

Modulus Quantization for Matching Pursuit Video Coding

Ralph Neff and Avideh Zakhor
Electrical Engineering and Computer Science
University of California, Berkeley

Abstract

Overcomplete signal decomposition using matching pursuits has been shown to be an efficient technique for coding motion residual images in a hybrid video coder. Unlike orthogonal decomposition, matching pursuit uses an in-the-loop modulus quantizer which must be specified before coding begins. This complicates the quantizer design, since the optimal quantizer depends on the statistics of the matching pursuit coefficients which in turn depend on the in-loop quantizer actually used. In this paper, we address the modulus quantizer design issue, specifically developing frame-adaptive quantization schemes for the matching pursuit video coder. Adaptive dead-zone subtraction is shown to reduce the information content of the modulus source, and a uniform threshold quantizer is shown to be optimal for the resulting source. Practical 2-pass and 1-pass algorithms are developed to jointly determine the quantizer parameters and the number of coded basis functions in order to minimize coding distortion for a given rate. The compromise 1-pass scheme performs nearly as well as the full 2-pass algorithm, but with the same complexity as a fixed quantizer design. The adaptive schemes are shown to outperform the fixed quantizer used in earlier works, especially at high bit rates where the gain is as high as 1.7 dB.

I. INTRODUCTION

Most video compression systems in use today are built on a hybrid structure. Motion compensation is used to predict the current frame from previously reconstructed neighboring frames, and quantized transform coding is used for the prediction error image. Most motion residual coding work has concentrated on orthogonal transforms. Of these, the block discrete cosine transform (DCT) is the most popular and is widely used in video coding standards, e.g. [1] [2] [3]. Although the DCT is fast and efficient, it tends to introduce coding artifacts, especially visible block edges at low bit rates. Other orthogonal transforms such as wavelets have also been used [4], but these may show ringing artifacts around high contrast edges.

Recent work has demonstrated that better performance may be achieved through the use of overcomplete transforms. In an earlier work [5], we presented a motion residual coding system based on the iterative matching pursuit algorithm [6]. The scheme uses a greedy matching technique to decompose the residual image into atoms, which are coded basis functions from an overcomplete Gabor dictionary. Because the dictionary is large and varied, each coded atom closely matches the local signal to be coded. The structure of the basis set is thus not imposed on the reconstructed

image, and systematic artifacts such as block edges and ringing are reduced. The result is both a visual and a PSNR improvement over standard DCT-based video coders [7].

A flurry of recent work has been done to improve matching pursuit video coding. Al-Shaykh et.al. [8] increased the efficiency of the basic algorithm with a heuristic search procedure, and also introduced an SNR-scalable version. Others have focused on improving the basis set. De Vleeschouwer and Macq [9] achieved faster encoding through the use of a lower-cost subband dictionary. Redmill et.al. [10] showed both increased coding speed and a slight quality improvement by using a 2-stage filtering technique in which separably computable coefficients from one stage are combined spatially in a second stage to produce a dictionary containing oriented nonseparable bases.

Quantization of orthogonal transforms has the important property that the error incurred in the quantization of one basis coefficient is orthogonal to all of the other basis functions. Introducing quantization error for one coefficient thus has no effect on the computation of the other coefficients. For this reason, the quantizer for orthogonal transforms need not be known in advance, and may in fact be designed and applied based on the already-computed transform coefficients. The situation is much different for non-orthogonal matching pursuit decomposition. Since the quantization error in a single basis element is generally not orthogonal to the other basis functions, this quantization error must be considered when computing the remaining coefficients. Quantization is thus performed in-the-loop. This means that each matching pursuit coefficient or modulus is quantized as it is found, and the resulting quantized basis function is subtracted from the remaining signal energy before the next atom can be found. In this way, the quantization error from one matching pursuit stage may be progressively coded in future stages. The set of atoms chosen thus depends on the design of the quantizer. Likewise the optimal quantizer design depends on the modulus statistics from the chosen atoms. This results in a chicken-and-egg problem which is absent in coefficient quantizer design for orthogonal transforms.

For simplicity, previously published matching pursuit compression systems have used a quantizer in which both the form and the stepsize are fixed in advance heuristically. Since the rate-distortion tradeoff for the system is sensitive to both quantization and the number of coded atoms, by fixing the quantizer a degree of freedom for rate-distortion optimization is removed.

TABLE I
COMMON MPEG-4 TEST CONDITIONS

Seq	Name	Kbit/s	frame/s	size	Seq	Name	Kbit/s	frame/s	size
1	Container	10	7.5	QCIF	8	Foreman	48	10	QCIF
2	Hall	10	7.5	QCIF	9	News	48	7.5	CIF
3	Mother	10	7.5	QCIF	10	Coast	112	15	CIF
4	Container	24	10	QCIF	11	Foreman	112	15	CIF
5	Silent	24	10	QCIF	12	News	112	15	CIF
6	Mother	24	10	QCIF	13	Mobile	1024	30	SIF
7	Coast	48	10	QCIF	14	Stefan	1024	30	SIF

In this paper, we improve the efficiency of the matching pursuit video coding algorithm presented in [5] by applying optimal adaptive modulus quantization at the frame level. Section II reviews the basic system and Section III defines the fixed quantizer used for comparison. Section IV considers the problem of optimal scalar quantization for the system. An adaptive dead zone subtraction method is applied to the modulus values before quantization, and a uniform threshold quantizer (UTQ) with an adaptive step size is shown to be optimal for quantization of the resulting source. The problem is then reduced to finding the best combination of dead zone size DZ , quantization step size QP , and number of coded atoms N in a rate-distortion sense. A 2-pass algorithm is defined to solve this problem and is shown to have near optimal performance. Since the 2-pass algorithm is computationally expensive, a compromise 1-pass algorithm is developed in Section V. This algorithm is shown to perform nearly as well as the 2-pass scheme, but with about the same complexity as the simple fixed quantizer. A comparison of results is presented in Section VI, along with final conclusions.

II. MATCHING PURSUIT VIDEO COMPRESSION

A. Preliminaries

Throughout this work, integer indices k and i are used to represent the current frame and the current atom stage, respectively. Vectors such as \vec{v} are generally defined on a Hilbert space, but may be considered for our purpose to exist in 2-D image space. For coding results, we adopt the sequences and test conditions which were used in the development of the MPEG-4 standard. These are summarized in Table I.

B. Matching Pursuit Theory

The Matching Pursuit algorithm of Mallat and Zhang [6] expands a signal \vec{f} using an overcomplete dictionary of functions \mathcal{G} . Since \vec{f} will be coded over multiple atom stages, define \vec{f}_i to be the remaining signal energy after i atoms have been coded, with $\vec{f}_0 = \vec{f}$. Individual dictionary

functions $\vec{b}_\gamma \in \mathcal{G}$ are assumed to have unit norm. Here γ is an indexing parameter associated with a particular dictionary function. Each single stage $i \in [1 \dots N]$ of the decomposition requires choosing γ to maximize the absolute value of the following inner product:

$$p = \langle \vec{f}_{i-1}, \vec{b}_\gamma \rangle \quad (1)$$

Denote the index chosen at stage i as γ_i and the associated inner product as the modulus value p_i . These two parameters define a coded basis function or *atom*, and the remaining signal energy for the next stage is computed as:

$$\vec{f}_i = \vec{f}_{i-1} - p_i \vec{b}_{\gamma_i} \quad (2)$$

The procedure continues iteratively until some stopping criteria is reached. For example, the bit rate available to encode \vec{f} may be exhausted, or some energy threshold for the residual may be reached. Each stage i yields a dictionary index γ_i , a modulus p_i , and a residual signal \vec{f}_i which is passed on to the next stage. After N atoms are coded, the signal can be approximated by a linear function of the dictionary elements:

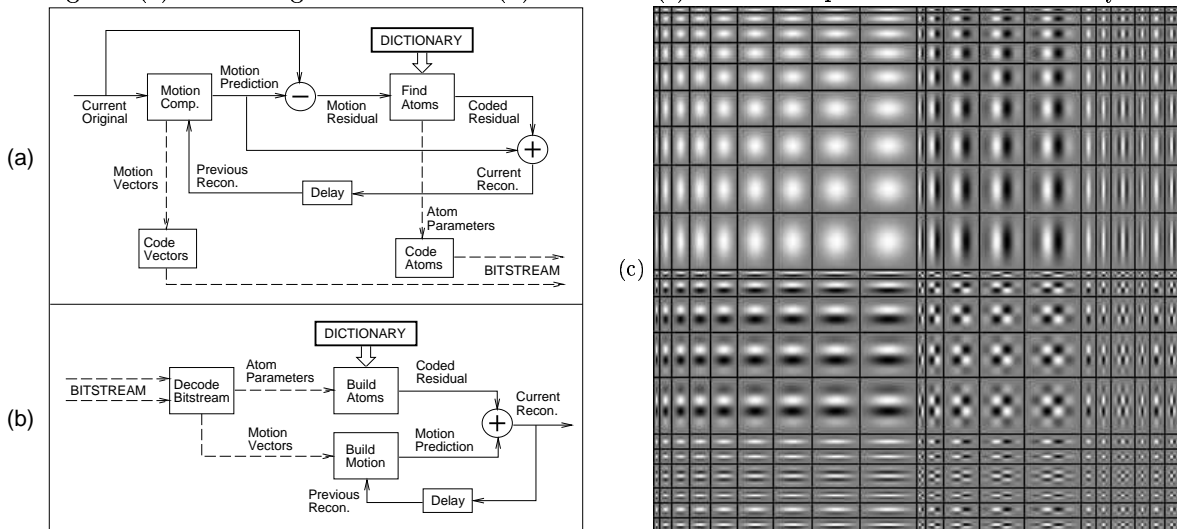
$$\vec{f} \approx \sum_{i=1}^N p_i \vec{b}_{\gamma_i} \quad (3)$$

To use this technique in a coding system requires modulus p_i to be quantized before the subtraction in Equation (2). Quantization must be done in the coding loop, since the residual signal \vec{f}_i is used to find the remaining atoms. A description of the resulting video coding system is given in Section II-C, and the fixed quantizer design used in previous works is reviewed in Section III.

C. System Description

The basic matching pursuit video coding system was defined in [5] and is shown in block form in Figure 1. Figure 1a shows the hybrid encoder in which motion compensation is used to predict the current frame from the previous reconstructed frame. The prediction image is subtracted from the current original image to get the motion residual image. In the notation of the previous section, the motion residual image is \vec{f} , the signal to which matching pursuit coding will be applied. This is done in the ‘‘Find Atoms’’ block of the figure, which decomposes \vec{f} into N atoms which are then coded and transmitted to the decoder. The bitstream consists mainly of the motion vectors and atom parameters, with a small amount of header information sent per frame. The decoder,

Fig. 1. (a) Block diagram of encoder. (b) Decoder. (c) The 2-D separable Gabor dictionary.



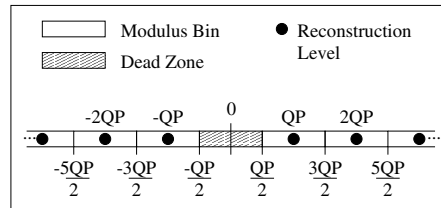
as shown in Figure 1b, recovers the motion vectors and atom parameters from the bitstream and uses them to reconstruct the motion prediction image and the atom-coded residual signal. These are then combined to produce the reconstructed image.

The basis set \mathcal{G} consists of the separable collection of 400 2-D Gabor functions shown in Figure 1c. An exact mathematical specification is given in [5]. An atom consists of a single basis shape from the figure placed at some spatial location (x, y) in the coded residual image. For transmission, the atoms are sorted into a position scanning order, and the positions are differentially coded along the scanning path. The differential position codes and the basis shape indices are entropy coded and transmitted without loss to the decoder. The atom modulus values p_i are also transmitted; these have been quantized in the coding loop, and the appropriate quantizer bin indices are entropy coded and transmitted.

DCT and wavelet based systems typically use quantization based rate control schemes. The quantizer stepsize is controlled in an attempt to match a target rate or other coding requirement. Previously published matching pursuit systems have instead fixed the quantizer and used the number of coded atoms N to vary the rate. Atoms are coded until a target bit rate R is exhausted or some quality constraint is reached. Very fine-grain control is possible, with the actual bit rate typically within 25 bits per frame of the intended target rate, or the visual quality within 0.1 dB of a specified target PSNR.

For a matching pursuit system, modulus quantization is only one of two sources of coding distor-

Fig. 2. Simple midtread linear symmetric quantizer with fixed stepsize QP , as used in previously published matching pursuit systems.



tion. The other arises because only a finite number of atoms may be sent, and the approximation in (3) only approaches equality when the number of atoms is very large. It is possible to trade off distortion between these two sources. For example, if more bits from a limited budget are assigned to represent the moduli, then quantization error is reduced, but the additional cost per atom reduces the number N that may be coded at the same target bit rate. Choosing the combination that minimizes the overall coding distortion is the central theme of this paper, and will be addressed in Section IV.

III. FIXED QUANTIZER DESIGN

Previously published matching pursuit video coders used a midtread linear symmetric quantizer with midpoint reconstruction, as shown in Figure 2. For the results shown in [5], [8] the quantization stepsize was fixed at $QP = 30$, a value heuristically chosen to perform well on the MPEG-4 test sequence set.

Since the focus of this paper is adaptive quantizer design, we now consider this fixed quantizer for use as a comparison base for our later adaptive designs. The advantages and disadvantages of this design may be summarized as follows:

Advantages:

- A1.* Fixed quantizer requires no knowledge or pre-analysis of the source.
- A2.* Provides for simple and effective rate control based on number of coded atoms N .
- A3.* Performs well for standard MPEG-4 test sequences and rates.

Disadvantages:

- D1.* Fixed quantizer removes a degree of optimization freedom. That is, for a given frame, a different quantizer with a different number of coded atoms may produce a lower distortion at the same target bit rate.
- D2.* Large fixed dead zone imposes a hard limit on attainable bit rate and image quality.

Of the two listed disadvantages, note that *D1* is intrinsic to any fixed quantizer design. That is, no fixed quantizer will be suitable for all frames at all possible bit rates. It is thus appropriate that

this disadvantage be present in our fixed quantizer comparison model. However, $D2$ is really a flaw in the specific quantizer design shown in Figure 2. This flaw will be explained and characterized in the remainder of this section. Sections III-A and III-B will then develop a modified fixed quantizer design which eliminates this flaw while preserving the listed advantages.

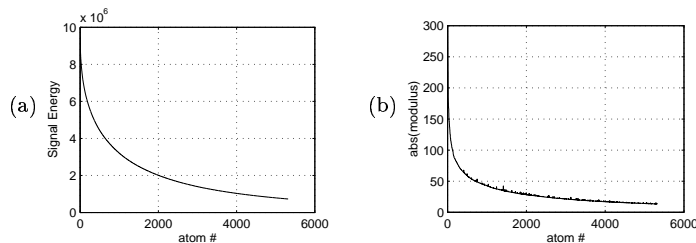
To see why $D2$ is true, consider the matching pursuit decomposition of a single frame of the Stefan sequence as shown in Figure 3. The sequence was coded at 30 fps at a constant bit rate of 3 Mbit/s. This particular frame reached the desired rate using 5318 atom stages. Figure 3a shows the remaining signal energy $\|\vec{f}_i\|^2$ after each atom iteration i . Figure 3b shows the absolute modulus values $|p_i|$ before quantization. Note that the signal energy is strictly decreasing, reflecting that each atom results in some positive energy reduction. The modulus values generally decrease with iteration, although not monotonically. If a given frame is coded to an arbitrarily high bit rate using the fixed quantizer of Figure 2, at some stage i the modulus value p_i will satisfy $|p_i| \leq (QP/2)$, and thus will fall into the dead zone of the quantizer. Since the resulting atom quantizes to zero, the matching pursuit update equation (2) gives us $\vec{f}_i = \vec{f}_{i-1}$. Because the residual energy was unaffected by stage i , the ‘‘Find Atom’’ block shown in Figure 1 will find the exact same atom for stage $i + 1$, and the algorithm will become deadlocked, unable to code additional nonzero coefficients. This effectively imposes a hard limit on both rate and quality, and makes the specific quantizer design of Figure 2 unusable at high bit rates. Specifically, it restricts the matching pursuit decomposition to operate only at rates low enough such that all moduli p_i satisfy:

$$|p_i| > QP/2 \quad (4)$$

Our experience has shown this condition to be violated for SIF sequences such as Mobile and Stefan when coded in the 1 to 2 Mbit/s range. It may also be violated at much lower rates for particular frames with low residual energy, which occur during low motion periods of the sequence. The absolute modulus values for such frames tend to be much lower than frames with higher levels of motion.

Restriction (4) is specific to the fixed quantizer design of Figure 2. In the next section, we develop the comparable restriction for a more general fixed quantizer design. A modified design which may be used at higher bit rates is then developed in Section III-B. This modified design will then be used as the comparison base for the adaptive quantizers developed in Section IV.

Fig. 3. Matching pursuit decomposition of first residual frame of the Stefan sequence. (a) Remaining signal energy vs. iteration. (b) Absolute modulus vs. iteration.



A. Restriction for a general fixed quantizer

Consider a general quantizer $Q(\cdot)$, and allow the quantization error to be represented by $d_i = p_i - Q(p_i)$. Suppose \vec{f}_i is the signal remaining after the atom for iteration i has been quantized and subtracted from the source energy. Then the Energy Reduction (ER) for stage i can be defined as:

$$ER_i \triangleq \|\vec{f}_{i-1}\|^2 - \|\vec{f}_i\|^2$$

As shown in Appendix A, the single-stage Energy Reduction can be represented as:

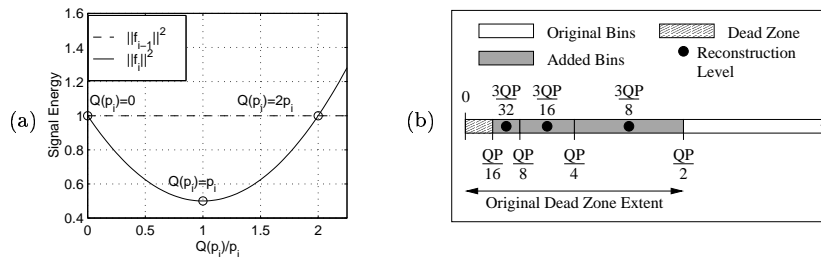
$$ER_i = p_i^2 - d_i^2 \quad (5)$$

To avoid the deadlocking problem seen in the earlier fixed quantizer design, we require positive energy reduction, specifically that $ER_i > 0$, for each stage i . For an even-symmetric quantizer design, consider without loss of generality the $p_i > 0$ case. By combining Equation (5) with the definition for d_i , we arrive at the following restriction in order to avoid deadlocking:

$$0 < Q(p_i) < 2p_i \quad (6)$$

A graphical interpretation of this restriction is shown in Figure 4a. The figure shows the remaining signal energy before and after stage i as a function of $Q(p_i)/p_i$. The dotted line represents $\|\vec{f}_{i-1}\|^2$, the energy before the atom is coded, and the solid line shows $\|\vec{f}_i\|^2$, the energy after the atom is coded. The plot is a direct consequence of Equation (5). The maximum energy reduction of p_i^2 occurs in the absence of quantization error, that is when $Q(p_i)/p_i = 1$. When this ratio falls outside the range $(0, 2)$, the coded atom has a destructive effect in that the remaining signal energy is increased. By designing a modified fixed quantizer to operate within this restricted range (6) for a wide range of interesting bit rates, we can avoid the deadlocking problem seen earlier.

Fig. 4. (a) Remaining signal energy before and after the current stage as a function of $Q(p_i)/p_i$. Positive energy reduction requires $Q(p_i) \in [0, 2p_i]$. (b) Modification of the fixed quantizer design. Original dead zone has been divided into several additional quantization bins.



B. Modified fixed quantizer design

Here we show a heuristically modified fixed quantizer design which allows high bit rate operation without changing the earlier documented low bit rate performance. Recall that the main disadvantage $D2$ of the original fixed quantizer design was a large fixed dead zone which resulted in the deadlocking problem at high rates. We modify this design by splitting the dead zone into additional bins to be used at high bit rates when the modulus values become small. The resulting nonlinear quantizer is shown in Figure 4b. Midpoint reconstruction is used, and even symmetry allows us to show only the positive half. Each new bin is formed by splitting the remaining dead zone in half. Reducing the dead zone in this way allows more non-zero atoms to be coded, and thus enables the encoder to work at higher bit rates without deadlocking. Any number of such bins may be added in order to achieve arbitrarily high rates. We use three added bins as shown in the figure, reducing the dead zone to $[0, QP/16]$. Our experiments show that this is sufficient to allow operation without deadlocking for the SIF MPEG-4 test sequences at 3 Mbit/s, the highest rate tested in this paper. Note however that this strategy merely postpones the deadlocking problem. That is, at some rate beyond 3 Mbit/s, the condition $p_i > QP/16$ will be violated, and deadlocking will again occur.

It is important to note that the added bins are consistent with the general quantizer restriction (6), and so positive energy reduction, that is $ER_i > 0$, is guaranteed for all moduli which fall into these bins. At low bit rates where moduli are larger than $QP/2$, the quantizer is identical to the earlier design, and so low bit rate performance is not affected. This is true for the standard MPEG-4 sequences and rates, for which the earlier restriction (4) holds.

A remaining question is how to entropy code the new high rate bins. The simplest approach

TABLE II
 VARIABLE LENGTH CODE (VLC) MAPPINGS FOR NEWVLC AND SHIFTVLC CODING SCHEMES.

NewVLC Scheme			ShiftVLC Scheme				
Index	Recon	Bin Type	Index	$shift = 0$	$shift = 1$	$shift = 2$	$shift = 3$
0	$3 \cdot QP/32$	Added	0	QP	$3 \cdot QP/8$	$3 \cdot QP/16$	$3 \cdot QP/32$
1	$3 \cdot QP/16$	Added	1	$2 \cdot QP$	QP	$3 \cdot QP/8$	$3 \cdot QP/16$
2	$3 \cdot QP/8$	Added	2	$3 \cdot QP$	$2 \cdot QP$	QP	$3 \cdot QP/8$
3	QP	Original	3	$4 \cdot QP$	$3 \cdot QP$	$2 \cdot QP$	QP
4	$2 \cdot QP$	Original	4	$5 \cdot QP$	$4 \cdot QP$	$3 \cdot QP$	$2 \cdot QP$
5	:	:	5	:	:	:	:

is to add a new VLC code for each new bin and retrain the system as before. This method increases the size of the modulus VLC table, and will be denoted “NewVLC”. We also present an alternate approach from our MPEG standards work [11]. Suppose the size of the modulus VLC table remains unchanged and instead a single *shift* codeword is sent at the frame level to specify how many of the new bins were used during coding. For three added bins, *shift* takes a value in $\{0, 1, 2, 3\}$, for an overhead of 2 bits per frame. The VLC table is then interpreted so that the first code in the table points to the lowest bin actually used in coding the current frame. The VLC table is effectively shifted down by *shift* bins. This idea is based on the empirical observation that the modulus distribution is effectively a shifted exponential distribution, and so the lowest bin used typically has a high probability of occurrence. We defer a detailed look at the modulus distribution until Section IV-A.

To clarify the difference between the two schemes, the VLC mappings are shown as Table II. Note that $shift = 0$ at low bit rates; in this case both the quantizer and VLC structure match those of the original fixed quantizer design shown in Figure 2. The 2-bit per frame *shift* codeword is insignificant, and so low bit rate performance is equivalent to earlier published results [5], [7].

Since “NewVLC” and “ShiftVLC” are based on the same quantizer design, they differ only by the resulting cost per atom. To compare them on this basis, the two schemes were implemented and separately trained on a set of sequences outside the MPEG-4 test set. The QCIF sequences from the MPEG-4 set were then run for a comparison, with the resulting cost per modulus summarized in the top section of Table III. The first three columns show that at these rates the two schemes produce a nearly identical cost per modulus. The remaining columns show that the ShiftVLC system uses a slightly lower percentage of its atom budget on modulus codewords, but the difference is very small. The two schemes are thus nearly identical at low bit rates where the additional “high-rate” bins aren’t used. The difference is much more apparent at higher rates, as shown in the bottom

TABLE III
AVERAGE MODULUS VLC COST AND MODULUS BITS AS A PERCENTAGE OF TOTAL ATOM BITS.

Sequence	Bits/Modulus			% Bits to Modulus		
	NewVLC	ShiftVLC	Diff	NewVLC	ShiftVLC	Diff
Container10	3.30	3.29	-0.01	14.97	14.90	-0.07
Hall10	4.36	4.36	0.00	18.79	18.78	-0.01
Mother10	3.10	3.11	0.01	12.26	12.21	-0.05
Container24	2.38	2.38	0.00	12.61	12.65	0.04
Silent24	3.62	3.63	0.01	15.73	15.71	-0.02
Mother24	2.34	2.34	-0.01	11.21	11.13	-0.07
Coast48	3.44	3.44	0.00	16.26	16.23	-0.03
Foreman48	3.25	3.25	0.00	15.38	15.32	-0.06
Mother24	2.34	2.34	-0.01	11.21	11.13	-0.07
Mother48	2.12	2.14	0.02	11.34	11.42	0.08
Mother72	2.93	2.29	-0.64	15.65	12.70	-2.95
Mother96	4.18	2.60	-1.58	21.37	14.62	-6.75

section of Table III. Here the QCIF Mother sequence is coded with bit rates from 24 to 96 kbit/s. At the highest rates, ShiftVLC shows a performance advantage over the simpler NewVLC system, saving up to 1.6 bits per modulus code and using a smaller percentage of atom bits for modulus codewords. Because of this advantage, ShiftVLC will be used for comparison to the adaptive systems developed in the next sections.

IV. OPTIMAL ADAPTIVE QUANTIZATION

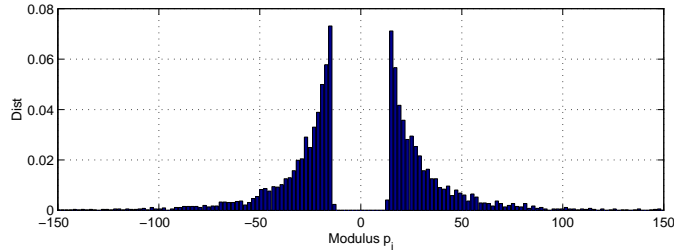
This section considers the problem of designing an optimal adaptive modulus quantizer for the matching pursuit video encoder. Section IV-A examines the modulus distribution and shows that an adaptive dead zone subtraction technique can be used to reduce the energy content before quantization. The optimal form and stepsize of the quantizer for the resulting source will be developed in Sections IV-B and IV-C, respectively, and a practical 2-pass algorithm will be shown to have near-optimal performance.

A. Adaptive Dead Zone Subtraction

A sample modulus distribution for the matching pursuit video coding system is shown in Figure 5. This is a normalized histogram of the moduli for a single frame of the Stefan sequence coded to 3 Mbit/s. In-loop quantization¹ is used in coding, but the values before quantization are used to produce the continuous distribution shown. The figure illustrates some properties that generally hold for the modulus distribution. Because the distribution is symmetric about zero, the quantized moduli may be coded as magnitude plus a sign bit with no loss of efficiency. This reduces the

¹ Figure 5 uses ShiftVLC as the in-loop quantizer. This choice is arbitrary, since the distribution changes very little if a different quantizer is used, or even no quantizer at all. This is further investigated in Figure 9.

Fig. 5. Modulus distribution for a single frame of Stefan (SIF) coded at about 3 Mbit/s.



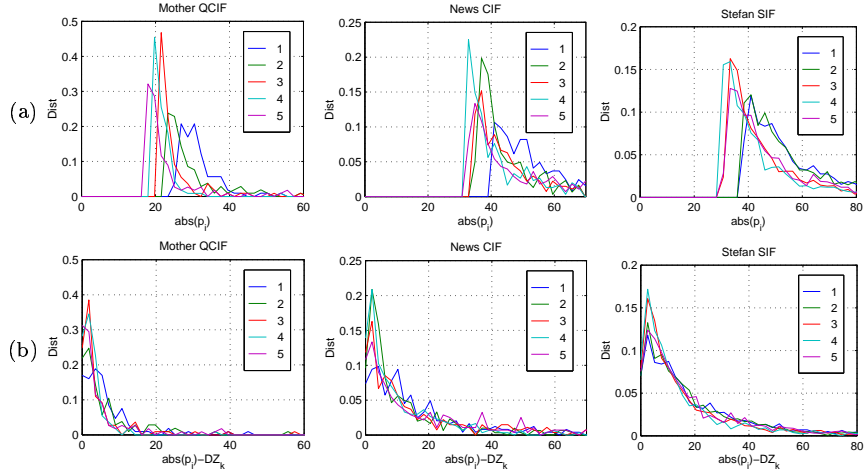
number of quantizer bins and VLC codes by half, and has an added benefit for adaptive quantizer design, since the data per bin available for run-time analysis is effectively doubled. A second interesting property is the presence of a large dead zone. Note that the cutoff around the dead zone is very sharp, with high probabilities near the dead zone decision border.

We now show how the modulus distribution changes from frame to frame. Figure 6a shows absolute modulus distributions for coded frames of three different MPEG-4 test sequences. For each sequence, five frames are coded sequentially without quantization and the resulting per-frame modulus distributions are superimposed. Since the in-loop quantizer is omitted, exact bit rates cannot be computed. Instead, the number of atoms per frame is matched to the corresponding standard MPEG-4 test run. The resulting plots show a poor match between individual distribution curves, which suggests the need for a frame adaptive quantizer design. Close examination of the plots shows that the shape of the distribution is fairly consistent, but that this shape shifts horizontally from frame to frame. This is seen most clearly at high bit rates such as Stefan, where a large number of atoms per frame allows more reliable computation of the distribution curves.

Because the distribution shifts between frames, we propose a simple transformation to the original modulus source. Define DZ_k as the minimum absolute modulus value found in encoding frame k . A “DZ-Subtraction” operation is then performed on the modulus source, that is DZ_k is subtracted from all absolute modulus values within the frame. The distributions of the resulting DZ-Subtracted source are more consistent from frame to frame, as illustrated in Figure 6b.

The effect is even more dramatic when modulus distributions are combined over multiple sequences or across multiple bit rates. Figure 7 shows the combination of multiple sequences at the same rate. Each plot is generated by encoding all of the MPEG-4 test sequences from Table I that correspond to the given bit rate categories (24 kbit/s, 112 kbit/s, and 1 Mbit/s). The sequences from each category were coded without quantization using the same conditions as in the

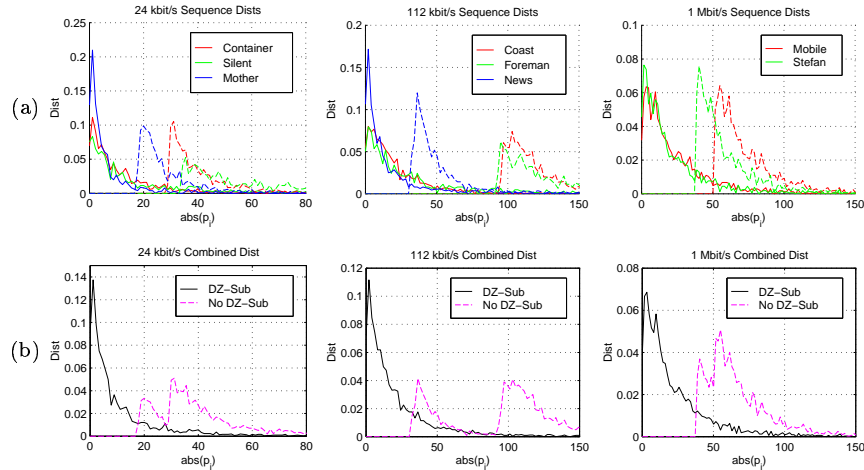
Fig. 6. Absolute modulus distributions for single matching pursuit residual frames coded sequentially without quantization. (a) Sequences coded without DZ-subtraction. (b) Same sequences with framewise DZ-subtraction. Note the number of atoms per frame match the standard MPEG-4 sequence runs of Mother at 24 kbit/s, News at 112 kbit/s, and Stefan at 1 Mbit/s.



previous figure. Modulus data was then taken from enough frames to produce about 3000 atoms per bit rate category. For example, the leftmost plots combine the 24 kbit/s QCIF encodings of Container, Silent and Mother with 10 frames used per sequence. Each trace in Figure 7(a) shows the modulus distribution combined across the frames of a single sequence. Note that the dotted lines, which represent the distributions without DZ-subtraction, match poorly across sequences. A much better match is seen in the solid lines, which show the distributions after DZ-subtraction. Figure 7(b) shows the combined distributions for each bit rate category. In each of these plots, a sharp, exponential-like distribution is produced by the DZ-subtraction method as shown by the solid line. This shape is well-matched to the single-frame distributions of Figure 6b, and is even consistent across bit rates. The distributions without DZ-subtraction are much less consistent, because the modulus source varies greatly across frames, sequences, and bit rates.

Figure 8 shows the effect of combining distributions across bit rate categories. To generate the figure, all of the standard MPEG-4 test sequences were coded in the same manner as used in the previous two figures, and enough frames were taken from each of the five bit rate categories (10, 24, 48, 112, and 1024 Kbit/s) to produce about 3000 atoms from each category. This was done to prevent the higher bit rate sequences, having more coded atoms per frame, from dominating over the lower bit rate data. Figure 8a shows a simple combination of the absolute moduli, while Figure 8b shows the DZ-subtracted result. Without the DZ-subtraction, the combined distribution is somewhat random and matches neither the individual frame distributions from Figure 6a nor

Fig. 7. Combined modulus distributions for each of three MPEG-4 test rate categories. (a) Separate distribution curves for each sequence in the given bit rate category. Solid lines show runs with DZ-subtraction. Dotted lines show runs without DZ-subtraction. (b) Combined distribution curve for all sequences in the given rate category.



the combined distributions for the various rate categories shown in Figure 7. With DZ-subtraction, the final combined distribution is again a sharp exponential, which well matches the corresponding DZ-subtracted distributions in Figures 6b and 7. To verify the exponential form of the distribution, a least squares best fit exponential curve is also shown in Figure 8b.

We compare the information content of the two distributions of Figure 8 by computing the Shannon differential entropy, $h(X)$, for each distribution.² Simple numerical integration [13] gives a value $h(X) = 6.914$ for the combined distribution in Figure 8a, and a value $h(X) = 5.682$ for the DZ-subtracted distribution in Figure 8b. Suppose the modulus source is to be quantized using the optimal entropy constrained scalar quantizer as defined in [14]. Dead-zone subtraction will then improve the rate-distortion performance $r(d)$ of the quantizer by the difference of 1.279 bits per symbol, a savings of about 18 percent. The dead zone parameter DZ_k must be transmitted once per frame, but this is a negligible cost even at low bit rates.

The above experiments were done in the absence of quantization. We now show that the addition of in-loop quantization does not significantly change the observed distribution. Figure 9 shows single-frame modulus distributions for two sample frames coded with and without in-loop

² Shannon differential entropy reflects the information content of a continuous source, similar to entropy for discrete valued sources [12]. It is defined in terms of the probability density function $p_x(x)$ as:

$$h(X) = \int_{-\infty}^{\infty} p_x(x) \log_2 p_x(x) dx$$

Fig. 8. Modulus distributions combined across all MPEG-4 sequences and rates. Enough coded frames were taken to provide approximately 3000 atoms from each of the five bit rate categories (10, 24, 48, 112, and 1024 Kbit/s). (a) Without dead-zone subtraction. (b) With frame-wise dead zone subtraction. Dotted line shows best exponential fit.

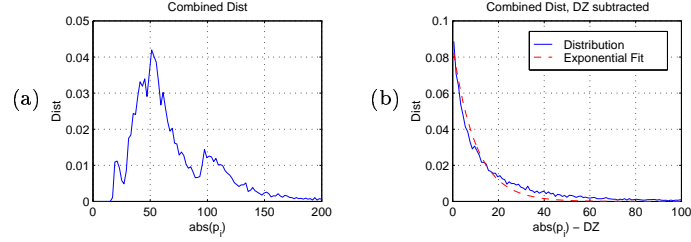
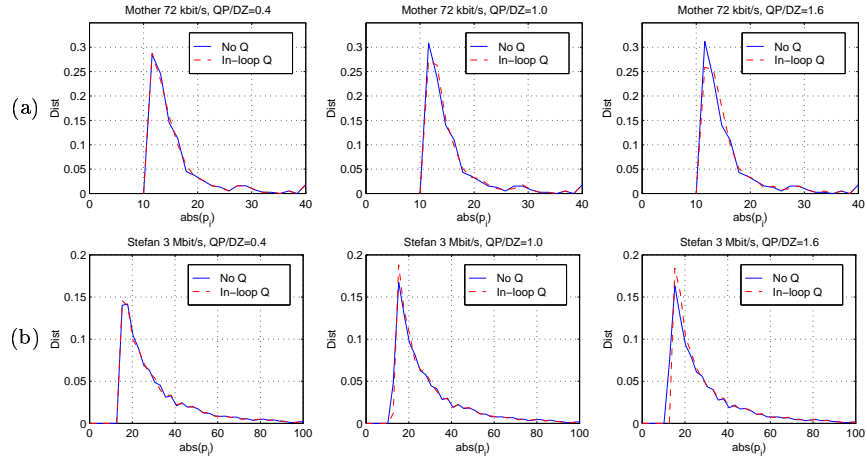


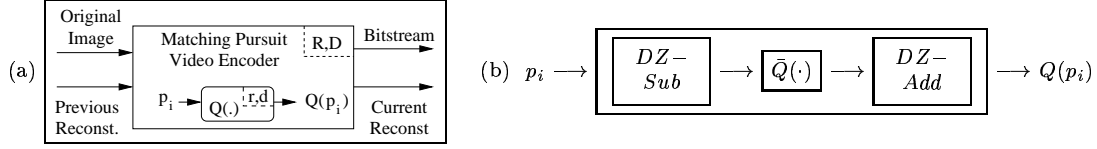
Fig. 9. The effect of in-loop quantization on the absolute modulus distribution for single coded frames. (a) Mother QCIF at 72 kbit/s. (b) Stefan SIF at 3 Mbit/s. Three different quantizers with progressively larger stepsize to deadzone ratios are used for each frame.



quantization. The top row of plots, (a), shows a frame from Mother QCIF at 24 kbit/s, while the bottom row, (b), shows a frame of Stefan SIF at 3 Mbit/s. For each plot, the solid line depicts matching pursuit without quantization, and the dotted line shows in-loop quantization. Clearly the amount of error introduced by the quantizer will determine the amount of change seen in the resulting distribution. For this reason, three quantizers are used with varying stepsizes. The quantizer is uniform with an expanded dead zone, which is equivalent to DZ-subtraction followed by a uniform threshold quantizer (UTQ). The dead zone size DZ is set to be the minimum absolute modulus value seen in the non-quantized encoding. The step size QP is then set so that the ratio QP/DZ takes on values of 0.4, 1.0 and 1.6. Note that this stepsize to dead-zone ratio provides a convenient description of the quantizer, and so will be used in the remainder of this paper.

Scanning the plots of Figure 9 from left to right, it is clear that for ratio values of 1.0 and below, in-loop quantization produces no noticeable change to the absolute modulus distribution. For more coarse quantization with $QP/DZ = 1.6$, some change is visible, but the distributions

Fig. 10. (a) Abstraction of the encoder, showing relation between overall rate and distortion (R,D) and that of the internal modulus quantizer (r,d). (b) Definition of quantizer $Q(\cdot)$, which combines DZ-subtraction and quantization.



with and without quantization are still nearly identical. We can assume that within this range of ratio values, in-loop quantization has little effect on the modulus distribution or the effective dead zone size. Under this assumption, the skeleton of a 2-pass adaptive quantization algorithm can be seen. A first pass pre-codes the residual signal without quantization to estimate the number of atoms and the size of the dead zone DZ_k . The second pass uses dead-zone subtraction with in-loop quantization of the resulting moduli. The details of this algorithm, along with a solution for the optimal quantizer form and stepsize, will be developed in the next two sections.

B. Optimal Quantizer for the Matching Pursuit System

In this section, we develop the optimal form of the quantizer to use in conjunction with the DZ-subtraction process. To clarify the problem, we first introduce some notation. Consider the abstracted view of the encoder shown in Figure 10a. The large box represents the entire system, which takes the current original and previous reconstructed images as input, performs both motion compensation and matching pursuit based residual coding, and outputs the coded bitstream and the current reconstruction to be used in encoding the next stage. Overall distortion D is defined as the mean square error in the coded image. Overall rate R is the number of bits used to code the current frame, which is the sum of rates to header R_{hdr} , motion vectors R_{mot} , and atoms R_{atm} . The atom budget may be further divided into bits used for position coding R_P , basis indices R_I , and quantized modulus R_Q . In summary,

$$R = R_{hdr} + R_{mot} + R_{atm} \quad (7)$$

$$= R_{hdr} + R_{mot} + (R_P + R_I + R_Q) \quad (8)$$

The modulus quantizer for the system is $Q(\cdot)$, shown as the smaller block in Figure 10a. The quantization process also has an associated rate and distortion. Quantizer rate r is defined as the

average number of bits spent to encode each modulus value, and quantizer distortion is:

$$d = \text{mean}((p_i - Q(p_i))^2) = \text{mean}(d_i^2)$$

Consider a 2-pass encoding algorithm in which a first pass without quantization provides an estimate of the dead-zone parameter DZ_k . The second pass then uses in-loop quantizer $Q(\cdot)$, which is defined to perform DZ-subtraction, apply some quantizer $\bar{Q}(\cdot)$ to the resulting moduli, and then add DZ_k to the result to obtain the quantized value. This is illustrated in block form as Figure 10b. A magnitude VLC code and separate sign-bit are transmitted for each modulus. The floating point representation of DZ_k is sent in the frame header to allow quantizer tracking at the decoder.

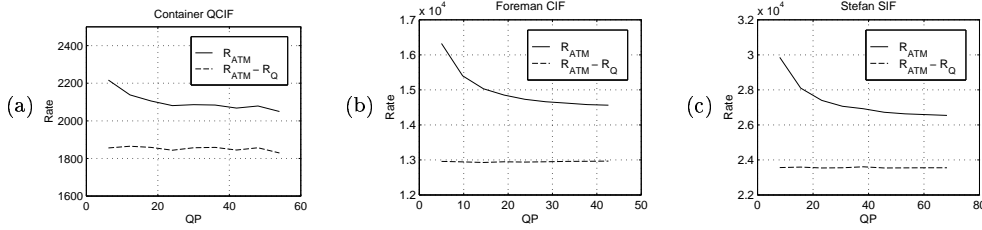
The specific problem is to design quantizer $\bar{Q}(\cdot)$ such that the overall distortion D for the matching pursuit encoder is minimized at a given overall rate R . Essentially the design of $\bar{Q}(\cdot)$ allows a tradeoff between modulus precision and number of coded atoms. A quantizer may code modulus values more precisely at the cost of increasing R_Q , however this decreases the total number of atoms which may be coded at an overall target rate R . We thus seek an optimal balance between quantizer precision and number of coded atoms in order to minimize distortion for a given rate.

In the previous section, we showed that the DZ-subtracted modulus source to which $\bar{Q}(\cdot)$ is to be applied has an exponential distribution. The optimal entropy constrained scalar quantizer (ECSQ) for this distribution is a UTQ [12], [15]. That is to say, using a UTQ for $\bar{Q}(\cdot)$ will provide the best rate-distortion performance for the DZ-subtracted moduli, and will in turn minimize d for any given quantizer rate r . Through the use of some simplifying assumptions, we will show in the remainder of this section that the UTQ is optimal in terms of overall (R, D) for the matching pursuit system. We begin by presenting a pair of assumptions to relate (r, d) for the quantizer to (R, D) for the overall system. These are presented without proof, but empirical evidence is provided to show that each holds for the matching pursuit video system.

Assumption 1: For a given number of atoms N , R_Q may vary as the quantizer changes. The other bit rate terms and their sum $R - R_Q = (R_{hdr} + R_{mot} + R_P + R_I)$ will remain roughly constant with changes in the quantizer.

The first part of the assumption is obvious, since $R_Q = Nr$, and r is certainly affected by changes in the design of $\bar{Q}(\cdot)$. As for the other bit rate terms, R_{hdr} and R_{mot} are determined before atom

Fig. 11. Bit rates used to code a set number of atoms as quantizer stepsize QP changes. Total rate R_{atm} changes with QP , while the rate excluding quantization $R_{atm} - R_Q$ stays approximately constant. Sample coded frames are from: (a) Container QCIF coded to 110 atoms, (b) Foreman CIF coded to 735 atoms, and (c) Stefan SIF coded to 1419 atoms.



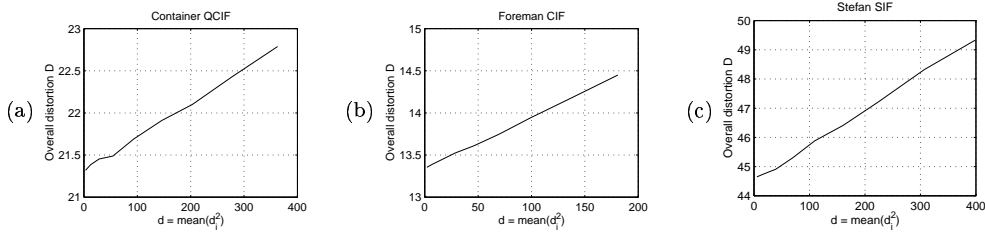
coding begins and thus are unaffected by the quantizer design. R_P depends on the spatial density of atoms, and thus should be approximately constant for a constant N . R_I depends on the set of basis shapes coded during the matching pursuit process. While this set does depend on the in-loop quantizer design, we expect the average basis index codeword length (and thus R_I) to be approximately constant if N is sufficiently large.

We are mainly interested in the last part of the assumption, that $R - R_Q$ is roughly constant with changes in the quantizer. Empirical evidence is provided in Figure 11. Three different sample frames with a wide range of coding complexities and effective bit rates are shown. Each frame is coded to a constant number of atoms using the 2-pass DZ-subtracted quantizer described earlier. For $\bar{Q}(\cdot)$, we use a UTQ with a variable stepsize QP . The total bit rate to atoms R_{atm} and the rate without quantization $R_{atm} - R_Q$ are plotted against stepsize QP . The plots show an increase in the total rate R_{atm} for the finer stepsizes, but the rate with R_Q subtracted is roughly constant as the quantizer changes. This implies that the change in rate with QP is entirely within the R_Q term, and that the sum of the remaining terms $R_P + R_I$ is approximately constant. This matches the expectations developed in the preceding paragraph, and is consistent with Assumption 1.

Assumption 2: Suppose a frame is coded using an initial quantizer $\bar{Q}(\cdot)$ to produce N atoms, and that the rate and distortion are (r, d) for the resulting modulus quantizer and (R, D) for the overall system. A second quantizer $\bar{Q}'(\cdot)$ produces rate and distortion pairs (r', d') and (R', D') for the same number of atoms N . If $d' < d$, then $D' < D$.

Although this assumption seems reasonable, it is not strictly true, particularly when the change $|d - d'|$ is small. However, for larger changes in distortion the assumption does approximately hold for our system, as demonstrated in Figure 12. Sample frames are coded each to a constant number of atoms N as the UTQ stepsize is varied. The same sample frames and stepsize values are used as

Fig. 12. Overall image distortion D plotted against mean square distortion of the modulus quantizer d for various quantization stepsizes. Conditions shown are: (a) Container QCIF coded to 110 atoms, (b) Foreman CIF coded to 735 atoms, and (c) Stefan SIF coded to 1419 atoms.



in Figure 11. Each stepsize value produces a mean quantizer distortion d and an overall distortion D which are plotted on the horizontal and vertical axes, respectively. Since each plot is strictly increasing, the assumption is shown to hold.

We are now ready to show that the optimal ECSQ for the DZ-subtracted modulus distribution is also optimal in terms of (R, D) for the overall system.

Claim 1: For a given frame, define the optimal scalar quantizer for the DZ-subtracted modulus source as the quantizer which minimizes overall distortion D for the given rate R . If $\bar{Q}_{opt}(\cdot)$ is the optimal quantizer, then it must take the form of the optimal ECSQ for the DZ-subtracted modulus distribution produced by the encoder. That is, for a given rate r , $\bar{Q}_{opt}(\cdot)$ should minimize d .

Proof: See Appendix B.

We conclude that $\bar{Q}_{opt}(\cdot)$ should be the optimal ECSQ for the DZ-subtracted modulus source, and thus takes the form of a UTQ. However, the stepsize QP must still be optimized for the given frame and target rate. This is the topic of the next section.

C. Optimal Stepsize

In the previous section, the optimal quantizer $\bar{Q}_{opt}(\cdot)$ for application to the DZ-subtracted modulus source was shown to be a UTQ. Although the form of the quantizer is known, the step size has not been determined, and a given target rate may be achieved through many combinations of quantizer step size and number of coded atoms. We now look at the problem of choosing the optimal combination of QP and N which minimize distortion D for a given target rate R .

Figure 13 shows the proposed 2-pass quantization algorithm in greater detail. At the start of matching pursuit residual coding, motion compensation has been performed and so R_{hdr} and R_{mot} are known. The first pass decomposes residual frame \vec{f} using matching pursuit without quantization. The goal of this pass is to determine the number of atoms N_k needed to hit the

Fig. 13. Basic 2-pass adaptive quantization algorithm for the matching pursuit system.

```

For each frame  $k$ :
• Compute header budget ( $R_{hdr}$ ).
• Perform motion estimation; Compute motion budget ( $R_{mot}$ ).
• Create motion residual frame  $\vec{f}$ .
• Pass 1: Decompose  $\vec{f}$  without quantization to determine the
            $N_k$  needed to hit target rate  $R$  and resulting dead
           zone parameter  $DZ_k$ .
• Pass 2: Decompose  $\vec{f}$  to  $N$  atoms using in-loop quantization.
           The quantizer is defined by  $(DZ_k, QP_k)$ .
• Transmit Pass 2 atoms ( $R_{atm}$ ).
• Increment  $k$ ; repeat.

```

target rate R and also approximate the dead zone parameter DZ_k for use in the second quantized pass. To determine the stopping point, a progressive estimate of R_{atm} is needed, since the exact value cannot be computed in the absence of quantization. The estimated value $\tilde{R}_{atm(i)}$ is computed in the following way. After each stage, the i atoms found so far are quantized as a post-operation and coded in order to estimate the rate. The quantizer for rate estimation is based on DZ-subtraction using the minimum $|p_i|$ found so far as the value for DZ . This is followed by uniform quantizer $\bar{Q}(\cdot)$ with stepsize QP . For the moment, consider QP to be arbitrarily set by the user. Pass 1 then continues until the following stopping criteria is reached:

$$R_{hdr} + R_{mot} + \tilde{R}_{atm(i)} > R \quad (9)$$

When condition 9 is satisfied, Pass 1 stops and the number of atoms N_k and minimum absolute modulus DZ_k are preserved for use in Pass 2.

The step size QP_k is the only degree of optimization freedom in the 2-pass algorithm described above. It is thus possible, although costly, to find the optimal QP_k using an exhaustive linear search. This is done for a few sample frames in Figure 14. For each coded frame, the overall distortion D for the 2-pass algorithm is plotted against step size QP . A distinct minimum is found in each case, however the optimal step size varies greatly with sequence and bit rate.

We have found experimentally that a much more consistent value is the ratio between the optimal step size QP_k and the associated dead zone value DZ_k . This is demonstrated in Figure 15. In this figure, the same three sample frames are coded multiple times with the 2-pass algorithm, each coding reflecting a different quantizer ratio QP/DZ . Note that the optimal ratio is much more consistent than the optimal step size, falling near 0.6 for each of the three sample frames. This suggests that we choose QP in our 2-pass algorithm so that the quantizer ratio is near this optimal

Fig. 14. Search for best quantizer step QP for single coded frames. Overall image distortion D is plotted against quantization step QP for the 2-pass algorithm. (a) Container QCIF coded to 24 kbit/s, (b) Foreman CIF coded to 336 kbit/s, and (c) Stefan SIF coded to 1 Mbit/s.

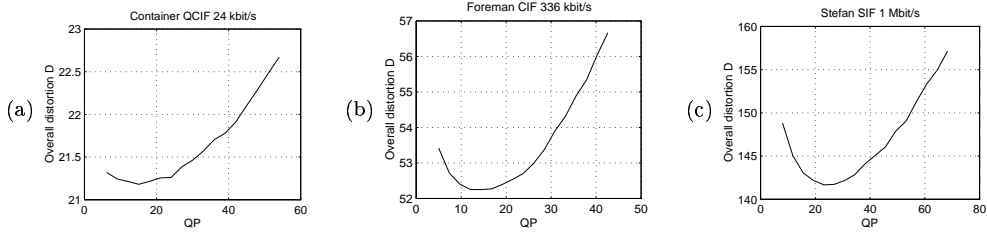
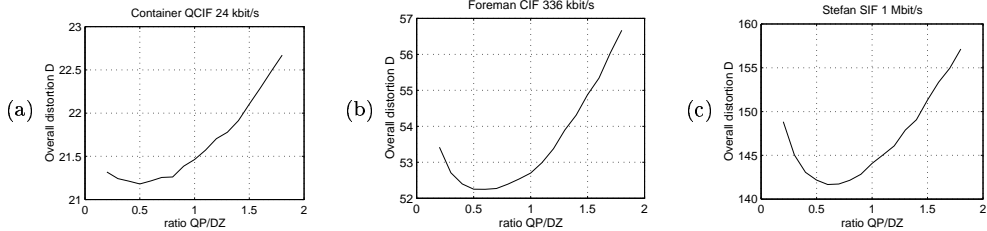


Fig. 15. Search for best ratio QP/DZ for single coded frames. Overall image distortion D is plotted against quantizer ratio for the 2-pass algorithm. (a) Container QCIF coded to 24 kbit/s, (b) Foreman CIF coded to 336 kbit/s, and (c) Stefan SIF coded to 1 Mbit/s.



value. To maintain a constant ratio QP_k/DZ_k in the 2-pass algorithm of Figure 13, we modify the rate estimation quantizer used to compute $\tilde{R}_{atm(i)}$ during Pass 1. As before, DZ is computed as the minimum $|p_i|$ found for the i atoms found so far. QP is then calculated to maintain the desired quantizer ratio. When stopping criteria (9) is reached, current values N_k , DZ_k and QP_k are passed on for use in the second quantized pass.

One convenient property seen in Figure 15 is that the curves are relatively flat in the vicinity of the optimal ratio. This suggests that the overall distortion D is not sensitive to small deviations from the optimal ratio, and so choosing a quantizer ratio near the optimal value will be sufficient. To see why the curves are flat near the optimal ratio, consider the effects on distortion due to changing the quantizer ratio. If the ratio is increased, then QP is increased for a constant DZ . This has two opposing effects on distortion. First, a larger QP increases the quantization error and tends to increase distortion D . However, the reduced bit cost to quantization allows more atoms to be coded which tends to reduce D . The optimal ratio occurs at the point when these opposing effects are equal and opposite, and hence the derivative of D with respect to ratio QP/DZ is zero. It is not hard to imagine a range of ratios around the optimal in which the opposing effects are nearly equal. This would explain the flat portions of the curves seen in Figure 15.

Table IV shows the result of the optimal ratio search procedure for a wider range of sample

TABLE IV

SEARCH FOR BEST RATIO QP/DZ FOR SINGLE CODED FRAMES. SAMPLE FRAMES CODED AT ONE AND THREE TIMES THE STANDARD MPEG TEST BIT RATE. OPTIMAL RATIO VALUE AND RESULTING PSNR ARE SHOWN FOR EACH OF 32 TEST POINTS.

Sequence	Base Rate				3x Base Rate			
	DZ	QP	Ratio	PSNR	DZ	QP	Ratio	PSNR
Cont10	54.32	27.16	0.5	31.33	23.15	11.58	0.5	35.96
Hall10	54.96	21.99	0.4	31.85	23.75	9.50	0.4	36.29
Mother10	43.92	13.18	0.3	34.19	18.89	9.44	0.5	38.34
Cont24	30.12	15.06	0.5	34.87	14.15	8.49	0.6	38.89
Silent24	39.10	15.64	0.4	32.60	17.84	7.14	0.4	36.90
Mother24	25.95	7.78	0.3	37.18	10.85	6.51	0.6	41.79
Coast48	40.97	20.49	0.5	31.36	19.58	11.75	0.6	35.44
Foreman48	45.05	27.03	0.6	32.92	14.89	11.91	0.8	38.93
News48	67.63	40.58	0.6	32.46	24.04	12.02	0.5	37.67
Coast112	97.84	29.35	0.3	27.60	40.22	24.13	0.6	32.21
Foreman112	106.95	53.47	0.5	30.09	24.08	14.45	0.6	36.97
News112	40.09	28.06	0.7	35.47	14.75	10.33	0.7	40.65
Mobile1024	51.27	30.76	0.6	28.00	20.66	18.59	0.9	33.87
Stefan1024	38.52	23.11	0.6	31.85	13.36	13.36	1.0	38.54
FunFair1024	46.78	37.42	0.8	30.54	16.78	16.78	1.0	36.45
Tunnel1024	34.09	17.05	0.5	31.58	13.24	9.27	0.7	37.35

frames and bit rates. A sample frame from each of the standard MPEG-4 test sequences is coded at the usual test rate and three times that rate. Two additional high bit rate sequences, FunFair and Tunnel, are also coded. The two pass algorithm using a linear search for best ratio QP_k/DZ_k is again performed, and the results summarized in the table. Note that the optimal ratio ranges from 0.3 to 1.0, with an average of 0.57 over all 32 tested frames. By our earlier logic, the shape of the curves in Figure 15 suggests that choosing the ratio to be near the optimal value should produce only small losses in distortion. This is verified for the test set in Table V, which shows the PSNR loss due to using constant ratio values of 0.6 or 1.0 instead of the optimal ratio found by linear search. For a constant ratio value of 0.6, the loss is extremely small, less than .02 dB difference for all sequences coded up to 1 Mbit/s. For the 3 Mbit/s sequences, the loss is a bit higher, with values in the .05 dB to .13 dB range. For such high coded bit rates, a constant ratio value of 1.0 is a better choice, as shown in the last column of the table.

We thus conclude that an adaptive 2-pass algorithm using dead zone subtraction and a constant value of step size to dead zone ratio achieves nearly optimal performance for the matching pursuit video coding system. For our experiments, a ratio value of 0.6 will be used for all rates up to 1 Mbit/s and a value of 1.0 will be used for higher rates. Note that our quantizer uses simple midpoint reconstruction. Experiments suggest that the more optimal centroid reconstruction [12] provides little or no gain [16].

TABLE V
 PERFORMANCE LOSS FOR SAME 32 SAMPLE FRAMES USING CONSTANT RATIO VALUES OF 0.6 AND 1.0.
 PSNR DIFFERENCES ARE RELATIVE TO THE BEST RATIO FOUND USING LINEAR SEARCH.

Sequence	Base Rate			3x Base Rate		
	Best Ratio	Ratio=0.6 PSNR Diff	Ratio=1.0 PSNR Diff	Best Ratio	Ratio=0.6 PSNR Diff	Ratio=1.0 PSNR Diff
Container10	0.5	0.0042	0.0448	0.5	0.0097	0.0486
Hall10	0.4	0.0152	0.0776	0.4	0.0145	0.0971
Mother10	0.3	0.0197	0.0853	0.5	0.0115	0.0793
Container24	0.5	0.0071	0.0576	0.6	0.0	0.0572
Silent24	0.4	0.0139	0.0842	0.4	0.0153	0.1265
Mother24	0.3	0.0023	0.0723	0.6	0.0	0.0547
Coast48	0.5	0.0141	0.0727	0.6	0.0	0.0641
Foreman48	0.6	0.0	0.0616	0.8	0.0201	0.0288
News48	0.6	0.0	0.0642	0.5	0.0097	0.0704
Coast112	0.3	0.0099	0.0638	0.6	0.0	0.0662
Foreman112	0.5	0.0061	0.0456	0.6	0.0	0.0373
News112	0.7	0.0006	0.0262	0.7	0.0085	0.0714
Mobile1024	0.6	0.0	0.0634	0.9	0.0573	0.0239
Stefan1024	0.6	0.0	0.0736	1.0	0.1287	0.0
FunFair1024	0.8	0.0073	0.0273	1.0	0.0711	0.0
Tunnel1024	0.5	0.0069	0.0866	0.7	0.0124	0.0575

V. SINGLE-PASS ADAPTIVE SOLUTION

The main disadvantage of the two-pass algorithm is the computational cost of the first pass. This initial pass performs a pre-analysis of the residual image in order to determine DZ and QP to achieve near-optimal quantizer performance. This doubles the complexity of the encoder, and so a less expensive estimate of these parameters is highly desirable. In this section, we present a method for predicting DZ for the current frame from known previous frame values. Prediction of DZ is sufficient, since QP follows from the constant quantizer ratio.

Consider the 2-pass algorithm and suppose that the encoded modulus values for frame k are $\{\hat{p}_{(k,i)}, i = 1 \dots N\}$ from the first pass and $\{\check{p}_{(k,i)}, i = 1 \dots N\}$ from the second pass. Recall that dead zone parameter DZ_k is normally determined from Pass 1 results as:

$$DZ_k = \min\{|\hat{p}_{(k,i)}|\} \quad (10)$$

To eliminate the analysis pass, we make two observations. First, quantization has little effect on the modulus distribution, as demonstrated earlier in Figure 9. This allows us to approximate:

$$\min\{|\hat{p}_{(k,i)}|\} \approx \min\{|\check{p}_{(k,i)}|\} \quad (11)$$

If we also assume that the dead zone parameter changes slowly from frame to frame, then we have the approximation:

$$DZ_k \approx DZ_{k-1} \quad (12)$$

By combining these with Equation (10), we may approximate DZ_k as:

$$DZ_k \approx \min\{|\tilde{p}_{(k-1,i)}|\} \quad (13)$$

The dead zone parameter for the current frame is thus defined by the Pass 2 decomposition of the previous frame, and so the analysis pass is eliminated altogether.

This scheme works to some extent, however it is not always true that DZ_k changes slowly from frame to frame. To compensate for this, we may consider the initial energy $E = \|\vec{f}_0\|^2$ of the current frame. If the residual for frame k has a higher initial energy E_k , all else equal, we would expect the distribution of $\{|p_i|, i = 1 \dots N\}$ for the frame to shift to the right and effectively increase the dead zone size. Figure 16a confirms this relationship for the Mother sequence coded using the 2-pass algorithm at 24 kbit/s. Each coded frame is represented by a circle, with the ratio of current to previous frame dead zone values on the horizontal axis and the corresponding ratio of initial frame energies on the vertical. A correlation between the two ratios can be seen, and the dotted line $y = x$ suggests the following energy-corrected prediction for the current frame dead zone:

$$DZ_k = DZ_{k-1} \cdot \left(\frac{E_k}{E_{k-1}} \right) \quad (14)$$

The effectiveness of these prediction methods is evaluated in Figure 16b. Data is from the 2-pass encoding of Mother at 24 kbit/s, and the value of DZ_k is predicted with energy correction according to Equation (14) and without correction as in Equation (13). The energy corrected prediction is seen to be better on average, however it is not guaranteed to be better for individual frames. If we compute mean square prediction error in DZ_k , then energy correction reduces the error from 15.03 to 5.45.

The improved DZ prediction method is thus used to define a one-pass adaptive quantization algorithm. This algorithm is similar to the two-pass algorithm of Figure 13, except Pass 1 is omitted and the dead zone parameter is predicted according to Equation (14). If there is no previous coded frame, the dead zone parameter is arbitrarily set to a reasonable value $DZ = 15$. The resulting 1-pass algorithm is summarized in Figure 17.

There is one potential problem of combining DZ -subtraction with a 1-pass algorithm. Without the benefit of a pre-analysis pass, the predicted DZ_k may be large enough to allow the single encoder pass to produce atoms which violate the condition $|p_i| > DZ_k$. When this occurs, we

Fig. 16. (a) Justification of the energy correction for DZ prediction. Horizontal axis shows ratio of current to previous frame DZ values as found by the full 2-pass adaptive algorithm. Vertical axis shows ratio of current to previous frame residual energy before atom coding. (b) Absolute error in predicting DZ_k for each frame of Mother at 24 kbit/s. Dotted line shows error for simple prediction. Solid line shows reduced error when energy corrected prediction is used.

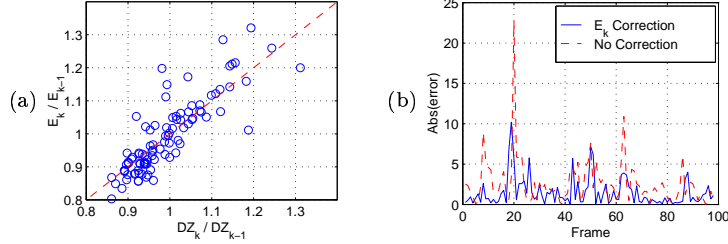


Fig. 17. The 1-pass adaptive quantization algorithm for the matching pursuit system.

```

For each frame  $k$ :
  • Compute header budget ( $R_{hdr}$ ).
  • Perform motion estimation; Compute motion budget ( $R_{mot}$ ).
  • Create motion residual frame  $\tilde{f}$ .
  ▷ If (previous frame data exists)
    Set  $DZ_k$  according to Equation (14).
  ▷ Else
    Set  $DZ_k = 15$ .
  • Set  $QP_k$  according to  $DZ_k$  and desired ratio.
  • Decompose  $\tilde{f}$  using in-loop quantization.
    Quantizer is defined by  $(DZ_k, QP_k)$ .
    Dead zone is divided into high-rate (ShiftVLC) bins.
    Code until progressive rate estimate equals target rate.
  • Transmit resulting atoms ( $R_{atm}$ ).
  • Increment  $k$ ; repeat.

```

revert back to the dead zone splitting technique used for the modified fixed quantizer of Figure 4b. That is, ShiftVLC-like bins are thus added to allow such small moduli to be coded. Note that these were not necessary for the full 2-pass algorithm, since the approximation in (11) holds very well in practice.

The resulting 1-pass system can thus be thought of as a combination of the fixed ShiftVLC and adaptive 2-pass algorithms. The advantages of each algorithm are retained, as will be observed in the next section.

VI. COMPARISON OF RESULTS

This section compares the three quantization schemes presented in the previous sections on the basis of PSNR performance and complexity. The schemes are summarized as:

ShiftVLC: Fixed quantization using midread linear quantizer with $QP = 30$ and three added high-rate quantization bins.

TABLE VI
AVERAGE LUMINANCE PSNR FOR THE STANDARD MPEG-4 SEQUENCES AND RATES.

Sequence	Shift VLC	Two Pass	One Pass	2Pass-1Pass	2Pass-ShiftVLC	1Pass-ShiftVLC
Container10	31.27	31.50	31.41	0.09	0.23	0.14
Hall10	32.12	32.23	32.29	-0.07	0.11	0.18
Mother10	33.28	33.52	33.47	0.05	0.23	0.19
Container24	34.39	34.55	34.46	0.08	0.16	0.08
Silent24	32.27	32.36	32.35	0.01	0.09	0.08
Mother24	36.54	36.69	36.63	0.05	0.14	0.09
Coast48	30.04	30.21	30.16	0.05	0.17	0.12
Foreman48	31.88	32.18	32.10	0.08	0.29	0.21
News48	33.02	33.25	33.20	0.05	0.23	0.18
Coast112	26.42	26.60	26.58	0.02	0.18	0.16
Foreman112	30.94	31.12	31.07	0.05	0.18	0.12
News112	36.06	36.15	36.12	0.03	0.09	0.06
Mobile1024	27.00	27.37	27.26	0.11	0.37	0.26
Stefan1024	32.04	32.25	32.15	0.10	0.22	0.12

2-Pass: Two-pass adaptive quantization. The first analysis pass sets the DZ and QP for the second quantized pass.

1-Pass: Single-pass adaptive quantization. Similar to 2-Pass, but analysis pass is omitted and DZ and QP are predicted from known values in the previous frame. Three high-rate quantization bins are added for safety, similar to the ShiftVLC case.

All three encoders were implemented and trained on a set of sequences outside the standard MPEG-4 test set. Coding was done on luminance only at a constant bit rate. Each ten-second sequence has a single leading intra frame, which was encoded using DCT followed by an in-the-loop deblocking filter. All other frames were encoded using motion compensated matching pursuit with one of the associated quantization schemes.

The PSNR results are summarized in Tables VI and VII. Table VI shows average PSNR results for standard MPEG-4 test sequences and rates. The rightmost three columns are the most interesting, since they show the PSNR differences between the three schemes. The first of these columns shows that the 1-Pass algorithm comes very close in performance to the 2-Pass algorithm, almost always within 0.1 dB. The second and third difference columns show the improvement of the adaptive algorithms over the fixed quantizer. These differences are also shown graphically in Figure 18. A consistent improvement over ShiftVLC is seen across all tested sequences. For the 2-Pass case, the improvement ranges from 0.09 to 0.37 dB with an average value of .19 dB. The 1-Pass algorithm also shows consistent improvement over ShiftVLC, with an average gain of .14 dB.

Fig. 18. Average Y-PSNR gains over ShiftVLC for the 1-Pass and 2-Pass algorithms. See Table I for sequence indices.

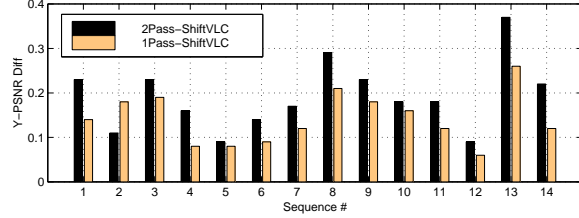
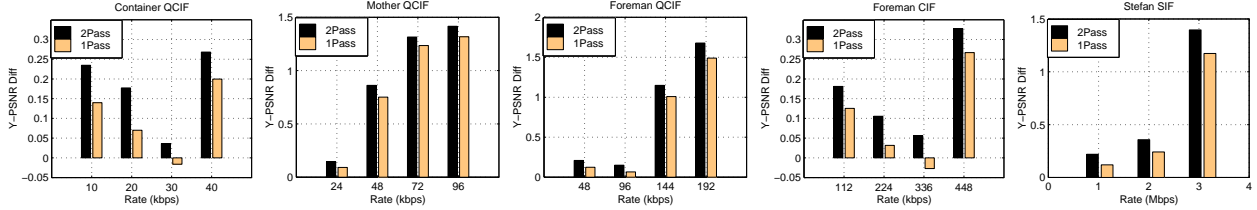


Fig. 19. Average Y-PSNR gains over ShiftVLC for five sequences at various bit rates.



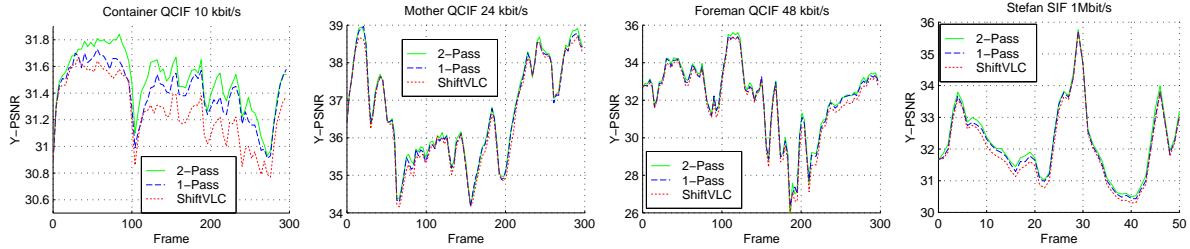
We expect the 2-Pass result to outperform the 1-Pass scheme, however such performance cannot be guaranteed. This is illustrated by the Table VI result for Hall at 10 kbit/s, where 1-Pass is 0.07 dB better than 2-Pass. This result can be explained by recalling that the 2-Pass scheme maintains a constant QP to DZ ratio, and thus chooses a stepsize that is only approximately optimal. It is thus possible for the 1-Pass scheme to do better, however this is seen to be the exception rather than the rule.

Performance at higher bit rates is shown in Table VII. Results are shown for five of the MPEG-4 test sequences chosen to represent a range of encoding conditions. For each sequence, the standard MPEG-4 test rate is used as a base, and several multiples of this rate are also encoded. Again it is useful to focus on the rightmost three columns which describe the PSNR differences. The first of these columns shows that 1-Pass performance is again close to that of the full 2-Pass algorithm, typically less than 0.10 dB. The second and third difference columns effectively show how the improvement over ShiftVLC due to the adaptive schemes varies with increasing bit rate. These differences are also shown graphically in Figure 19. It is seen that the adaptive schemes perform much better than the fixed quantizer at higher bit rates, and that for three of the five sequences the improvements fall in the 1 to 2 dB range for the highest tested bit rates. The largest improvements are 1.68 dB for the QCIF Foreman sequence coded at 192 kbit/s, and 1.40 dB for the SIF Stefan sequence at 3 Mbit/s. For these high rates, the visual improvement is mainly seen as a reduction of ringing around sharp edges.

TABLE VII
AVERAGE LUMINANCE PSNR FOR ASCENDING BIT RATES.

Sequence	Shift VLC	Two Pass	One Pass	2Pass-1Pass	2Pass-ShiftVLC	1Pass-ShiftVLC
Container10	31.27	31.51	31.41	0.09	0.23	0.14
Container20	34.19	34.36	34.26	0.11	0.18	0.07
Container30	35.91	35.94	35.89	0.05	0.04	-0.02
Container40	36.74	37.01	36.94	0.07	0.27	0.20
Mother24	36.54	36.69	36.64	0.05	0.15	0.09
Mother48	38.77	39.63	39.52	0.11	0.86	0.75
Mother72	40.26	41.58	41.50	0.08	1.31	1.24
Mother96	41.51	42.93	42.83	0.10	1.42	1.32
Foreman48	33.15	33.36	33.27	0.09	0.21	0.12
Foreman96	37.09	37.24	37.16	0.08	0.15	0.07
Foreman144	38.38	39.53	39.39	0.14	1.15	1.01
Foreman192	39.55	41.22	41.04	0.19	1.68	1.49
Foreman112	30.96	31.14	31.08	0.06	0.18	0.13
Foreman224	35.06	35.16	35.09	0.07	0.11	0.03
Foreman336	36.93	36.99	36.90	0.08	0.06	-0.03
Foreman448	37.92	38.24	38.18	0.06	0.33	0.27
Stefan1024	32.04	32.26	32.16	0.10	0.22	0.12
Stefan2048	35.99	36.35	36.23	0.12	0.36	0.24
Stefan3072	37.75	39.14	38.92	0.22	1.40	1.17

Fig. 20. Luminance PSNR vs. frame for four sample sequences at standard MPEG-4 test rates.



Figures 20 and 21 show a more detailed look at PSNR vs. Frame for specific sequences. Four sequences coded at standard MPEG-4 test rates are depicted in Figure 20. In each case, it is evident that the adaptive schemes show a consistent improvement over the fixed quantizer, and that the 1-Pass algorithm has performance only slightly below that of the 2-Pass algorithm. Higher encoding rates for the same sequences are shown in Figure 21. It is clear from these plots that the gain for the adaptive schemes is much higher at these rates, and that again the 1-Pass and 2-Pass schemes are nearly equivalent.

Fig. 21. Luminance PSNR vs. frame for four sample sequences at higher bit rates.

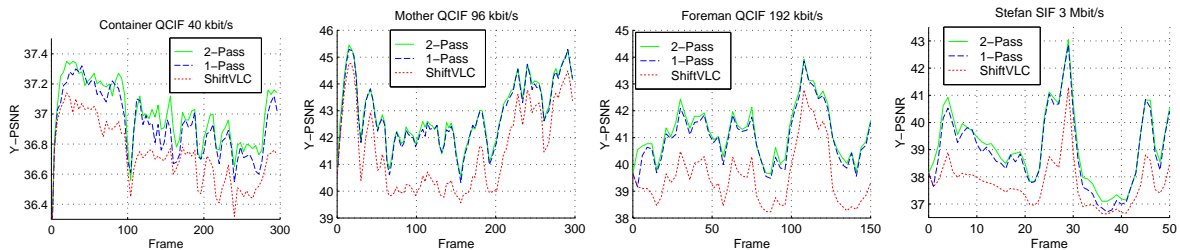
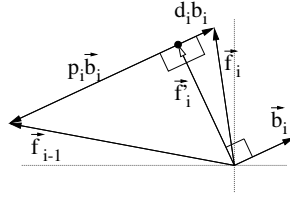


Fig. 22. Two dimensional vector representation of a single quantized matching pursuit stage. This is used in the proof of Equation (5), as shown in Appendix A.



VII. CONCLUSIONS

Frame adaptive quantizer designs have been shown to outperform the fixed quantizer used in earlier matching pursuit video coding work. Our design uses a combination of adaptive dead zone subtraction and a UTQ with adaptive stepsize. A 2-pass algorithm was shown to choose a nearly optimal balance between the quantizer precision and the number of coded atoms, and hence comes close to minimizing overall coding distortion D for a given target rate R . A 1-pass algorithm was shown to have nearly the same performance, but with complexity similar to that of the original fixed quantizer design. Small improvements of .1 to .4 dB are shown for the MPEG sequences at standard rates, and larger gains of up to 1.7 dB can be shown for the same sequences at higher bit rates.

APPENDIX A

Equation 5 stated that the single-stage energy reduction for matching pursuit can be represented as:

$$ER_i = p_i^2 - d_i^2$$

Proof: Since p_i is the projection of \vec{f}_{i-1} onto basis element \vec{b}_i , the difference $\vec{f}_i^r = \vec{f}_{i-1} - p_i \vec{b}_i$ will be orthogonal to \vec{b}_i , as shown in Figure 22. \vec{f}_i^r would be the remaining signal energy for the current stage in the absence of quantization, and is also orthogonal to the quantization error $d_i \vec{b}_i$. We can thus apply the pythagorean theorem twice to get:

$$\|\vec{f}_{i-1}\|^2 = \|p_i \vec{b}_i\|^2 + \|\vec{f}_i^r\|^2 \quad (15)$$

$$\|\vec{f}_i^r\|^2 = \|\vec{f}_i\|^2 + \|d_i \vec{b}_i\|^2 \quad (16)$$

Since \vec{f}_i is the remaining signal energy with quantization, we can combine the two above equations to see that the single stage energy reduction is:

$$\|\vec{f}_{i-1}\|^2 - \|\vec{f}_i\|^2 = \|p_i \vec{b}_i\|^2 - \|d_i \vec{b}_i\|^2 \quad (17)$$

$$= (|p_i|^2 - |d_i|^2) \|\vec{b}_i\|^2 \quad (18)$$

which is equal to $p_i^2 - d_i^2$ since \vec{b}_i is unit-norm. \square

APPENDIX B

Proof of Claim 1: $\bar{Q}_{opt}(\cdot)$ is the optimal scalar quantizer for the system, which therefore codes the given residual signal to rate R using some number of atoms N to achieve the minimum overall distortion D . Take (r, d) to be the quantizer rate and distortion produced by $\bar{Q}_{opt}(\cdot)$. For the moment, assume that $\bar{Q}_{opt}(\cdot)$ is not the optimal ECSQ for the DZ-subtracted modulus distribution produced in coding. Then there exists a quantizer which achieves a lower scalar quantization error at the same average rate, that is

$$\exists \bar{Q}'(\cdot) \text{ s.t. } r' = r \text{ and } d' < d$$

Note that (r', d') is the rate-distortion point produced by quantizer $\bar{Q}'(\cdot)$. We are implicitly assuming that the DZ-subtracted modulus distribution is unaffected by application of either in-loop quantizer. This is a reasonable assumption considering the results shown earlier in Figure 9.

Assumption 1 shows that this new quantizer $\bar{Q}'(\cdot)$ will hit the target rate R with the same number of atoms as the optimal quantizer, since for coding N atoms,

$$R'_Q = Nr' = Nr = R_Q$$

and the other contributions $(R_{hdr} + R_{mot} + R_P + R_I)$ are approximately constant with change in quantizer.

Since the number of atoms N is constant and $d' < d$, Assumption 2 implies that $D' < D$. This is a contradiction, since $\bar{Q}_{opt}(\cdot)$ is given to be the optimal quantizer for the system. Therefore our momentary assumption must be incorrect, and $\bar{Q}_{opt}(\cdot)$ must be the optimal ECSQ for the given modulus distribution. \square

REFERENCES

- [1] ISO/IEC 13818-2, "MPEG-2 video coding standard," March 1995.
- [2] ITU-T Recommendation H.263, "Video coding for low bit rate communication," ITU Document, September 1997.
- [3] ISO/IEC JTC1/SC29/WG11, "MPEG-4 visual coding standard," Final Draft of International Standard, December 1998.
- [4] S. A. Martucci, Iraj Sodagar, T. Chiang, and Y.-Q. Zhang, "A zerotree wavelet coder," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 7, no. 1, pp. 109–118, February 1997.
- [5] Ralph Neff and Avidesh Zakhor, "Matching pursuit video coding at very low bit rates," *IEEE Transactions on Circuits and Systems for Video Technology*, pp. 158–171, February 1997.
- [6] Stéphane Mallat and Zhifeng Zhang, "Matching pursuits with time-frequency dictionaries," *IEEE Transactions on Signal Processing*, vol. 41, no. 12, pp. 3397–3415, December 1993.
- [7] Ralph Neff, Toshio Nomura, and Avidesh Zakhor, "Decoder complexity and performance comparison of matching pursuit and MPEG-4 video codecs," in *Proceedings of the IEEE International Conference on Image Processing*, 1998.
- [8] Osama Al-Shaykh, Eugene Miloslavsky, Toshio Nomura, Ralph Neff, and Avidesh Zakhor, "Video compression using matching pursuits," *IEEE Transactions on Circuits and Systems for Video Technology*, September 1998.
- [9] C. de Vleeschouwer and B. Macq, "New dictionaries for matching pursuit video coding," in *Proceedings of the IEEE International Conference on Image Processing*, 1998, pp. 764–768.
- [10] D. Redmill, D. Bull, and P. Czrepiński, "Video coding using a fast non-separable matching pursuits algorithm," in *Proceedings of the IEEE International Conference on Image Processing*, 1998, pp. 769–773.
- [11] ISO/IEC JTC1/SC29/WG11, "MPEG-4 video verification model 11.0," MPEG Document N2172, March 1998.
- [12] N.S. Jayant and Peter Noll, *Digital Coding of Waveforms*, Prentice-Hall, Englewood Cliffs, NJ, 1984, Chapter 4.
- [13] Kendall Atkinson, *Elementary Numerical Analysis*, John Wiley and Sons, New York, 1985, Chapter 7.
- [14] Alan Gersho and Robert Gray, *Vector Quantization and Signal Compression*, Kluwer Academic Publishers, Boston, 1991, Chapter 9.
- [15] G. Sullivan, "Efficient scalar quantization of exponential and laplacian random variables," *IEEE Transactions on Information Theory*, vol. 42, no. 5, pp. 1365–1374, September 1996.
- [16] Ralph Neff, *Practical Methods for Matching Pursuit Video Compression*, Ph.D. thesis, U.C. Berkeley, December 1999.