# Distributed Grayscale Stereo Image Coding with Unsupervised Learning of Disparity

David Varodayan, Aditya Mavlankar, Markus Flierl and Bernd Girod

Max Planck Center for Visual Computing and Communication

Stanford University, Stanford, CA 94305

Email: {varodayan, maditya, mflierl, bgirod}@stanford.edu

**Abstract**

Distributed compression is particularly attractive for stereo images since it avoids communication between cameras. Since compression performance depends on exploiting the redundancy between images, knowing the disparity is important at the decoder. Unfortunately, distributed encoders cannot calculate this disparity and communicate it. We consider the compression of grayscale stereo images, and develop an Expectation Maximization algorithm to perform unsupervised learning of disparity during the decoding procedure. Towards this, we devise a novel method for joint bitplane distributed source coding of grayscale images. Our experiments with both natural and synthetic 8-bit images show that the unsupervised disparity learning algorithm outperforms a system which does no disparity compensation by between 1 and more than 3 bits/pixel and performs nearly as well as a system which knows the disparity through an oracle.

## I. Introduction

Colocated pixels from pairs of stereo images are strongly statistically dependent after compensation for disparity induced by the geometry of the scene. Much of the disparity between these images can be characterized as shifts of foreground objects relative to the background. Assuming that the disparity information and occlusions can be coded compactly, joint compression is much more efficient than separate encoding and decoding. Surprisingly, distributed lossless encoding combined with joint decoding can be just as efficient as the wholly joint system, according to the Slepian-Wolf theorem [1]. Distributed compression is preferred because it avoids communication between the stereo cameras. The difficulty, however, lies in discovering and exploiting the scene-dependent disparity at the decoder, while keeping the transmission rate low.

A similar situation arises in low-complexity Wyner-Ziv encoding of video captured by a single camera [2] [3] [4]. These systems encode frames of video separately and decode them jointly, so discovering the motion between successive frames at the decoder is helpful. A very computationally burdensome way to learn the motion is to run the decoding algorithm with every motion realization [3]. Another approach requires the encoder to transmit additional hash information, so the decoder can perform suitable motion compensation before running the decoding algorithm [5]. Since the encoder transmits the hashes at a constant rate, it wastes bits when the motion is small. On the other hand, if there is too much change between frames, the fixed-rate hash may be insufficient for reliable motion search. Due to the drawbacks of excessive computation and difficulty of rate allocation for the hash, we use neither of these approaches towards compression of stereo images.

Fig. 1. (a) Source image $X$ (8-bit 72-by-88 pixels), (b) horizontal disparity legend, (c)-(f) source images $Y$ (8-bit 72-by-88 pixels) with respective 8-by-8 block-wise horizontal disparity fields $D$ of $X$ with respect to $Y$

In Section II, we review our unsupervised disparity learning algorithm [6] for distributed lossless compression of pairs of binary random dot stereograms [7]. In Section III, we extend the system to distributed lossless compression of grayscale natural stereo images. Note that, in this work, we only exploit pixel-wise redundancy between stereo images and not the spatial redundancy within a single image; that is the topic of future work. We then describe the algorithm formally within the framework of Expectation Maximization (EM) [8] in Section IV. Section V reports our simulation results for natural and synthetic images.

## II. BACKGROUND

The relationship between a pair of stereo images $X$ and $Y$ in terms of their disparity $D$ is illustrated in Fig. 1. A sample source image $X$ is depicted in Fig. 1(a). Fig. 1(c)-(f) show four realizations of source image $Y$ taken from different viewpoints [9]. For each pair, the respective block-wise horizontal disparity field $D$ indicates which 8-by-8 block of $Y$ (among candidates shifted up to 5 pixels horizontally) best matches each 8-by-8 block of $X$ (in terms of mean square error). These stereo image pairs, when viewed stereoscopically, create an illusion of depth: various parts of the scene appear on different planes, according to the value of the disparity field in those parts.

The compression setup used both in [6] for the binary case and the current paper is shown in Fig. 2. Images $X$ and $Y$ are encoded separately and decoded jointly. For simplicity, we assume that $Y$ is conventionally coded and is available at the decoder. The challenge is to encode $X$ efficiently in the absence of $Y$ so that it can be reliably decoded in the presence of $Y$. The Slepian-Wolf theorem states that $X$ can be communicated losslessly to the decoder using $R$ bits on average as long as $R > H(X|Y)$ [1].

We now review the binary case, where coding was performed on the pixel values of $X$ directly [6]. Fig. 3 depicts three compression systems that can be applied to the binary

Fig. 2. Distributed compression: separate encoding and joint decoding

random dot stereogram problem. The baseline system in Fig. 3(a) is due to [10] and performs compression of $X$ with respect to the colocated pixels of $Y$ without disparity compensation. The encoder computes the syndrome $S$ (of length $R$ bits) of $X$ with respect to a low-density parity-check (LDPC) code [11]. The decoder initially estimates $X$ statistically using the colocated pixels of $Y$ and refines these estimates using $S$ via an iterative belief propagation algorithm. When disparity is introduced between $X$ and $Y$, this scheme performs badly because the estimates of $X$ are poor in shifted regions. For comparison, Fig. 3(b) shows an impractical scheme in which the decoder is endowed with a disparity oracle. The oracle informs the decoder which pixels of $Y$ should be used to inform the estimates of the pixels of $X$ during LDPC decoding. Finally, Fig. 3(c) depicts the practical decoder of [6] that learns disparity $D$ via EM. In place of the disparity oracle, a disparity estimator maintains an *a posteriori* probability distribution on $D$. Every iteration of LDPC decoding sends the disparity estimator a soft estimate of $X$ (denoted by $\theta$) in order to refine the distribution on $D$. In return, the disparity estimator updates the side information $\psi$ for the LDPC decoder by blending information from the pixels of $Y$ according to the refined distribution on $D$. A formal EM treatment of this algorithm for the case of binary random dot stereograms is included in [6].

## III. DISTRIBUTED COMPRESSION OF GRAYSCALE STEREO IMAGES

For grayscale pairs $X$ and $Y$, the compression systems of Fig. 3(a) and (b) may be applied to each bitplane of $X$ separately, which are then used as additional side information for the decoding of subsequent bitplanes [4] [12]. Unfortunately, conditional bitplane distributed source coding of this type cannot be applied with the disparity learning decoder of Fig. 3(c) because efficient disparity estimation requires the intermediate soft estimate $\theta$ to be calculated using all bitplanes at once. A possible alternative is to perform symbol encoding and decoding using LDPC codes in a high-order Galois field. Instead, in order to leverage the better understanding of binary LDPC codes, we devise a novel joint bitplane distributed source coding scheme and apply it to all three systems in Fig. 3. In joint bitplane distributed source coding, all bitplanes are encoded and decoded together, permitting the computation of $\theta$ over gray levels.

The joint bitplane LDPC encoder is simply the LDPC encoder of [10] applied to the bit representation of $X$. The belief propagation decoding graph of the joint bitplane LDPC decoder is shown in Fig. 4. Like the LDPC decoder of [10], it decodes $X$ at bit nodes subject to constraints set at syndrome nodes, by propagating log likelihood ratios of the

Fig. 3. Distributed compression for stereo images with (a) no disparity compensation, (b) a disparity oracle, and (c) unsupervised learning of disparity $D$ via EM

bit beliefs along the edges of the graph. (Note that a log likelihood ratio is the logarithm of the ratio of the likelihood of a certain bit being 1 to the likelihood of it being 0.) The difference is that the soft side information $\psi$ no longer supplies log likelihood ratios to the bit nodes directly, but instead feeds the new symbol nodes. As depicted in Fig. 4, each symbol node (one per pixel) aggregates $\psi_{pixel}$ together with log likelihood ratios $\log \frac{\alpha_g}{1-\alpha_g}$ for $g \in \{1, \ldots, b\}$ from each bit node associated with that pixel. After local computation, the symbol node sends log likelihood ratios $\log \frac{\beta_h}{1-\beta_{h_j}}$ for $h \in \{1, \ldots, b\}$ to each of those bit nodes and outputs a soft estimate $\theta_{pixel}$ of the pixel.

In the case of Fig. 4, $b = 3$ so the side information distribution is 8-valued, say,

$$\psi_{pixel} = (p_{000}, p_{001}, p_{010}, p_{011}, p_{100}, p_{101}, p_{110}, p_{111}).$$

Using most-significant-bit first notation, we choose the log likelihood ratio sent to bit node 2 (for example) to be

$$\log \frac{\beta_2}{1-\beta_2} = \log \frac{p_{010}(1-\alpha_1)(1-\alpha_3) + p_{011}(1-\alpha_1)\alpha_3 + p_{110}\alpha_1(1-\alpha_3) + p_{111}\alpha_1\alpha_3}{p_{000}(1-\alpha_1)(1-\alpha_3) + p_{001}(1-\alpha_1)\alpha_3 + p_{100}\alpha_1(1-\alpha_3) + p_{101}\alpha_1\alpha_3}.$$

This formula begins by multiplying the elements of $\psi_{pixel}$ by elements of the binary distributions $(\alpha_1, 1-\alpha_1)$ and $(\alpha_3, 1-\alpha_3)$ according to the values of the first and third bits,

Fig. 4. Belief propagation decoding graph of the joint bitplane LDPC decoder

respectively, of their indices. Then these products are summed in two groups according to the value of the second bit, before the log likelihood ratio is taken. For this calculation of $\log \frac{\beta_2}{1-\beta_2}$, we avoid using $(\alpha_2, 1 - \alpha_2)$ to prevent recycling information to bit node 2. The calculation of $\log \frac{\beta_1}{1-\beta_1}$ and $\log \frac{\beta_3}{1-\beta_3}$ follow analogously by shuffling the roles of the first, second and third bits of the index. Similarly, we compute the output soft estimate $\theta_{pixel}$ by multiplying the elements of $\psi_{pixel}$ by elements of all three binary distributions $(\alpha_1, 1 - \alpha_1)$, $(\alpha_2, 1 - \alpha_2)$ and $(\alpha_3, 1 - \alpha_3)$ according to the values of the first, second and third bits, respectively, of their indices. After normalization,

$$
\begin{aligned}
\theta_{pixel} &= (q_{000}, q_{001}, q_{010}, q_{011}, q_{100}, q_{101}, q_{110}, q_{111}) \\
q_{000} &\propto p_{000}(1 - \alpha_1)(1 - \alpha_2)(1 - \alpha_3) \\
q_{001} &\propto p_{001}(1 - \alpha_1)(1 - \alpha_2)\alpha_3 \\
q_{010} &\propto p_{010}(1 - \alpha_1)\alpha_2(1 - \alpha_3) \\
q_{011} &\propto p_{011}(1 - \alpha_1)\alpha_2\alpha_3 \\
q_{100} &\propto p_{100}\alpha_1(1 - \alpha_2)(1 - \alpha_3) \\
q_{101} &\propto p_{101}\alpha_1(1 - \alpha_2)\alpha_3 \\
q_{110} &\propto p_{110}\alpha_1\alpha_2(1 - \alpha_3) \\
q_{111} &\propto p_{111}\alpha_1\alpha_2\alpha_3
\end{aligned}
$$

With these modifications, we extend distributed lossless compression of binary random dot stereogram to grayscale stereo images, by replacing the LDPC encoder and decoder blocks in Fig. 3 by the joint bitplane LDPC encoder and decoder described here. Section IV presents a formal EM treatment of the proposed disparity learning system.

## IV. EXPECTATION MAXIMIZATION ALGORITHM

### A. Model

Let $Y$ be a $b$-bit grayscale image of size $m$-by-$n$. Let $D$ represent an $m$-by-$n$ horizontal disparity field with $|D(i, j)| \leq l$, where $l \ll n$ is the maximum possible magnitude.

Then $X$ is also a $b$-bit grayscale image of size $m$-by-$n$, disparity-compensated from $Y$ through disparity $D$ with independent Laplacian noise $Z$ added. We model the decoder's *a posteriori* probability distribution of source $X$ based on parameters $\theta$ as

$$
\begin{aligned}
P_{app}\{X\} &= P\{X; \theta\} \\
&= \prod_{i,j} \theta(i, j, X(i, j))
\end{aligned}
$$

where $\theta(i, j, w) = P_{app}\{X(i, j) = w\}$ defines a soft estimate of $X(i, j)$ over gray values $w \in \{0, \ldots, 2^b - 1\}$. The restriction that the disparity field $D$ have small maximum magnitude and be in one dimension is reasonable for a pair of closely-spaced cameras.

*B. Problem*

The decoder aims to calculate the *a posteriori* probability distribution of the disparity $D$,

$$
\begin{aligned}
P_{app}\{D\} &:= P\{D|Y, S; \theta\} \\
&\propto P\{D\}P\{Y, S|D; \theta\},
\end{aligned}
$$

with the second step by Bayes' Law. The form of this expression suggests an iterative EM solution. The E-step updates the disparity field distribution with reference to the source model parameters, while the M-step updates the source model parameters with reference to the disparity field distribution.

*C. E-step Algorithm*

The E-step updates the estimated distribution on $D$ and before renormalization is written as

$$
P_{app}^{(t+1)}\{D\} := P_{app}^{(t)}\{D\}P\{Y, S|D; \theta^{(t+1)}\}.
$$

But this operation is expensive due to the large number of possible values of $D$. We simplify in two ways. First, we ignore knowledge of the syndrome $S$ since it is exploited in the M-step of LDPC decoding. Second, we permit the estimation of the horizontal disparity field $D$ as block-by-block disparity shifts $L_{u,v}$. For a specified blocksize $k$, every $k$-by-$k$ block of $\theta$ is compared to the colocated block of $Y$ as well as all those shifted between $-l$ and $l$ pixels horizontally. For a block $\theta_{u,v}$ with top left pixel located at $(u, v)$, the distribution on the shift $L_{u,v}$ is updated as below and normalized:

$$
P_{app}^{(t+1)}\{L_{u,v}\} := P_{app}^{(t)}\{L_{u,v}\}P\{Y_{u,v+L_{u,v}}|L_{u,v}; \theta_{u,v}^{(t+1)}\},
$$

where $Y_{u,v+L_{u,v}}$ is the $k$-by-$k$ block of $Y$ with top left pixel at $(u, v + L_{u,v})$. Note that $P\{Y_{u,v+L_{u,v}}|L_{u,v}; \theta_{u,v}\}$ is the probability of observing $Y_{u,v+L_{u,v}}$ given that it was generated through shift $L_{u,v}$ from $X_{u,v}$ as parameterized by $\theta_{u,v}$. This procedure, shown in the left hand side of Fig. 5, occurs in the disparity estimator.

Fig. 5. E-step disparity field estimation (left) and side information blending (right)

## D. M-step Algorithm

The M-step updates the model parameters $\theta$ by maximizing the likelihood of $Y$ and the syndrome $S$.

$$\theta^{(t+1)} \quad := \quad \arg\max_\theta P\{Y, S; \theta^{(t)}\}$$

$$= \quad \arg\max_\theta \sum_d P_{app}^{(t)}\{D = d\} P\{Y, S|D = d; \theta^{(t)}\}$$

True maximization is intractable, so we approximate it with an iteration of joint bitplane LDPC decoding, introduced in Section III. The joint bitplane LDPC decoder's input side information $\psi_{u,v}$ is created by blending estimates from each of the blocks $Y_{u,v+L_{u,v}}$ according to $P_{app}^{(t)}\{L_{u,v}\}$, as shown in the right hand side of Fig. 5. More generally, the probability that the blended side information has value $w$ at pixel $(i, j)$ is

$$\psi(i, j, w) \quad = \quad \sum_d P_{app}^{(t)}\{D = d\} P\{X(i, j) = w|D = d, Y\}$$

$$= \quad \sum_d P_{app}^{(t)}\{D = d\} p_Z(w - Y(i, j + d)),$$

where $p_Z(z)$ is the probability mass function of the independent additive noise $Z$. The symbol nodes in the joint bitplane LDPC decoder combine this side information distribution with incoming log likelihood ratios $\log \frac{\alpha_g}{1-\alpha_g}$ for $g \in \{1, \ldots, b\}$ from their connected bit nodes. They send log likelihood ratios $\log \frac{\beta_h}{1-\beta_h}$ for $h \in \{1, \ldots, b\}$ to their connected

bit nodes, according to

$$\log \frac{\beta_h}{1-\beta_h} = \log \frac{\sum_{w:w_b[h]=1} \psi(i,j,w) \prod_{g \neq h} (\alpha_g \mathbf{1}_{[w_b[g]=1]} + (1-\alpha_g)\mathbf{1}_{[w_b[g]=0]})}{\sum_{w:w_b[h]=0} \psi(i,j,w) \prod_{g \neq h} (\alpha_g \mathbf{1}_{[w_b[g]=1]} + (1-\alpha_g)\mathbf{1}_{[w_b[g]=0]})},$$

where $w_b[h]$ denotes the $h$th most-significant-bit in the binary representation of gray value $w$ and $\mathbf{1}_{[.]}$ denotes the indicator function. The symbol nodes also produce the next soft estimate of the source $X$, which before normalization over $w \in \{0, \ldots, 2^b - 1\}$ is computed as

$$\theta^{(t+1)}(i,j,w) := \psi(i,j,w) \prod_{g=1}^{b} (\alpha_g \mathbf{1}_{[w_b[g]=1]} + (1-\alpha_g)\mathbf{1}_{[w_b[g]=0]}).$$

*E. Termination*

Iterating between the E-step and the M-step in this way provides a profile of the disparity at the granularity of $k$-by-$k$ blocks. The decoding algorithm terminates successfully when the hard estimates $\hat{X}(i,j) = \arg\max_w \theta(i,j,w)$ yield a syndrome equal to $S$.

## V. SIMULATION RESULTS

We compare the performance of the joint bitplane systems in Fig. 3 and the Slepian-Wolf bound $H(X|Y)$ for both natural and synthetic grayscale stereo images, using constants: image height $m = 72$, image width $n = 88$, number of bitplanes $b = 8$, maximum horizontal shift $l = 5$, blocksize $k = 8$. Rate control is implemented by using rate-adaptive regular degree 3 LDPC accumulate codes of length 50688 bits [13] as a platform for the joint bitplane systems. In these experiments, the decoder is provided with a good value for the variance of the Laplacian noise $Z$. After 150 decoding iterations, if $\hat{X}$ still does not satisfy the syndrome condition, the decoder requests additional incremental transmission from the encoder via a feedback channel. For the disparity learning algorithm, the distributions of $L_{u,v}$ are initialized to

$$P_{app}^{(0)}\{L_{u,v}\} := \begin{cases} 0.75, & \text{if } L_{u,v} = 0; \\ 0.025, & \text{if } L_{u,v} \neq 0. \end{cases}$$

Table I compares the compression bit-rates of the joint bitplane systems and the pixel-wise Slepian-Wolf bounds for the combinations of natural stereo images $X$ and $Y$ in Fig. 1. The system that learns disparity unsupervised outperforms the system that allows no disparity compensation by between 1 and more than 3 bits/pixel for these stereo images. As the average magnitude of the disparity field grows, the rate for no compensation increases but the rate for unsupervised learning is robust and closely follows the rate for the oracle-assisted scheme. Fig. 6 shows a sample evolution of disparity probability distribution for an 8-by-8 block.

We also perform simulations with synthetic grayscale random dot stereograms, which have well-defined statistics and easy-to-compute bounds. Each image by itself resists compression, but substantial savings are possible by exploiting the disparity field. A pair $X$ and $Y$ has disparity field $D$ constant in some rectangular region and equal to zero elsewhere. Fig. 7 shows realizations of $X$ and $Y$ (with $D$ equal to 5 in a 32-by-32 region and zero elsewhere) and their log absolute differences under relative shifts of 0 and -5.

Fig. 8 shows the compression performance for 8-bit random dot stereograms and a lower bound on $H(X|Y)$. The entropy $H(Z)$ of the additive Laplacian noise ranges

| Side information image $Y$ | Fig. 1(c) | Fig. 1(d) | Fig. 1(e) | Fig. 1(f) |
|---|---|---|---|---|
| Pixel-wise $H(X|Y)$ (bits/pixel) | 3.93 | 3.88 | 3.84 | 3.91 |
| Oracle-assisted rate (bits/pixel) | 4.48 | 4.48 | 4.48 | 4.48 |
| Unsupervised rate (bits/pixel) | 4.48 | 4.48 | 4.48 | 4.61 |
| No compensation rate (bits/pixel) | 5.58 | 6.06 | 7.51 | 8 |

TABLE I

BIT-RATE OF $X$ (IN BITS/PIXEL) WITH $Y$ IN FIG. 1 FOR JOINT BITPLANE SYSTEMS IN FIG. 3.



Fig. 6. Evolution of a disparity probability distribution for a sample 8-by-8 block with true disparity value $L_{u,v} = 5$.

from $0.8$ to $4$ bits/pixel. The disparity field $D$ takes a constant value drawn uniformly from $\{-5, \ldots, 5\}$ in a 32-by-32 pixel region; elsewhere $D = 0$. For the proposed disparity learning scheme, we show results when the 32-by-32 disparity region is aligned with the 8-by-8 block grid (best case) and when it is offset from the grid by 4 pixels horizontally and vertically (worst case). These results show that best case unsupervised learning of disparity can save 2 bits/pixel compared to the system that allows no disparity compensation for this setup. Moreover, it incurs negligible performance loss with respect to the impractical oracle-assisted scheme because, in the best case, the 8-by-8 decoder granularity matches the effective granularity of the oracle. This gap being relatively smaller than the gap in the binary case [6] is due to more data being used in the grayscale case to learn the same disparity field. The further gap from the oracle-assisted scheme to the lower bound on $H(X|Y)$ is due to the inefficiency of moderate length LDPC codes.

## VI. CONCLUSIONS

We extend the iterative EM algorithm [6] for distributed lossless stereo image compression with disparity learning at the decoder to the case of grayscale stereo images, by inventing a method for joint bitplane distributed source coding. For natural 8-bit stereo images, our proposed scheme demonstrates compression savings of between 1 and more than 3 bits/pixel compared to a system that does no disparity compensation, and performs only



Fig. 7. (a) Source image $X$ (b) Source image $Y$ (c) Log absolute difference of $X$ and $Y$ (d) Log absolute difference of $X$ and $Y$ (shifted to realign the rectangular nonzero disparity region)

Fig. 8. Rate (in bits/pixel) required to communicate grayscale random dot stereogram $X$ for the joint bitplane systems shown in Fig. 3, for learning $D$ constant drawn uniformly $\{-5, \ldots, 5\}$ in a 32-by-32 region, and $D = 0$ elsewhere.

negligibly worse than an oracle-assisted scheme. Simulations with synthetic grayscale random dot stereograms confirm these findings, especially in the best case, where decoder granularity matches the effective oracle granularity. In future work, we intend to extend the system to exploit spatial redundancy within stereo images as well.

## REFERENCES

[1] D. Slepian and J. K. Wolf, "Noiseless coding of correlated information sources," *IEEE Trans. Inform. Theory*, vol. 19, no. 4, pp. 471–480, July 1973.

[2] A. Aaron, R. Zhang, and B. Girod, "Wyner-Ziv coding of motion video," in *Proc. Asilomar Conf. on Signals, Syst., Comput.*, Pacific Grove, CA, 2002.

[3] R. Puri and K. Ramchandran, "PRISM: A new robust video coding architecture based on distributed compression principles," in *Proc. Allerton Conf. Commun., Contr. and Comput.*, Allerton, IL, 2002.

[4] A. Aaron, S. Rane, E. Setton, and B. Girod, "Transform-domain Wyner-Ziv codec for video," in *Proc. Visual Commun. and Image Processing*, San Jose, CA, 2004.

[5] A. Aaron, S. Rane, and B. Girod, "Wyner-Ziv video coding with hash-based motion compensation at the receiver," in *Proc. IEEE International Conf. on Image Processing*, Singapore, 2004.

[6] D. Varodayan, A. Mavlankar, M. Flierl, and B. Girod, "Distributed coding of random dot stereograms with unsupervised learning of disparity," in *Proc. IEEE International Workshop on Multimedia Signal Processing*, Victoria, BC, Canada, 2006.

[7] B. Julesz, "Binocular depth perception of computer generated patterns," *Bell Sys. Tech. J*, vol. 38, pp. 1001–1020, 1960.

[8] A. Dempster, N. Laird, and D. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *J. Royal Stat. Soc., Series B*, vol. 39, no. 1, pp. 1–38, 1977.

[9] D. Scharstein and R. Szeliski, "High-accuracy stereo depth maps using structured light," in *Proc. Comput. Vision and Pattern Recog.*, Madison, WI, 2003.

[10] A. Liveris, Z. Xiong, and C. Georghiades, "Compression of binary sources with side information at the decoder using LDPC codes," *IEEE Commun. Lett.*, vol. 6, no. 10, pp. 440–442, Oct. 2002.

[11] R. G. Gallager, "Low-density parity-check codes," *Cambridge MA: MIT Press*, 1963.

[12] J. Chen, A. Khisti, D. M. Malioutov, and J. S. Yedidia, "Distributed source coding using serially-concatenated-accumulate codes," in *IEEE Inform. Theory Workshop*, San Antonio, TX, 2004.

[13] D. Varodayan, A. Aaron, and B. Girod, "Rate-adaptive distributed source coding using low-density parity-check codes," in *Proc. Asilomar Conf. on Signals, Syst., Comput.*, Pacific Grove, CA, 2005.