

$$T(1 + m_1) = T(2 + m_1) = \cdots = T(0) = 0. \quad (44)$$

The remaining $m_N$ initial conditions, namely,

$$T(1), T(2), \cdots, T(m_N),$$

are harder to obtain, and it is necessary to look more deeply into the dynamics of the process to obtain these quantities.

For every time instant, we define the $L$-dimensional vector $\mathbf{z}(n)$ with components $z(j; n)$, $1 \leq j \leq L$, where

$$z(j; n) \triangleq \Pr[\omega(n) = j \text{ and } \omega(s) \geq 1 \text{ for all } s \leq n]. \quad (45)$$

We show in Appendix D that the vectors $\mathbf{z}(n)$ evolve in time according to the homogeneous equation

$$\mathbf{z}(n + 1) = \mathbf{D}\mathbf{z}(n), \quad n \geq 0, \quad (46)$$

where $\mathbf{D}$ is an $L \times L$ matrix. Also, in Appendix D we prove the following: For $i \geq 1$,

$$\begin{array}{c} T(i) = \sum_{j \geq 1} x_j^{(i)}, \\ \text{where} \\ [\mathbf{I} - \mathbf{D}]\mathbf{x}^{(i)} = \mathbf{e}^{(i)} \end{array} \quad , \quad (47)$$

and the elements of the $L$-vector $\mathbf{e}^{(i)}$ are zero everywhere except at the $i$th location where the element is unity. It is shown in Appendix D that $[\mathbf{I} - \mathbf{D}]$ is nonsingular.

The simple recursion in (43) may be used to generate the sequence $\{T(i)\}$ after obtaining the nonzero initial conditions via $m_N$ inversions, as in (47). Alternatively, if $T(i)$ is required for only a few particular values of $i$, it may be easier to obtain them via the inversions in (47).

The bulk of the equations in (47) [see eq. (72)] are in the form encountered in the analysis of the stationary distribution, eq. (21).

Also, the elements of the vectors $x_j^{(i)}$ are all nonnegative. Hence, by applying the techniques and results of the preceding section, we may draw certain conclusions about eq. (47).

First, the bandwidth of the matrix $[I - D]$ may be reduced by 1 by carrying out the reduction of the equations described in Section 2.1. For $m_1 = -1$ and arbitrary values of $m_2, \cdots, m_N$, this step is enough to triangularize the matrix $[I - D]$ for any countable $L$ and thus substantially simplify the computations. Second, we may conclude from Section 2.2 that, with increasing $j$, the solution elements $x_j^{(i)}$ decrease at least geometrically. This is a very useful property from the point of view of numerical inversion of $[I - D]$ for $L$ large and the approximation of the solution for $L = \infty$ by finite $L$.

### 3.2 A bound on the mean first passage time

Let $T(i, j)$, $0 \leq i < j$, denote the following mean first passage time: the initial state $\omega(0) = j$, first crossing occurs after $\tau$ transitions if $\omega(\tau) \leq i$, and $\omega(n) > i$ for all $n < \tau$, and $T(i, j) = E(\tau)$. The quantity $T(j)$ of the preceding section is equivalent in our present notation to $T(0, j)$. We now give an explicit bound on $T(i, j)$ that provides some insight into the dependence of $T(i, j)$ on the multipliers.

For both the ordinary and saturating adaptive quantizer,

$$
T(i, j) \leq \frac{1}{C(i + 1)} \left[ (j - i) - (m_1 + 1) \right] \quad 0 \leq i < j,
$$

where

$$
C(i) = \sum_{r=1}^{N-1} (m_{r+1} - m_r)\psi^{(r)}(i) - m_N \tag{48}
$$

From the definition of the central state, eq. (18), and the monotonicity of $\psi^{(r)}(i)$ with respect to $i$, we observe that for $i > 0$, $C(i)$ is positive, monotonic, increasing with $i$. We only sketch the proof of (48) because the method of the proof is contained in the proof of the bound that we gave in Ref. 1 for the two-bit quantizer. First, recall [eq. (19)] that a supermartingale property exists that holds for both types of quantizers, according to which there is a net drift to the left from all states $j > 0$. Second, we define a new process in which $\omega'(n) = \omega(n) + nC(i + 1)$ and show that the supermartingale property, i.e., $E[\omega'(n + 1) | \omega'(n)] \leq \omega'(n)$, is preserved for the range of $n$ of interest. Finally, an application of Doob's theorem on optional stopping of supermartingales[10] on the new process yields the bound in eq. (48).

The bound provides some insight into the dependence of the mean adaptation times on the multipliers, and $\gamma$ in particular, when the

initial and final step sizes are $C\gamma^j$ and $C$, respectively. Briefly, consider the effect of making $\gamma' = \sqrt{\gamma}$, i.e., $M_i' = \sqrt{M_i}$ and the spread of the multipliers is reduced. The number of states between the states corresponding to $C\gamma^j$ and $C$ is doubled. Now $C(1)$ is hardly affected by the transformation and, as a consequence of the linear dependence of the bound on $T(i, j)$ on the distance $(j - i)$, we have the bound on the mean adaptation time approximately doubled. For $i = 0$ and $j \gg (-m_1)$, computations amply corroborate this conclusion.

## IV. COMPUTATIONAL RESULTS

We present here a sampling of rather extensive computations done on three- and four-bit adaptive quantizers ($N = 4$ and 8, respectively) for independent identically distributed input sequences with gaussian distributions. Both uniform, i.e., $\xi_i = i$, and nonuniform quantizers were considered. Max[8] has shown in the nonadaptive framework that optimal nonuniform quantizers can yield an improvement in the signal-to-noise ratio of about 20 percent over optimal uniform quantizers with the number of bits in the range of interest here. We note that four-bit adaptive quantizers have been breadboarded in Bell Laboratories,[3] and that Jayant's[2] systematic numerical study is restricted to uniform quantizers up to three bits. We also observe that a simple search procedure of the "optimal" set of multipliers grows to be almost unmanageable and expensive when the dimension of the parameter spaces is 8.

Table I lists five quantizers with their respective parameters $\{m_i\}$. The parameter $\gamma$ is not considered part of the characterization of the quantizer type. Among the quantizers investigated, the following five proved to be the most interesting in their respective classes, specified by number of bits and uniform or nonuniform. The first of the five, with $\gamma \approx 1.12$, is close to what Jayant calls the optimal, three-bit quantizer. The parameters $\{m_i\}$ were arrived at by the procedure described in remark $(i)$, Section 1.2.

## Table I — Five quantizers

| Specifications | | | Designation |
|---|---|---|---|
| Uniform or Nonuniform | Number of Bits | $\{\log_\gamma (M_i)\} : m_1, \cdots, m_N$ | |
| Uniform | 3 | $-1, -1, 2, 5$ | $UQ$, 3 bits, No. 1 |
| Uniform | 3 | $-1, 0, 1, 4$ | $UQ$, 3 bits, No. 2 |
| Nonuniform | 3 | $-2, -1, 2, 8$ | $NUQ$, 3 bits |
| Uniform | 4 | $-2, -2, 0, 0, 2, 5, 10, 17$ | $UQ$, 4 bits |
| Nonuniform | 4 | $-2, -2, 0, 0, 1, 2, 5, 16$ | $NUQ$, 4 bits |

The optimum division of the horizontal axis in Fig. 1, given by $\xi_i$, $i = 1, 2, \cdots, (N - 1)$, was obtained from Max,[8] and we reproduce these parameters for the reader's benefit.

*NUQ, 3 bits.* $\{\xi_i\} = \{1.0, 2.097, 3.492\}$.
*NUQ, 4 bits.* $\{\xi_i\} = \{1.0, 2.023, 3.097, 4.256, 5.565, 7.142, 9.299\}$.

Table II lists some statistics of the stationary step-size distribution for unit variance of the input distribution. The stationary distribution was obtained by solving the stationary equations of the saturating adaptive quantizers with suitably large saturating levels $(K + L \approx 100)$. We also give the stationary step-occupancy probabilities $q_i$, where $q_i = \text{Pr}_s\left[\xi_{i-1}\Delta(n) \leqq |x(n)| < \xi_i\Delta(n)\right]$, as in eq. (26). Table II also gives, for purposes of comparison, corresponding quantities of the optimal nonadaptive quantizer obtained from Max.[8] In particular, $\hat{\Delta}$ is the optimal, nonadaptive step size.

Figures 2 to 5 show the mean adaptation times for inputs with unit variance. Figures 2 and 3 are concerned with the three types of three-bit quantizers for various values of $\gamma$. These figures plot the mean time taken by the quantizers to adapt to the central, and optimal, step size for various values of the initial step size. In Fig. 2, the initial step size exceeds the central step size, while the reverse case is considered in Fig. 3. Similarly, Figs. 4 and 5 plot data on the mean adaptation times for the uniform and nonuniform four-bit quantizers.

The purpose of the remaining tables (III to V) is to give the reader a feel for the relative performance of the five quantizers. We measure performance by the ratio of the input signal energy to the quantization

## Table II — Statistics of the stationary step-size distributions

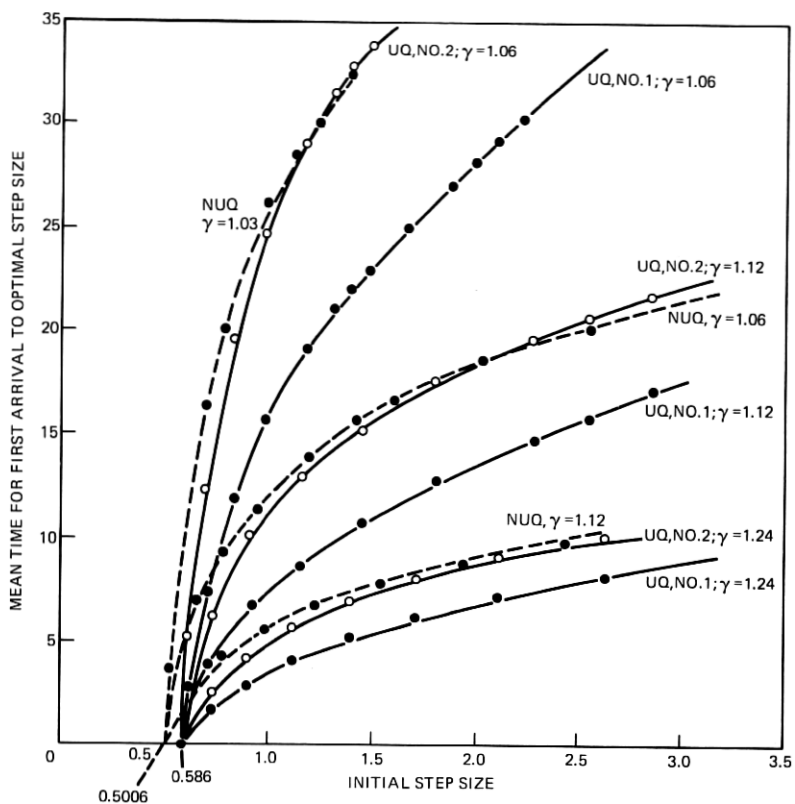| Type | $\gamma$ | $\hat{\Delta}$ (Max) | $E(\Delta)$ | $\sigma(\Delta)$ | Step Occupancy Probabilities {adaptive quantizer} {optimal nonadaptive quantizer} |
|---|---|---|---|---|---|
| *UQ*, 3 bits No. 1 | 1.04 | 0.586 | 0.594 | 0.105 | {0.445, 0.310, 0.156, 0.089} {0.442, 0.317, 0.162, 0.078} |
| *UQ*, 3 bits No. 2 | 1.04 | 0.586 | 0.613 | 0.089 | {0.458, 0.314, 0.152, 0.075} {0.442, 0.317, 0.162, 0.078} |
| *NUQ*, 3 bits | 1.04 | 0.501 | 0.522 | 0.114 | {0.396, 0.317, 0.198, 0.088} {0.383, 0.323, 0.213, 0.081} |
| *UQ*, 4 bits | 1.04 | 0.335 | 0.366 | 0.095 | {0.285, 0.244, 0.182, 0.121, 0.075, 0.043, 0.024, 0.027} {0.263, 0.235, 0.188, 0.135, 0.086, 0.049, 0.025, 0.019} |
| *NUQ*, 4 bits | 1.04 | 0.258 | 0.279 | 0.066 | {0.219, 0.205, 0.178, 0.145, 0.110, 0.076, 0.045, 0.022} {0.204, 0.195, 0.177, 0.152, 0.121, 0.086, 0.049, 0.016} |

Fig. 2—Transient response of three three-bit quantizers.

error energy. Unlike all previous data, the data for these tables were obtained by Monte Carlo simulation. The interval of time over which performance was monitored is denoted by $NA$. Thus, signal energy is $\sum_{n=1}^{NA} x^2(n)$. The remaining parameter in the tables is the initial step size, $\Delta$ (initial). However, we do not list the raw initial step size, but

Table III* — S/N performance of two uniform three-bit quantizers
(Main numbers are for *UQ*, three bits, No. 1; numbers in ( )
for *UQ*, three bits, No. 2)

| Log $\{\Delta(\text{initial})/\hat{\Delta}\}$ | $NA = 10$ | $NA = 100$ | $NA = 1000$ | $NA = 10,000$ |
|---|---|---|---|---|
| $-1$ | 6.92  (5.84) | 14.4 (14.8) | 17.4 (19.3) | 17.7 (20.1) |
| 0 | 25.7  (27.6) | 19.1 (21.4) | 17.9 (20.4) | 17.8 (20.2) |
| 1 | 0.549 (0.549) | 3.94 (3.99) | 13.1 (14.3) | 17.1 (19.2) |

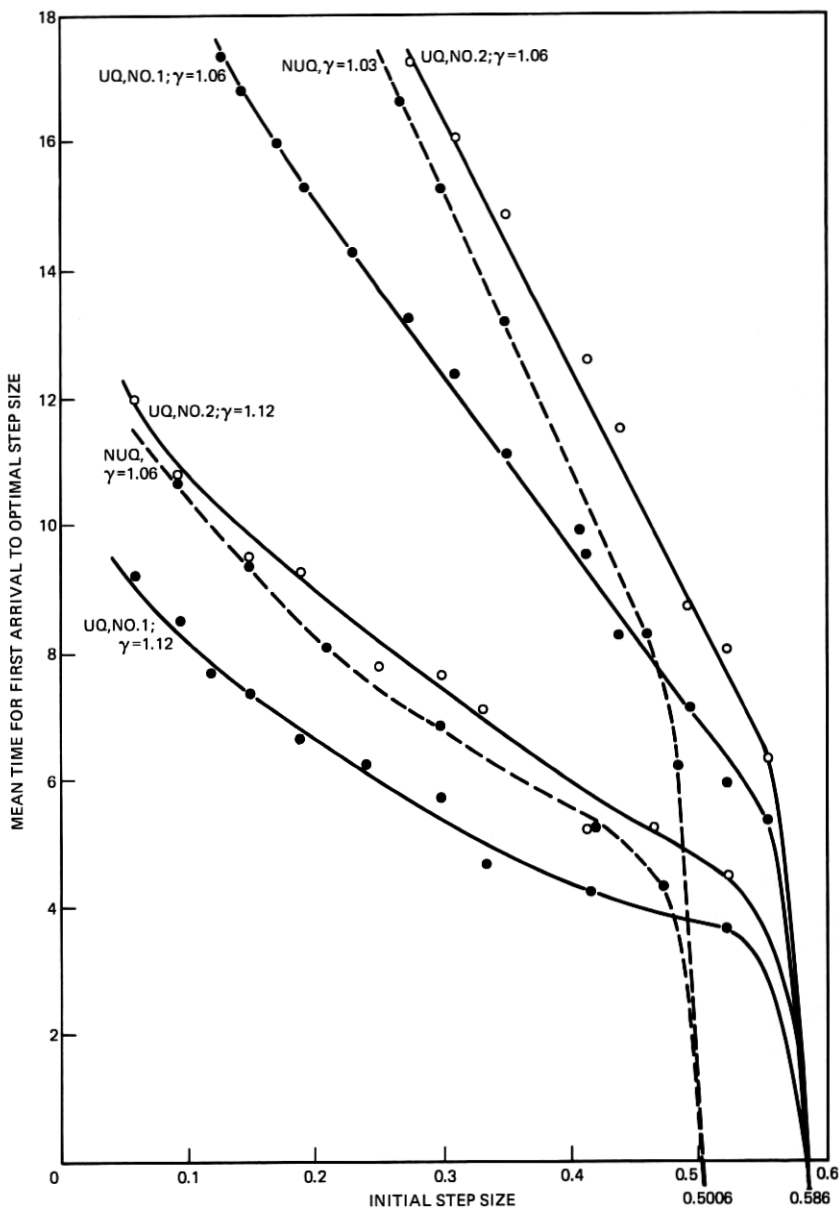* All logarithms in Tables III, IV, and V have base 10.

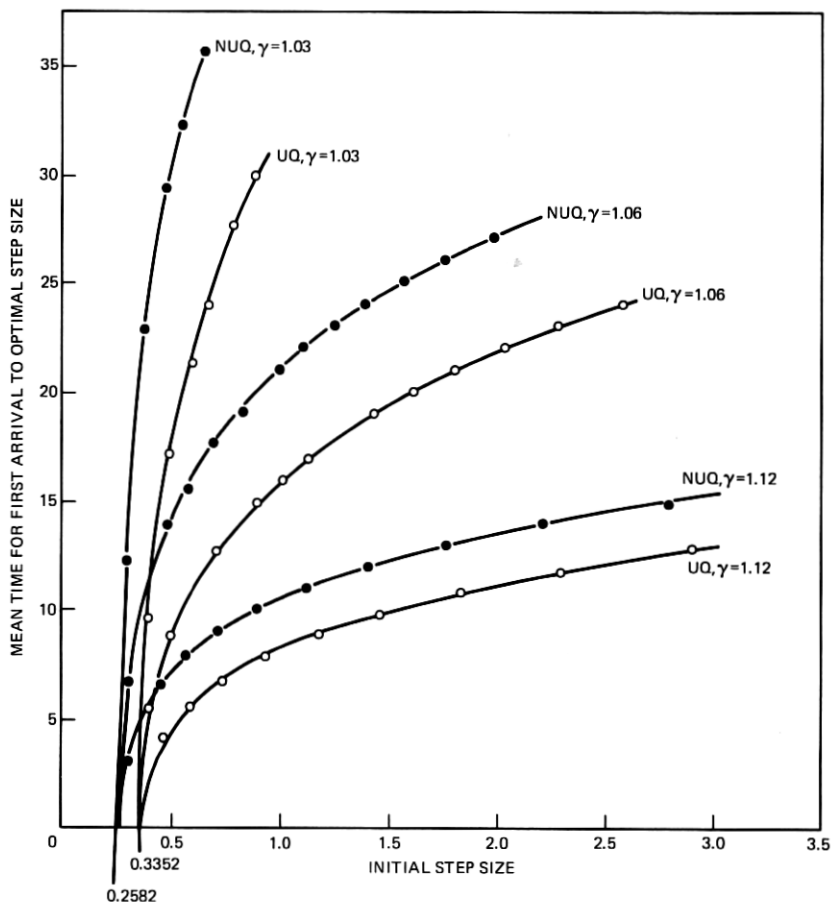Fig. 3—Transient response of three three-bit quantizers.

Fig. 4—Transient response of two four-bit quantizers.

the more relevant quantity $\Delta(\text{initial})/\hat{\Delta}$ where $\hat{\Delta}$ is, as usual, the optimal nonadaptive step size. After experimenting, we arrived at the following values of $\gamma$ for the five quantizers, since they gave a suitable mix of performances over short ($NA$ small) and long ($NA$ large) runs.

## Table IV — S/N performance of nonuniform three-bit quantizer (*NUQ*, three bits)

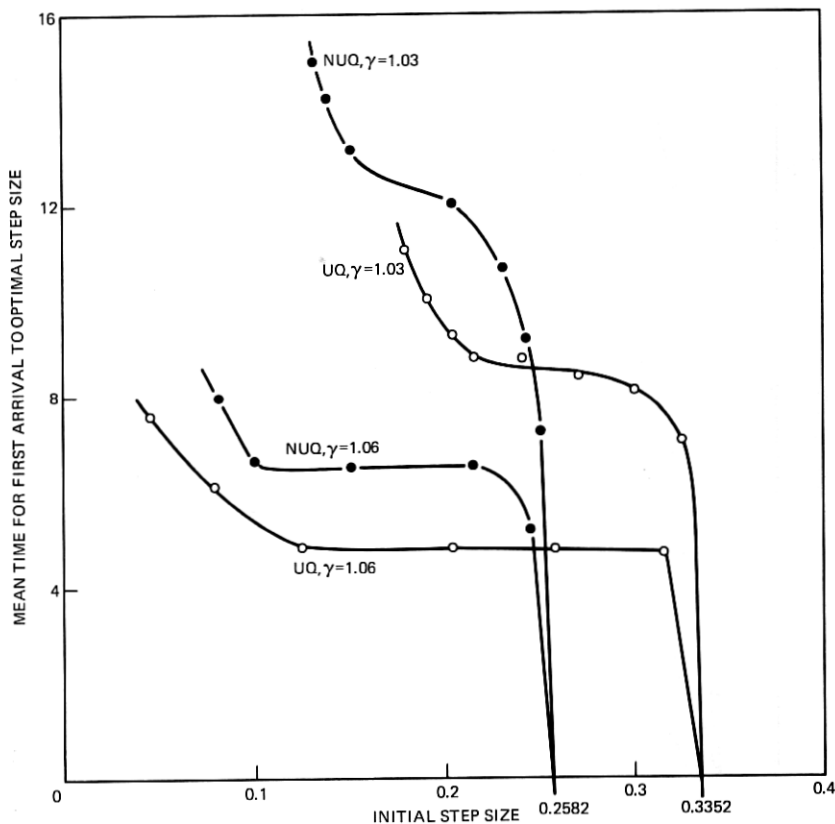| Log $\{\Delta(\text{initial})/\hat{\Delta}\}$ | $NA = 10$ | $NA = 100$ | $NA = 1000$ | $NA = 10{,}000$ |
|---|---|---|---|---|
| −1 | 5.81 | 16.0 | 21.2 | 22.0 |
| 0 | 29.8 | 23.8 | 22.4 | 22.1 |
| 1 | 1.12 | 7.00 | 18.2 | 21.6 |

Fig. 5—Transient response of two four-bit quantizers.

For a particular input process, the relative weightings may be quite different, and $\gamma$ may then be tuned accordingly.

| Quantizer | $\gamma$ |
|---|---|
| $UQ$, 3 bits, No. 1 | 1.12 |
| $UQ$, 3 bits, No. 2 | 1.12 |
| $NUQ$, 3 bits | 1.06 |
| $UQ$, 4 bits | 1.06 |
| $NUQ$, 4 bits | 1.06 |

The following observations may be made on the above results. There is a pronounced asymmetry in performance with respect to $\log \{\Delta(\text{initial})/\hat{\Delta}\}$ over short runs ($NA = 10$ or $100$). This is, of

## Table V — S/N performance of uniform and nonuniform four-bit quantizers
### (Main numbers are for *UQ*, four bits; numbers in ( ) for *NUQ*, four bits)

| Log $\{\Delta(\text{initial})/\hat{\Delta}\}$ | $NA = 10$ | $NA = 100$ | $NA = 1000$ | $NA = 10,000$ |
|---|---|---|---|---|
| $-1$ | 19.62 (21.65) | 36.98 (47.30) | 48.22 (67.35) | 48.97 (71.50) |
| 0 | 86.2 (111.0) | 56.0 (80.1) | 50.60 (72.50) | 49.20 (71.90) |
| 1 | 2.97 (4.86) | 17.7 (27.6) | 42.00 (62.00) | 48.10 (70.30) |

course, related to the contraction multipliers being grossly smaller than the expansion multipliers in all the quantizers considered (Table I). The s/n when $\Delta(\text{initial})/\hat{\Delta} = 1$ and $NA = 10$ is close to the s/n obtained with the step size optimally tuned to the known level of scaling of the input sequence. The steady but not excessive deterioration in performance with increasing $NA$ is the price paid for adaptability: it is due to the fluctuations in step size arising from the random walk. Finally, we observe from Table V that there is a striking gain from nonuniform quantization, the extent of the gain being somewhat greater than what may be expected from previous results on nonadaptive quantizers.

## APPENDIX A
### Existence and Uniqueness of the Stationary Distribution

We establish in this appendix that, for independent identically distributed inputs, there exists a unique, finite stationary step-size distribution (invariant measure). The proof given here is via the construction of a stochastic Liapunov function, and it relies on a standard, unified theory of stochastic stability[11,12] that is well-known. The stochastic stability of the adaptive quantizer has been proved by Goodman and Gersho,[4] and the prime reason for including an alternative proof is our belief that familiarity with the method followed here may be beneficial to future workers in adaptive processes. The positive function that is proved to be a stochastic Liapunov function here is identical to the function that worked in Ref. 1 for the two-bit quantizer, and the proof is a straightforward generalization.

We consider in turn two properties of well-behaved Markov chains, namely, irreducibility and recurrence.

### A.1 Irreducibility

The Markov chain is irreducible if and only if every state communicates with both neighboring states. This occurs if and only if

there exist nonnegative integers $n_i$ and $n_i'$, $1 \leq i \leq N$, such that

$$\Sigma m_i n_i = 1 \tag{49}$$

and

$$\Sigma m_i n_i' = -1. \tag{50}$$

It is an elementary fact from Euclid's theory that this occurs if and only if the integers $\{m_i\}$ are relatively prime, i.e., their greatest common divisor is unity.

### A.2 Recurrence

Consider the following nonnegative function of the states

$$V(i) \triangleq |i| \quad i = 0, \pm 1, \pm 2, \cdots. \tag{51}$$

Let $D(i)$ be defined as

$$D(i) \triangleq E[V\{\omega(n+1)\} | \omega(n) = i] - V(i). \tag{52}$$

Now $D(i)$ is uniformly bounded from above. By the monotonicity of $\psi^{(r)}(i)$ with respect to $i$ and the definition of the central state, (18), we obtain, for all $i \geq (-m_1)$,

$$D(i) = m_N - \sum_{r=1}^{N-1} (m_{r+1} - m_r)\psi^{(r)}(i)$$

$$\leq m_N - \sum_{r=1}^{N-1} (m_{r+1} - m_r)\psi^{(r)}(-m_1) < 0 \tag{53}$$

and, for all $i \leq -m_N$,

$$D(i) = -m_N + \sum_{r=1}^{N-1} (m_{r+1} - m_r)\psi^{(r)}(i)$$

$$\leq -m_N + \sum_{r=1}^{N-1} (m_{r+1} - m_r)\psi^{(r)}(-m_N) < 0, \tag{54}$$

where, as in eq. (23), $\psi^{(r)}(i)$ denotes $F(\xi, C\gamma^i)$. Hence, by virtue of eqs. (53) and (54), $D(i) \leq -\epsilon < 0$ for all but a finite set of states $i$, and $V(i)$ is a stochastic Liapunov function for the process.

From Kushner's Theorem 7,[11] we have recurrence and we can infer further, from Theorem 4, that there exists at least one finite invariant measure, i.e., stationary distribution. Also, since we have shown earlier that two or more disjoint self-contained subsets of the state space do not exist, we have, from Theorem 5, at most one invariant probability measure. The existence and uniqueness of a finite stationary

distribution for the step size of the ordinary adaptive quantizer is therefore established.

### A.3 The saturating adaptive quantizer

The argument leading to irreducibility is intact. In addition, we have here that the end states $(-K)$ and $L$ have period 1 and, since periodicity is a class concept (i.e., every state in a particular communicating class has the same periodicity), the entire Markov chain is aperiodic and, consequently, there is a single ergodic class that includes every state in the chain. Hence, the distribution at time $n$, $p(n)$ approaches $p$, the stationary distribution for all initial distributions, and furthermore every component probability of $p$ is strictly positive.

### APPENDIX B

#### The Saturating Adaptive Quantizer

We give in this appendix a set of equations satisfied by the stationary probabilities of the states in the saturating adaptive quantizer. These equations are complete and reduced by the method described in Section 2.1.

Let $\mu$ denote the number of contraction multipliers, i.e., multipliers having values less than 1, so that

$$m_1, \cdots, m_\mu < 0 < m_{\mu+1}, \cdots, m_N. \tag{55}$$

The tacit assumption that there are no multipliers exactly equal to unity is by no means necessary, but does lead to a simpler presentation.

The main set of equations is

$$\sum_{i=j-m_N+1}^{j} p(i) = \sum_{r=1}^{N-1} \sum_{i=j-m_r+1+1}^{j-m_r} \psi^{(r)}(i)p(i),$$
$$-K + m_N - 1 \leq j \leq L + m_1. \tag{56}$$

The lower boundary equations are*

$$\sum_{i=-K}^{j-1} p(i) = \sum_{r=1}^{s-1} \sum_{i=-K \wedge (j-m_{r+1})}^{j-m_r-1} \psi^{(r)}(i)p(i), \tag{57}$$

where $\mu + 1 \leq s \leq N$ and $-K + m_{s-1} + 1 \leq j \leq -K + m_s$. Finally,

---

* $x \wedge y = \text{Max}\,[x, y]$ and $x \vee y = \text{Min}\,[x, y]$.

the upper boundary equations are

$$\sum_{i=j-m_N+1}^{j} p(i) = \sum_{r=s}^{N-1} \sum_{i=j-m_{r+1}+1}^{L \vee (j-m_r)} \psi^{(r)}(i)p(i), \qquad (58)$$

where $1 \leqq s \leqq \mu$ and $L + m_s \leqq j \leqq L + m_{s+1} - 1$.

## APPENDIX C

### Proofs of Lemmas 1 and 2

### C.1 Proof of Lemma 1

($i$) It can be shown that the determinant of the matrix $\mathbf{A}_i$,

$$\det [\mathbf{A}_i] = (-1)^{m_N - m_1}[1 - \psi^{(N-1)'}(i)]/\psi^{(1)}(i + m_N - m_1 - 1).$$

As $\det [\mathbf{A}_i] > 0$, $\mathbf{A}_i^{-1}$ exists.

Since $\mathbf{P}_i = \mathbf{A}_i^{-1}\mathbf{P}_{i+1}$, we observe from the structures of $\mathbf{P}_i$ and $\mathbf{P}_{i+1}$ that the matrix $\mathbf{A}_i^{-1}$ is in companion form in that all rows except the first reflect shift operations, i.e., for $k \geqq 2$,

$$\begin{aligned} [\mathbf{A}_i^{-1}]_{k,l} &= 0 \quad \text{if} \quad l \neq (k-1) \\ &= 1 \quad \text{if} \quad l = (k-1). \end{aligned} \qquad (59)$$

The elements of the first row of $\mathbf{A}_i^{-1}$ are obtained from the equation

$$\sum_{l=0}^{m_N-1} p(i+l) - \sum_{r=1}^{N-1} \sum_{l=m_N-m_{r+1}}^{m_N-m_r-1} \psi^{(r)}(i+l)p(i+l) = 0. \qquad (24)$$

As the matrix $\mathbf{A}_i^{-1}$ is in companion form, we know that its characteristic polynomial is equal to within a constant of proportionality to the polynomial obtained by replacing, in eq. (24), $p(i+l)$ by $\mu^{m_N-m_1-1-l}$. That is, where

$$C(\mu) \triangleq (-1)^{m_N - m_1 - 1} \det [\mathbf{A}_i^{-1} - \mu \mathbf{I}],$$

we have

$$[1 - \psi^{(N-1)}(i)]C(\mu) = \sum_{l=0}^{m_N-1} \mu^{m_N-m_1-1-l}$$

$$- \sum_{r=1}^{N-1} \sum_{l=m_N-m_{r+1}}^{m_N-m_r-1} \psi^{(r)}(i+l)\mu^{m_N-m_1-1-l}. \qquad (60)$$

The quantity $[1 - \psi^{(N-1)}(i)]$ is merely the coefficient of $p(i)$ in eq. (24).

Scanning the coefficients of the polynomial $C(\mu)$, we observe that there is a single-sign alternation and, hence, by Descartes' rule, $C(\mu)$

has at most one real positive root. Since

$$C(0) = -\psi^{(1)}{}'(i + m_N - m_1 - 1)/[1 - \psi^{(N-1)}(i)] < 0$$

and $C(\mu) \to \infty$ as $\mu \to \infty$, there exists exactly one real positive root. Let $r$ denote this root.

Now

$$[1 - \psi^{(N-1)}(i)]C(1) = m_N - \sum_{r=1}^{N-1} \sum_{l=m_N-m_r+1}^{m_N-m_r-1} \psi^{(r)}(i + l)$$

$$< m_N - \sum_{r=1}^{N-1} \sum_{l=m_N-m_r+1}^{m_N-m_r-1} \psi^{(r)}(i)$$

$$= \sum_{r=1}^{N} m_r\{\psi^{(r)}(i) - \psi^{(r-1)}(i)\}, \qquad (61)$$

where we have followed the usual convention in setting $\psi^{(N)}(i) = 1$ and $\psi^{(0)}(i) = 0$. So $C(1) < 0$ if $\sum_{r=1}^{N} m_r\{\psi^{(r)}(i) - \psi^{(r-1)}(i)\} \leq 0$. The latter condition holds for all $i \geq 0$ [see eqs. (17) and (18)]. Hence, $r > 1$.

(ii) Let us denote the elements of the first row of $\mathbf{A}_i^{-1}$ by $\{\alpha_l\}$ and $\{\beta_l\}$ so that the row appears as

$$[-\alpha_1 - \alpha_2 \cdots - \alpha_{m_N-1}\beta_1\beta_2 \cdots \beta_{-m_1}]. \qquad (62)$$

One reason for expressing the row in this manner is that every $\alpha_l$ and $\beta_l$ is strictly positive by eq. (24).

The left eigenvector $\lambda$ of $\mathbf{A}_i^{-1}$ corresponding to the eigenvalue $r$ satisfies, by definition, $\lambda^t \mathbf{A}_i^{-1} = r\lambda^t$. Examining the component equations, we find that

$$\lambda_{l+1} = (r^l + \alpha_1 r^{l-1} + \cdots + \alpha_l)\lambda_1 \quad 1 \leq l \leq (m_N - 1). \qquad (63)$$

Also, for $1 \leq l \leq (-m_1)$,

$$\lambda_{m_N-m_1-l} = \frac{\lambda_{m_N-m_1-1}}{\beta_{-m_1}r^{l-1}}\left[\beta_{-m_1-l+1}r^{l-1} + \beta_{-m_1-l+2}r^{l-2} + \cdots + \beta_{-m_1}\right]. \qquad (64)$$

Finally,

$$\lambda_{m_N-m_1-1} = \frac{\beta_{-m_1}}{r}\lambda_1. \qquad (65)$$

Since the $\alpha$'s and $\beta$'s are positive quantities, statement (ii) of the lemma is true.

(iii) The statement may be verified by examining the characteristic polynomial $C(\mu)$ in eq. (60) and observing that the quantities $\psi^{(r)}(i)$ are monotonic, increasing with $i$.

## C.2 *Proof of Lemma 2*

It is required to prove that, for $j \geq i > 0$,

$$\lambda_t [\mathbf{A}_j^{-1} - \mathbf{A}_i^{-1}] \mathbf{P}_{j+1} \geq 0. \tag{66}$$

The matrices $\mathbf{A}_j^{-1}$ and $\mathbf{A}_i^{-1}$ are identical in all except the first row and also $\lambda^1 > 0$. Equation (66) is therefore equivalent to*

$$\mathbf{e}_1^t \mathbf{A}_j^{-1} \mathbf{P}_{j+1} \geq \mathbf{e}_1^t \mathbf{A}_i^{-1} \mathbf{P}_{j+1}. \tag{67}$$

We prefer to show that

$$\theta^{(N-1)}(i) p(j) \geq \theta^{(N-1)}(i) \mathbf{e}_1^t \mathbf{A}_i^{-1} \mathbf{P}_{j+1}, \tag{68}$$

where $\theta^{(N-1)}(i) \triangleq \{1 - \psi^{(N-1)}(i)\} > 0$. As $\mathbf{e}_1^t \mathbf{A}_j^{-1} \mathbf{P}_{j+1} = p(j)$, the lemma will then have been proved.

From eq. (24),

$$\theta^{(N-1)}(j) p(j) = p(j) - \psi^{(N-1)}(j) p(j)$$

$$= - \sum_{l=j+1}^{j+m_N-1} p(l) + \sum_{r=1}^{N-1} \sum_{l=j+m_N-m_r+1}^{j+m_N-m_r-1} \psi^{(r)}(l) p(l)$$

$$- \psi^{(N-1)}(j) p(j) \tag{69}$$

and

$$\theta^{(N-1)}(i) \mathbf{e}_1^t \mathbf{A}_i^{-1} \mathbf{P}_{j+1} = - \sum_{l=j+1}^{j+m_N-1} p(l)$$

$$+ \sum_{r=1}^{N-1} \sum_{l=j+m_N-m_r+1}^{j+m_N-m_r-1} \psi^{(r)}(l - j + i) p(l) - \psi^{(N-1)}(i) p(j). \tag{70}$$

Now

$$\theta^{(N-1)}(i) p(j) - \theta^{(N-1)}(i) \mathbf{e}_1^t \mathbf{A}_i^{-1} \mathbf{P}_{j+1}$$

$$\geq \theta^{(N-1)}(j) p(j) - \theta^{(N-1)}(i) \mathbf{e}_1^t \mathbf{A}_i^{-1} \mathbf{P}_{j+1}$$

$$\geq \sum_{r=1}^{N-1} \sum_{l=j+m_N-m_r+1}^{j+m_N-m_r-1} \{\psi^{(r)}(l) - \psi^{(r)}(l - j + i)\} p(l)$$

$$- \{\psi^{(N-1)}(j) - \psi^{(N-1)}(i)\} p(j) \geq 0, \tag{71}$$

because of the monotonicity of $\psi^{(r)}(l)$, and the final term in the expression on the right-hand side of (71) is cancelled by an identical component $(r = N - 1, l = j + m_N - m_{r+1})$ of the leading part. The lemma is proved.

---

* The column vector with the leading element equal to unity and all other elements equal to zero is denoted by $\mathbf{e}_1$.

## APPENDIX D

### Two Equations Concerning Mean First-Passage Times

We prove two assertions made in Section 3.1, eqs. (46) and (47), concerning (i) the homogeneous evolution of the vectors $\{z(n)\}$ via the matrix $D$ and (ii) the explicit formula for the mean first-passage time, $T(i)$.

#### D.1 Derivation of eq. (46)

Let $X(n)$ denote the event $1 \leqq \omega(\tau) \leqq L$ for all $\tau$, $0 \leqq \tau \leqq n$. Then, by definition,

$$z(j; n) = \Pr\left[\omega(n) = j \text{ and } X_n\right] \quad 1 \leqq j \leqq L.$$

Since it is also true that

$$z(j; n) = \Pr\left[\omega(n) = j \text{ and } X_{n-1}\right],$$

we have

$$z(j; n) = \sum_{i=1}^{L} \Pr\left[\omega(n) = j \,|\, \omega(n-1) = i, X_{n-1}\right] z(i; n-1).$$

We have obtained the quantities $\Pr\left[\omega(n) = j \,|\, \omega(n-1) = i, X_{n-1}\right]$ for $1 \leqq i, j \leqq L$ and, thereby, the following equations. In the following, $\mu$ denotes the number of contraction multipliers, that is,

$$m_1, m_2, \cdots, m_\mu < 0 < m_{\mu+1}, \cdots, m_N.$$

The basic recursion is, for $m_N + 1 \leqq j \leqq L + m_1$,

$$z(j; n) = \sum_{r=1}^{N} b^{(r)}(j - m_r) z(j - m_r; n-1). \tag{72}$$

The initial boundary equations are

$$z(j; n) = \sum_{r=1}^{\mu} b^{(r)}(j - m_r) z(j - m_r; n-1) \quad 1 \leqq j \leqq m_{\mu+1} \tag{73}$$

$$= \sum_{r=1}^{s} b^{(r)}(j - m_r) z(j - m_r; n-1) \quad m_s + 1 \leqq j \leqq m_{s+1}$$

$$s = \mu + 1, \mu + 2, \cdots, (N-1). \tag{74}$$

The final boundary equations are

$$z(j; n) = \sum_{r=s}^{N} b^{(r)}(j - m_r) z(j - m_r; n-1)$$

$$L + m_{s-1} + 1 \leqq j \leqq L + m_s, \ s = 2, 3, \cdots, \mu, \tag{75}$$

$$= \sum_{r=\mu+1}^{N} b^{(r)}(j - m_r)z(j - m_r; n - 1)$$

$$L + m_\mu + 1 \leqq j \leqq L - 1, \qquad (76)$$

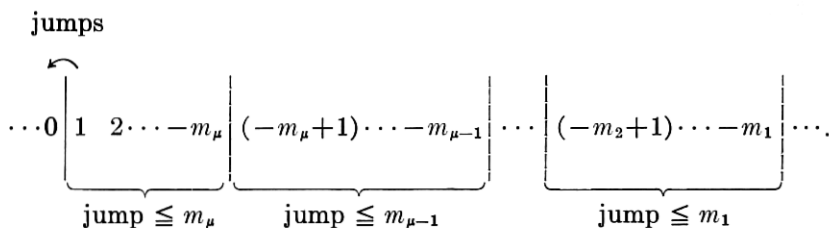$$= \sum_{r=\mu+1}^{N} \sum_{i=L-m_r}^{L} b^{(r)}(i)z(i; n - 1) \quad j = L. \qquad (77)$$

Equations (72) to (77) define the matrix $\mathbf{D}$ stated in the main text.

### D.2 Derivation of eq. (47)

For $i = 1, 2, \cdots, L$, let

$$f(i; n + 1) \triangleq \Pr [\text{first passage occurs at } (n + 1) | \omega(0) = i]$$

$$= \Pr [\omega(n + 1) \leq 0, X_n | \omega(0) = i]$$

$$= \sum_{j=1}^{-m_1} \Pr [\omega(n + 1) \leq 0 | \omega(n) = j]z(j; n), \qquad (78)$$

with $\mathbf{z}(0) = \mathbf{e}^{(i)}$, the vector with every element equal to zero except for the $i$th element, which is unity. The event $\omega(n + 1) = k \leq 0$ conditioned on $\omega(n) = j$ is associated with a jump $= k - j$. The following diagram illustrates the magnitudes of the jumps required for passage.

jumps



Equation (78) can be explicitly stated, thus,

$$f(i; n + 1) = \sum_{j=1}^{-m_\mu} \psi^{(\mu)}(j)z(j; n) + \sum_{j=-m_\mu+1}^{-m_{\mu-1}} \psi^{(\mu-1)}(j)z(j; n)$$

$$+ \cdots + \sum_{j=-m_2+1}^{-m_1} \psi^{(1)}(j)z(j; n). \qquad (79)$$

In the more convenient vector form,

$$f(i; n + 1) = \mathbf{c}^t\mathbf{z}(n), \qquad (80)$$

where the coefficients of the $L$-dimensional column vector $\mathbf{c}$ is obtained from (79), and we observe that only the leading $(-m_1)$ elements of $\mathbf{c}$ are nonzero.

The important fact about the vector **c** is that

$$\mathbf{c}^t = \mathbf{1}^t[\mathbf{I} - \mathbf{D}],\tag{81}$$

where **1** is the vector with every element equal to unity. Equation (81) may be established by either direct verification or by probabilistic reasoning. Now

$$
\begin{aligned}
T(i) &= \sum_{n \geqq 0} (n+1)f(i; n+1),\\
&= \mathbf{c}^t \sum_{n \geqq 0} n\mathbf{z}(n) + \sum_{n \geqq 0} f(i; n+1),\\
&= \mathbf{c}^t \sum_{n \geqq 0} n\mathbf{z}(n) + 1,\tag{82}\\
&= \mathbf{1}^t[\mathbf{I} - \mathbf{D}] \sum_{n \geqq 0} n\mathbf{z}(n) + 1 \quad \text{from (81),}\\
&= \mathbf{1}^t \sum_{n \geqq 0} \mathbf{z}(n),\tag{83}\\
&= \mathbf{1}^t\left[\sum_{n \geqq 0} \mathbf{D}^n\right]\mathbf{z}(0),\tag{84}\\
&= \mathbf{1}^t[\mathbf{I} - \mathbf{D}]^{-1}\mathbf{z}(0).\tag{85}
\end{aligned}
$$

Equation (82) is obtained by noting that the probability that passage occurs at finite time is unity. In obtaining Eq. (83), we have used $\mathbf{z}(n+1) = \mathbf{D}\mathbf{z}(n)$ and that $\mathbf{1}^t\mathbf{z}(0) = 1$. The convergence of the series $\Sigma\mathbf{D}^n$ is a consequence of the fact that every eigenvalue of the matrix **D** is strictly inside the unit circle. We omit the proof of this assertion, as it is similar to the proof given in Ref. 1 in connection with the matrix **D** for two-bit quantizers.

Equation (85) with $\mathbf{z}(0) = \mathbf{e}^{(i)}$ is the same as eq. (47) in the main text.

### REFERENCES

1. D. Mitra, "Mathematical Analysis of an Adaptive Quantizer," B.S.T.J., *53*, No. 5 (May-June 1974), pp. 867–898.
2. N. S. Jayant, "Adaptive Quantization with a One-Word Memory," B.S.T.J., *52*, No. 7 (September 1973), pp. 1119–1144.
3. P. Cummiskey, N. S. Jayant, and J. L. Flanagan, "Adaptive Quantization in Differential PCM Coding of Speech," B.S.T.J., *52*, No. 7 (September 1973), pp. 1105–1118.
4. D. J. Goodman and A. Gersho, "Theory of an Adaptive Quantizer," Proc. of December 1973 IEEE Symposium on Adaptive Processes, Decision and Control, pp. 361–365.
5. W. J. Dixon and A. M. Mood, "A Method for Obtaining an Analyzing Sensitivity Data," J. Amer. Statist. Assoc., *43*, 1948, pp. 109–126.
6. G. B. Wetherill, H. Chen, and R. B. Vasudeva, "Sequential Estimation of Quantal Response Curves: A New Method of Estimation," Biometrica, *53*, 1966, pp. 439–454.

7. Herman Chernoff, "Approaches in Sequential Design of Experiments," Technical Report, Stanford University, May 1973.
8. J. Max, "Quantization for Minimum Distortion," Trans. IRE, *IT-6*, March 1960, pp. 7–12.
9. J. C. Candy, "Limiting the Propagation of Errors in 1-Bit Differential Codecs," B.S.T.J., *53*, No. 8 (October 1974), pp. 1667–1676.
10. J. L. Doob, "Stochastic Processes," New York: John Wiley and Sons, 1953, pp. 300–301.
11. H. Kushner, "Stochastic Control," New York: Holt, Rinehart and Winston, 1971, pp. 188–224.
12. R. S. Bucy, "Stability and Positive Super-Martingales," J. on Differential Equations, *1*, 1965, pp. 151–155.