

End-to-end optimization of nonlinear transform codes for perceptual quality

Johannes Ballé, Valero Laparra, Eero P. Simoncelli

Center for Neural Science and Courant Institute of Mathematical Sciences

New York University, New York, NY, USA

{johannes.balle,valero,eero.simoncelli}@nyu.edu

Abstract—We introduce a general framework for end-to-end optimization of the rate–distortion performance of nonlinear transform codes assuming scalar quantization. The framework can be used to optimize any differentiable pair of analysis and synthesis transforms in combination with any differentiable perceptual metric. As an example, we consider a code built from a linear transform followed by a form of multi-dimensional local gain control. Distortion is measured with a state-of-the-art perceptual metric. When optimized over a large database of images, this representation offers substantial improvements in bitrate and perceptual appearance over fixed (DCT) codes, and over linear transform codes optimized for mean squared error.

I. INTRODUCTION

Transform coding [1] is one of the most successful areas of signal processing. Virtually all modern image and video compression standards operate by applying an invertible transformation to the signal, quantizing the transformed data to achieve a compact representation, and inverting the transform to recover an approximation of the original signal.

Generally, these transforms have been linear. Non-Gaussian/nonlinear aspects of signal statistics are typically handled by augmenting the linear system with carefully selected nonlinearities (for example, companding nonlinearities to enable non-uniform quantization, prediction for hybrid compression, etc.). Deciding which combination of these operations, also known as “coding tools,” are ultimately useful is a cumbersome process. The operations are generally studied and optimized individually, with different objectives, and any proposed combination of coding tools must then be empirically validated in terms of average code rate and distortion.

This is reminiscent of the state of affairs in the field of object and pattern recognition about a decade ago. As in the compression community, most solutions were built by manually combining a sequence of individually designed and optimized processing stages. In recent years, that field has seen remarkable performance gains [2], which have arisen primarily because of end-to-end system optimization. Specifically, researchers have chosen architectures that consist of a cascade of transformations that are differentiable with respect to their parameters, and then used modern optimization tools to jointly optimize the full system over large databases of images.

Here, we take a step toward using such end-to-end optimization in the context of compression. We develop an optimization framework for nonlinear transform coding (fig. 1), which generalizes the traditional transform coding paradigm.

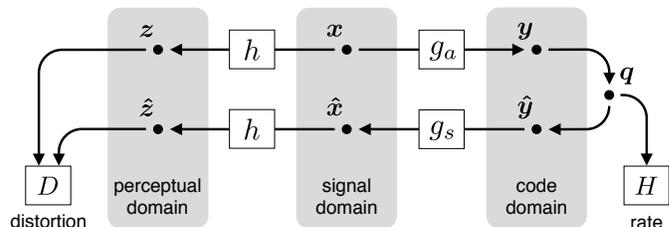


Fig. 1. Nonlinear transform coding optimization framework. See text.

An image vector x is transformed to a code domain vector using a differentiable function $y = g_a(x; \phi)$ (the analysis transform), parameterized by a vector ϕ (containing linear filter coefficients, for example). The transformed y is subjected to scalar quantization, yielding a vector of integer indices q and a reconstructed vector \hat{y} . The latter is then nonlinearly transformed back to the signal domain to obtain the reconstructed image $\hat{x} = g_s(\hat{y}; \theta)$, where this synthesis transform g_s is parameterized by vector θ .

The code rate is assessed by measuring the entropy, H , of the discrete probability distribution P_q of the quantization indices over an ensemble of images. Traditionally, the distortion is assessed directly in the image domain by taking the squared Euclidean norm of the difference between x and \hat{x} (or equivalently, the peak signal-to-noise ratio, PSNR). However, it is well known that PSNR is not well-aligned with human perception [3]. To alleviate this problem, we allow an additional “perceptual” transform of both vectors $z = h(x)$ and $\hat{z} = h(\hat{x})$, on which we then compute distortion using a suitable norm. A well-chosen transform h can provide a significantly better approximation of subjective visual distortion than PSNR (e.g., [4]).

II. OPTIMIZATION FRAMEWORK

In the transform coding framework given above, we seek to adjust the analysis and synthesis transforms g_a and g_s so as to minimize the rate–distortion functional:

$$L[g_a, g_s] = H[P_q] + \lambda \mathbb{E} \|z - \hat{z}\|. \quad (1)$$

The first term denotes the discrete entropy of the vector of quantization indices q . The second term measures the distortion between the reference image z and its reconstruction \hat{z} in a perceptual representation. Note that both terms are expectations taken over an ensemble of images.

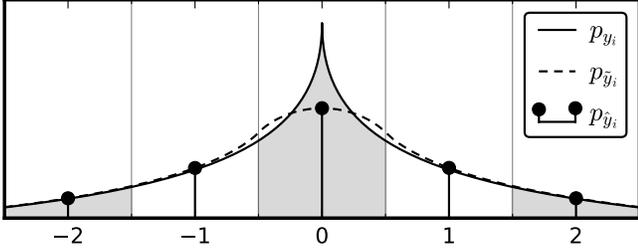


Fig. 2. Relationship between densities of y_i , \tilde{y}_i , and \hat{y}_i . $p_{\tilde{y}_i}$ is a continuous relaxation of the probability masses in each of the quantization bins.

We wish to minimize this objective over the continuous parameters $\{\theta, \phi\}$. Most optimization methods rely on differentiability, but both terms in the objective depend on the quantized values in \mathbf{q} , and the derivative of the quantizer is discontinuous (specifically, it is zero or infinite everywhere). To resolve this, we propose to approximate the objective function with one that is continuously differentiable, by replacing the deterministic quantizer with an additive uniform noise source.

A uniform scalar quantizer is a piecewise constant function applied to each of the elements of \mathbf{y} : $\hat{y}_i = \text{round}(y_i)$.¹ The marginal density of the quantized values is given by:

$$p_{\hat{y}_i}(t) = \sum_{n=-\infty}^{\infty} P_{q_i}(n) \delta(t - n), \quad (2)$$

where

$$P_{q_i}(n) = (p_{y_i} * \text{rect})(n), \text{ for all } n \in \mathbb{Z}, \quad (3)$$

is the probability mass in the n th quantization bin. Here, ‘*’ represents continuous convolution, and rect is a uniform distribution on $(-\frac{1}{2}, \frac{1}{2})$. If we add independent uniform noise to y_i , i.e., form the signal $\tilde{y}_i = y_i + \Delta y_i$, with $\Delta y_i \sim \text{rect}$, then the density of that signal is $p_{\tilde{y}_i} = p_{y_i} * \text{rect}$. $p_{\tilde{y}_i}$ is identical to P_{q_i} at all integer locations, and provides a *continuous relaxation* for intermediate values (fig. 2). We propose to optimize the differential entropy $h[p_{\tilde{y}_i}]$ as a proxy for the discrete entropy $H[P_{q_i}]$. To optimize it, we need a running estimate of $p_{\tilde{y}_i}$. This estimate need not be arbitrarily precise, since $p_{\tilde{y}_i}$ is band-limited by convolution with rect . Here, we simply use a non-parametric, piecewise linear function (a first-order spline approximation). We also use \tilde{y}_i rather than \hat{y}_i to obtain gradients of the distortion term. The overall objective can be written as:

$$L(\theta, \phi) = \mathbb{E}_{\mathbf{x}, \Delta \mathbf{y}} \left(-\log_2 p_{\tilde{\mathbf{y}}}(g_a(\mathbf{x}; \phi) + \Delta \mathbf{y}) \right. \\ \left. + \lambda \|h(g_s(g_a(\mathbf{x}; \phi) + \Delta \mathbf{y}; \theta)) - h(\mathbf{x})\| \right), \quad (4)$$

where $p_{\tilde{\mathbf{y}}}(\tilde{\mathbf{y}}) = \prod_i p_{\tilde{y}_i}(\tilde{y}_i)$. This is differentiable with respect to θ and ϕ and thus suited for stochastic gradient descent. Although finding a global optimum is not guaranteed, this optimization problem is similar to others which are encountered

¹Without loss of generality, we assume that the quantization bin size is 1, since we can always modify the analysis/synthesis transforms to include a rescaling. Further, we can implement non-uniform quantization by using nonlinear transforms (as in *companding*).

when optimizing deep neural networks, and which have been found to behave well in practice.

III. CHOICE OF PARAMETRIC TRANSFORMS

In a traditional transform code, both analysis and synthesis transforms are linear, and exact inverses of each other. In general, this need not be the case, so long as the overall system minimizes the rate–distortion functional. We have previously shown that a linear transform followed by a particular form of joint local gain control (generalized divisive normalization, GDN) is well-matched to the local probability structure of photographic images [5]. This suggests that jointly normalized representations might also prove useful for compression. To demonstrate the use of our optimization framework, we examine GDN as a candidate analysis transform, and introduce an approximate inverse as the corresponding synthesis transform. For the perceptual transform, we use the normalized Laplacian pyramid [4] (NLP), which mimics the local luminance and contrast behaviors of the human visual system.

A. Generalized divisive normalization (GDN)

The GDN transform consists of a linear decomposition \mathbf{H} followed by a joint nonlinearity, which divides each linear filter output by a measure of overall filter activity:

$$\mathbf{y} = g_a(\mathbf{x}; \phi) \quad \text{s.t.} \quad y_i = \frac{v_i}{(\beta_i + \sum_j \gamma_{ij} |v_j|^{\alpha_{ij}})^{\varepsilon_i}} \\ \text{and} \quad \mathbf{v} = \mathbf{H}\mathbf{x}, \quad (5)$$

with parameter vector $\phi = \{\alpha, \beta, \gamma, \varepsilon, \mathbf{H}\}$.

B. Approximate inverse of GDN

The approximate inverse we introduce here is based on the fixed point iteration for inversion of GDN introduced in [5]. It is similar in spirit to the LISTA algorithm [6], in that it uses the parametric form of the inversion iteration, but unties its parameters from their original values for faster convergence. We find that for purposes of image compression, one iteration is sufficient:

$$\hat{\mathbf{x}} = g_s(\hat{\mathbf{y}}; \theta) \quad \text{s.t.} \quad \hat{\mathbf{x}} = \mathbf{H}'\mathbf{w} \\ \text{and} \quad w_i = \hat{y}_i \cdot (\beta'_i + \sum_j \gamma'_{ij} |\hat{y}_j|^{\alpha'_{ij}})^{\varepsilon'_i}, \quad (6)$$

where the parameter vector consists of a distinct set of parameters: $\theta = \{\alpha', \beta', \gamma', \varepsilon', \mathbf{H}'\}$.

C. Normalized Laplacian pyramid (NLP)

The NLP imitates the transformations associated with the early visual system: local luminance subtraction and local gain control [4]. Images are decomposed using a Laplacian pyramid [7], which subtracts a local estimate of the mean luminance at multiple scales. Each pyramid coefficient is then divided by a local estimate of amplitude (a constant plus the weighted sum of absolute values of neighbors). Perceptual quality is assessed by evaluating a norm of the difference between reference and reconstruction in this normalized domain. The parameters (constant and weights used for amplitudes)

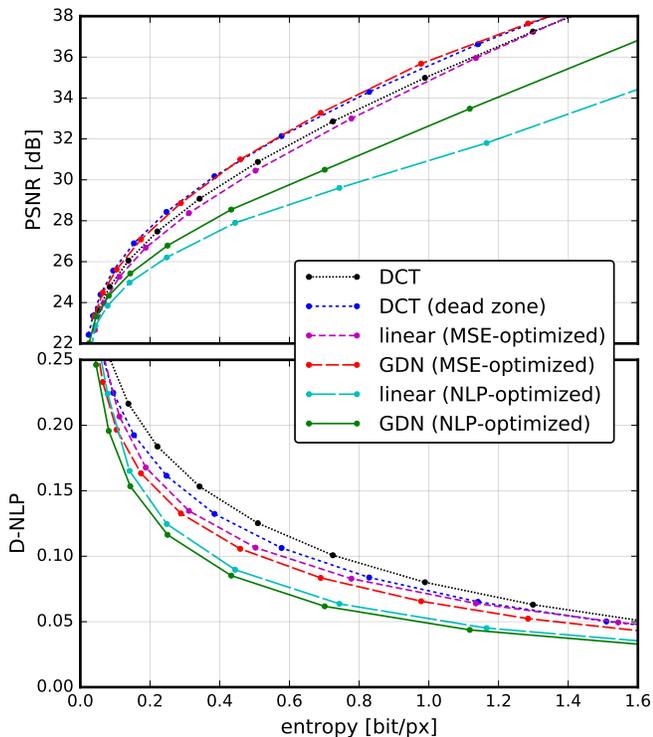


Fig. 3. Rate–distortion results averaged over the entire Kodak image set (24 images, 752×496 pixels each). Reported rates are average discrete entropy estimates. Reported distortion is average PSNR (top) and distance in the normalized Laplacian pyramid domain (bottom – see text).

are optimized to best fit perceptual data in the TID2008 database [8], which includes images corrupted by artifacts arising from compression with block transforms. This simple distortion measure provides a near-linear fit to the human perceptual judgments in the database, outperforming the widely-used SSIM [9] and MS-SSIM [10] quality metrics [4]. Examples and code are available at <http://www.cns.nyu.edu/~lcv/NLPyr>.

IV. EXPERIMENTAL RESULTS

The proposed framework can be used to optimize any differentiable pair of analysis and synthesis transforms in combination with any differentiable perceptual metric. Here, we consider two types of transform: a linear analysis and synthesis transform operating on 16×16 pixel blocks (in this case, θ and ϕ each only consist of 256×256 filter coefficients), and a 16×16 block GDN transform with the approximate inverse defined above (with θ and ϕ consisting of filter coefficients and normalization parameters, as defined above). We optimized each of these for two distortion metrics: mean squared error (MSE) and distance in the NLP domain [4]. Each combination of transform and distortion metric was optimized for different values of λ . We also include a fixed linear transform, the 16×16 discrete cosine transform (DCT), with or without dead-zone quantization, serving as a baseline. All other codes (i.e., those optimized in our framework) are constrained to use uniform quantization.

We used the Adam algorithm [11], a variant of stochastic gradient descent, to optimize all codes over a large collection

of images from the ImageNet database [2], initializing the parameters randomly. For each optimization step, we used a randomly selected mini-batch of 4 images of 128×128 pixels. To prevent overfitting to the training database, we performed all evaluations on a separate set of test images.²

We measured rate–distortion performance for all four transform/distortion metric combinations, along with the DCT transform. Note that the additive noise approximation was used only for optimization, not for evaluation: We evaluated rates by estimating discrete entropy of the quantized code vector \mathbf{q} . For each choice of λ , we optimized a separate set of transform parameters, which could be stored in the encoder and decoder. The only side information a real-world codec would need to transmit is the choice of λ and the image size (although it would be desirable to reduce the storage requirements by storing parameters jointly for different λ).

For evaluation of the distortion, we first computed the mean squared error (MSE) over the entire image set for each λ , and then converted these values into PSNRs (fig. 3, top panel). In terms of PSNR, the optimized linear transform is slightly worse than the DCT, because the statistics of the ImageNet database are slightly different from the Kodak set.³ The DCT with dead-zone quantization is better, but doesn’t outperform the MSE-optimized GDN transform, which uses only uniform quantization. The NLP-optimized transforms don’t perform well in terms of PSNR.

The situation is reversed, however, when we examine performance in terms of perceptual distortion (fig. 3, bottom). Here, we evaluated the norm in the NLP domain (D-NLP) for each image in the set, and then averaged across images. Note that this norm is almost (inversely) proportional to human perceptual mean opinion scores (MOS) across several image databases [4]. Overall, the combination of NLP and GDN achieves an impressive rate savings at similar quality when compared with MSE-optimized methods, and with the DCT (both uniform and dead-zone quantizers). It is also interesting to note that in terms of NLP distance, the optimized linear transform with uniform quantization outperforms both versions of the DCT. This may be because the optimized filters tend to be spatially localized (and oriented/bandpass), which possibly leads to visually less disturbing artifacts (not shown).

For visual evaluation, we show results on two example images (figs. 4 and 5). Results for the entire test set are available at <http://www.cns.nyu.edu/~balle/nlpgdn>. The figures serve to illustrate the main effect of using a perceptual metric that is aware of local *relative* contrast. Traditional, linear systems optimized for MSE give too much preference to high-contrast regions (e.g., the snow-covered mountains in the background, or the pebbles/debris in the foreground; fig. 4, center image). By performing joint normalization before quantization, the NLP-optimized GDN transform allocates more

²The Kodak image set, downloaded from <http://www.cipr.rpi.edu/resource/stills/kodak.html>. We converted the images to grayscale and discarded 8 pixels from each side to eliminate boundary artifacts.

³If we validate on a held-out test set of images from ImageNet, the two transforms perform equally well.



NLP-GDN, 0.190 bit/px. PSNR: 20.95 D-NLP: 0.21 MS-SSIM: 0.868



DCT (dead z.), 0.204 bit/px. PSNR: 21.81 D-NLP: 0.28 MS-SSIM: 0.827



Original, 8 bit/px. PSNR: ∞ D-NLP: 0 MS-SSIM: 1

Fig. 4. Example image from Kodak set (bottom), compressed with DCT and hand-optimized (for MSE) dead-zone quantization (middle), and GDN with uniform quantization optimized in the NLP domain (top). Cropped to fit page.

bits to represent detail in low-contrast regions (such as the forest in the depicted scene; top image). Overall, the rate allocation is perceptually more balanced, which leads to a more appealing visual appearance.

V. DISCUSSION

We have introduced a framework for end-to-end optimization of nonlinear transform codes, which can be applied to any set of parametric, differentiable analysis and synthesis transforms. By optimizing a nonlinear transform for a perceptual metric over a database of photographs, we obtained a



NLP-GDN, 0.044 bit/px. PSNR: 26.37 D-NLP: 0.21 MS-SSIM: 0.881



DCT (dead z.), 0.044 bit/px. PSNR: 26.93 D-NLP: 0.24 MS-SSIM: 0.857



Original, 8 bit/px. PSNR: ∞ D-NLP: 0 MS-SSIM: 1

Fig. 5. A second example image from the Kodak set (see caption for fig. 4).

nonlinear code that respects the perception of local luminance and contrast errors, allowing for significant rate savings.

The earliest instance of a linear transform optimized for signal properties may be the Karhunen–Loève transform (KLT), or principal components analysis (PCA). The DCT was originally introduced as an efficient approximation to the KLT for a separable autoregressive process of order 1 [12]. Other studies have optimized transform parameters specifically for perceptual compression (e.g., [13, 14]), but these were generally limited to optimizing weighting matrices for the DCT. Models that use “matched” nonlinear transformations as a means of converting to/from a more desirable representation

of the data are known in the machine learning literature as *autoencoders* [15]. However, we are unaware of any work that directly aims to optimize discrete entropy.

We assume uniform quantization in the transform domain, and replace quantization with additive uniform noise to relax the discontinuous problem into a differentiable one. We find that this method empirically outperforms results obtained by ignoring the effects of the discontinuous quantizer on the backpropagation of gradients (not shown). Under the presented conditions, adding noise is equivalent to performing dithered quantization [16]. Dithering was used in some early image coders [17], but generally does not improve rate–distortion performance. To our knowledge, it has not been used as a form of continuous relaxation for optimization purposes. While uniform quantization has been shown to be asymptotically optimal [18], it is well known that dead-zone quantization generally performs better for linear transform coding of images. Here, we demonstrate empirically that the use of nonlinear transforms with uniform quantization allows equivalent or better solutions, and our framework provides a means of finding these transforms.

Divisive normalization has previously been used in DCT-based image compression, e.g., [19, 20]. These approaches use the normalized representation both for coding and distortion estimation, reasoning that this domain is both perceptually and statistically uniform, and thus well-suited for both. The framework introduced here offers more flexibility, by allowing the perceptual domain and the code domain to be distinct (fig. 1). Further, previous methods required the decoder to invert the normalization transform, either by solving an iterative set of linear equations for every block [19], estimating the multipliers (i.e., the values of the denominators) from neighboring blocks [20], or embedding the multipliers into the code as side information. Our framework eliminates this problem by introducing a highly efficient approximate inverse transform, which is jointly optimized along with the normalization transform.

There are several directions in which to proceed with this work. Well-known techniques to improve performance of linear transform codes, such as run-length encoding, adaptive entropy coding, and signal-adaptive techniques in general, should be investigated in the context of nonlinear transform coding. Furthermore, our framework offers a means for exploring much more sophisticated nonlinear analysis/synthesis transforms as well as perceptual metrics, since it is built on the highly successful paradigm of end-to-end optimization over training data.

ACKNOWLEDGEMENT

JB and EPS are supported by the Howard Hughes Medical Institute. VL is supported by the APOSTD/2014/095 Generalitat Valenciana grant (Spain).

REFERENCES

[1] V. K. Goyal, “Theoretical foundations of transform coding,” *IEEE Signal Processing Magazine*, vol. 18, no. 5, 2001. DOI: 10.1109/79.952802.

[2] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “ImageNet: a large-scale hierarchical image database,” in *IEEE Conf. on Computer Vision and Pattern Recognition*, 2009. DOI: 10.1109/CVPR.2009.5206848.

[3] B. Girod, “What’s wrong with mean squared error?,” in *Digital Images and Human Vision*. M.I.T. Press, 1993, ISBN: 0-262-23171-9.

[4] V. Laparra, J. Ballé, A. Berardino, and E. P. Simoncelli, “Perceptual image quality assessment using a normalized Laplacian pyramid,” in *Proceedings of SPIE, Human Vision and Electronic Imaging XXI*, 2016.

[5] J. Ballé, V. Laparra, and E. P. Simoncelli, “Density modeling of images using a generalized normalization transformation,” *arXiv e-prints*, 2015, Presented at the 4th Int. Conf. for Learning Representations, 2016. arXiv:1511.06281.

[6] K. Gregor and Y. LeCun, “Learning fast approximations of sparse coding,” in *Proceedings of the 27th International Conference on Machine Learning*, 2010.

[7] P. J. Burt and E. H. Adelson, “The Laplacian pyramid as a compact image code,” *IEEE Transactions on Communications*, vol. 31, no. 4, Apr. 1983. DOI: 10.1109/TCOM.1983.1095851.

[8] N. Ponomarenko, V. Lukin, A. Zelensky, K. Egiazarian, M. Carli, and F. Battisti, “TID2008 – a database for evaluation of full-reference visual quality assessment metrics,” *Advances of Modern Radioelectronics*, vol. 10, pp. 30–45, 2009.

[9] Z. Wang, E. P. Simoncelli, and A. C. Bovik, “Multi-scale structural similarity for image quality assessment,” in *Conf. Rec. of the 37th Asilomar Conf. on Signals, Systems and Computers, 2004*, 2003. DOI: 10.1109/ACSSC.2003.1292216.

[10] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, “Image quality assessment: from error visibility to structural similarity,” *IEEE Transactions on Image Processing*, vol. 13, no. 4, Apr. 2004. DOI: 10.1109/TIP.2003.819861.

[11] D. P. Kingma and J. L. Ba, “Adam: a method for stochastic optimization,” *arXiv e-prints*, 2014, Presented at the 3rd Int. Conf. for Learning Representations, 2015. arXiv:1412.6980.

[12] N. Ahmed, T. Natarajan, and K. R. Rao, “Discrete cosine transform,” *IEEE Transactions on Computers*, vol. C-23, no. 1, 1974. DOI: 10.1109/T-C.1974.223784.

[13] A. J. Ahumada and H. A. Peterson, “Luminance-model-based DCT quantization for color image compression,” in *Proc. SPIE 1666, Human Vision, Visual Processing, and Digital Display III*, 1992. DOI: 10.1117/12.135982.

[14] A. B. Watson, “DCTune: a technique for visual optimization of DCT quantization matrices for individual images,” *Society for Information Display Digest of Technical Papers*, vol. 24, pp. 946–949, 1993.

[15] G. E. Hinton and R. R. Salakhutdinov, “Reducing the dimensionality of data with neural networks,” *Science*, vol. 313, no. 5786, 2006. DOI: 10.1126/science.1127647.

[16] L. Schuchman, “Dither signals and their effect on quantization noise,” *IEEE Transactions on Communication Technology*, vol. 12, no. 4, 1964. DOI: 10.1109/TCOM.1964.1088973.

[17] L. G. Roberts, “Picture coding using pseudo-random noise,” *IRE Transactions on Information Theory*, vol. 8, no. 2, 1962. DOI: 10.1109/TIT.1962.1057702.

[18] H. Gish and J. N. Pierce, “Asymptotically efficient quantizing,” *IEEE Transactions on Information Theory*, vol. 14, no. 5, 1968. DOI: 10.1109/TIT.1968.1054193.

[19] J. Malo, I. Epifanio, R. Navarro, and E. P. Simoncelli, “Non-linear image representation for efficient perceptual coding,” *IEEE Transactions on Image Processing*, vol. 15, no. 1, Jan. 2006. DOI: 10.1109/TIP.2005.860325.

[20] S. Wang, A. Rehman, Z. Wang, S. Ma, and W. Gao, “Perceptual video coding based on SSIM-inspired divisive normalization,” *IEEE Transactions on Image Processing*, vol. 22, no. 4, 2013. DOI: 10.1109/TIP.2012.2231090.