

Method, apparatus and software for color image compression based on non-linear perceptual representations and machine learning

Jesús Malo, Juan Gutiérrez, Gustavo Camps Valls and M^a José Luque

July 21, 2015

Abstract

It is provided a method for color image compression based on the following key steps: (1) expressing the coefficients of local space-frequency representation of achromatic and chromatic opponent channels in perceptually meaningful contrast units, (2) applying divisive normalization non-linearities including relations among the coefficients of the achromatic and chromatic contrast channels, and (3) using machine learning algorithms to select relevant data from the non-linearly normalized channels.

Psychophysical and numerical experiments have been conducted to estimate and validate the parameters of the contrast and divisive normalization transforms. Besides, an implementation of the proposed method has been experimentally compared to a JPEG implementation. In both cases 16×16 DCT blocks were used, and when applicable, equivalent chromatic and frequency dependent parameters were taken for a fair comparison. In these conditions, experimental results on a 25 color image database show that the proposed method gives rise to substantial compression gain over JPEG at the same (RMSE, SSIM and S-CIELab) distortion in the commercially meaningful bit-rate range [1, 2.2] bits/pix. Depending on the distortion measure applied, the compression gain is about 120% for the same RMSE, 45% for the same SSIM, and 170% for the same S-CIELab.

Apparatus for implementing the method and software product are also claimed.

Keywords: Color Image coding/compression, Support Vector Machine (SVM), Support Vector Regression (SVR), Regression, Function Approximation, Adaptive Insensitivity, Discrete Cosine Transform (DCT), Perceptual Distortion, Achromatic Perceptual Non-linearity, Red-Green Non-linearity, Yellow-Blue Non-Linearity, Gain Control, Color Contrast, Sparse Coding, Sparseness/Compactedness, Quantization, JPEG.

Contents

1	Technical and scientific field of the invention	3
2	Description of the prior art	3
3	Proposed alternatives	4
4	Summary and general overview of the invention	5
4.1	Scheme of the invention	5
4.2	Further capabilities	8
5	Detailed description to carry out the invention	9
5.1	Contrast transforms	10
5.1.1	Formulation: Contrast of achromatic and chromatic gratings from the transform coefficients of images in opponent color spaces	11
5.1.2	Experiment 1: Blue-Yellow and Red-Green luminance proportions in uniform brightness gratings	15
5.1.3	Experiment 2: Maximum modulation in opponent spaces in natural images	16
5.1.4	Summary: achromatic and chromatic contrast definition	17
5.2	Non-linear perceptual transforms	19
5.2.1	Divisive normalization transforms (Achromatic, Yellow-Blue and Red-Green)	19
5.2.2	Experiment 3: checking the inversion of the non-linear transform	21
5.3	Support Vector Regression (SVR) with Adaptive Profile	22
5.3.1	Standard SVR	23
5.3.2	SVR with Adaptive Insensitivity Profile	24
5.3.3	Remarks: adaptive insensitivity and the working domain	25
6	Experimental results to assess the performance of the method	25
6.1	Experimental setup	25
6.2	Numerical comparison with JPEG: rate-distortion and compression gain	26
6.3	Visual comparison	29

1 Technical and scientific field of the invention

The present invention is related to a method, apparatus and software for being used in coding or compression of color images. The method uses as basic tools a proper linear (or non-linear) spatio-frequency transformation, followed by perceptually-adapted contrast functions *per* color channel, different non-linear perceptual transforms *per* color channel, and the application of a machine learning method for adaptive selection of the most informative data, which are finally properly encoded. The method relates to, but not exclusively, still and moving color image compression for transmission and/or storage.

2 Description of the prior art

Nowadays, the volume of imaging data increases exponentially in a very wide variety of applications, such as remote sensing, digital photography and consumer camcorders, medical imaging, digital libraries and documents, movies, and video-conferences. This poses several problems and needs for transmitting, storing and retrieving images. As a consequence, digital image compression is becoming a crucial technology. However, compressing an image is significantly different than compressing raw binary data, given their particular statistical properties, and thus the direct application of general-purpose compression methods is far from being optimal. Therefore, statistical knowledge about the problem becomes extremely important to develop efficient coding schemes.

The ISO/IEC 10918-1 standard, commonly referred to as JPEG [1], has become the most popular image compression tool nowadays. Several methods have been proposed to increase compression at the same image quality level, either based on machine learning techniques such as neural networks, or under the inclusion of perceptual knowledge in the coding steps.

Neural networks have been used and patented for image coding in **Patent No. 5005206** [2]. In the proposed scheme, the image is finally defined by the weights associated to a trained neural network in a particular linear domain. Another patent using neural networks for image compression is filed in **Patent No. 6798914** [3]. The method combines artificial intelligence and neural networks to convert digital image data into symbolic data which is further compressed with run-length encoding. In all these schemes, the problems of using neural networks are present; training instability, selection of a proper training algorithm, structure of the network, and the choice of the free parameters. Lately, the introduction of another machine learning algorithm, the support vector machine (SVM), working in a linear transformation (DCT) has been reported in **Patent No. W0/2003/050959** [4] for lossy data compression. The method is used in particular for compressing grayscale images. The introduction of SVM offers some advantages and control on the learning process: interpretability of the free parameters, sparsity in the application space, and lack of heuristics in training as the optimization problem is convex and the solution is unique. However, there is not a reported application of the method in colorful images or how the scheme could deal with them. In addition, the direct application of the standard SVM in the pre-specified (linear) working domain (the DCT domain) is not a good choice either. In [5], the latter problem is solved by introducing an *adaptive SVM* that takes into account the different perceptual relevance of the coefficients in the linear domain. Certainly, inclusion of perceptual (or statistical) *prior* knowledge is of paramount relevance for image coding. In this sense, DCTune is a technology for optimizing JPEG still image compression, and is covered by **Patent No. 5426512** [6]. DCTune calculates the best JPEG quantization matrices to achieve the maximum possible compression for a specified perceptual error, given a particular image and

a particular set of viewing conditions. This method is based on a point-wise non-linear model of masking (auto-masking), but it does not consider relations among coefficients (cross-masking) nor any machine learning method.

These methods produce good and competitive results compared to the existing standard JPEG method. However, in all patents and method descriptions, some limitations are observed: (i) there is not a single method that integrates perceptual knowledge and machine learning methods for colorful image compression so far; (ii) the existing methods do not give proper transformation schemes to deal with contrast in color images, or how could they be computed and introduced in the already patented schemes; (iii) the existing methods do not introduce non-linear perceptual transformations with interactions among coefficients before an eventual learning machine is applied; and (iv) the learning machine is not formulated to be independent on the domain of application. Alternatives to all these shortcomings are given by the presented method, apparatus and/or software.

3 Proposed alternatives

The invention presented hereafter consists of three essential blocks, which generalize the previous color image coding algorithms. Essentially, the scheme contains the following steps: (1) color contrast transforms, (2) non-linear divisive normalization perceptual transform, and (3) a general-purpose machine learning algorithm designed for the particular application domain at hand. A final quantization process is generally applied to the image description provided by the learning algorithm. These key processes are summarized in Fig. 1.

The method therefore extends the work in [7, 6, 8, 9, 10] by applying a non-linear perceptual transform that does include interactions among the contrast coefficients after the common linear transformation (e.g. DCT). Besides, it also uses a more general divisive normalization than JPEG2000 [11, 12] because it applies the fast recursive inversion method introduced in [13, 14] thus allowing a general (non-causal) interaction matrix eventually including interactions among chromatic channels. Moreover, the proposed method solves the problems found in [10, 13] since in those cases the color contrast preprocessing was not applied with the corresponding inaccuracy of the non-linear (point-wise or divisive normalization) transform. The method also differs from the neural-network based algorithms proposed so far in US Patents No. 5005206 [2], and No. 6798914 [3], since a general-purpose learning algorithm can be applied, which must be necessarily adapted to the relative perceptual or statistical relevance of the coefficients to be encoded in the working domain. In [5] the support vector regression (SVR) method was applied on the DCT domain, as in [15, 4]. However, given the special requirements of the domain, the algorithm was adapted to incorporate the perceptual relevance of the coefficients in this domain, which was not previously considered in [15, 4]. SVR was applied in a non-linear perceptual domain for achromatic image coding in [16], and due to the characteristics of the used (perceptually Euclidean) domain (as analyzed in [17]), a static SVR could be correctly applied. However, when working in other domains, the use of an adaptive SVR is necessary as illustrated in [5], and further analyzed in [17]. The presented method is aimed at defining a general purpose framework for color image coding, hence constituting a clear innovation in the color image compression field.

It is an object and goal of the present invention to provide a method, apparatus and/or software product/tool that yields improved performance in color image compression over existing methods, apparatus and software for comparative processing effort, or at least to provide the

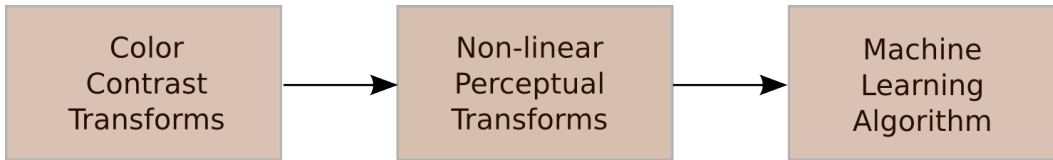


Figure 1: Key blocks used in the invention.

public with a useful alternative.

Further objects of the present invention may become apparent from the following description. Any discussion of the prior art throughout the specification provided here should in no way and by no means be considered as an admission that such prior art is widely known or forms part of common general knowledge in the field.

4 Summary and general overview of the invention

In this section, we present the method and analyze each of the blocks of the scheme for efficient coding of color images.

4.1 Scheme of the invention

A general scheme of the invention is illustrated in figure 2. The scheme for color images coding includes the following blocks:

Block 1 Initialization of parameters:

Block 1.A The user selects the desired rate (size of the coded image) or distortion (quality level of the coded image).

Block 1.B Initialization/modification of the learning and quantization parameters.

In the initial step, 1.B selects the parameters of the machine learning procedure and the quantization resolution to be applied to the description obtained by the learning process. To do so, the entropy or distortion choices in 1.A are considered. This can be done through a look-up table (LUT) relating free parameters and the obtained distortion/rate over an appropriate image set.

In further steps, 1.B modifies the preset parameters taking into account the deviation from the requested rate/distortion and the current values coming from block 16 and 17.

For example, if SVR is the selected learning procedure, 1.B would set/modify (1) the insensitivity values, ε_f , the penalization factor, λ , and the parameters of the kernel K (see block 9 and section 5.3); and (2) the number of bits used in the quantization to be applied to the weights associated to the selected support vectors (see block 10).

Block 2 Receiving the color image data expressed in Red-Green-Blue digital channels (hereafter RGB).

Block 3 Transforming the RGB signals to a color opponent space, that is, starting from the digital device-dependent color characterization, first obtain the color description in a device-independent tristimulus space (such as CIE XYZ). This requires the calibration parameters

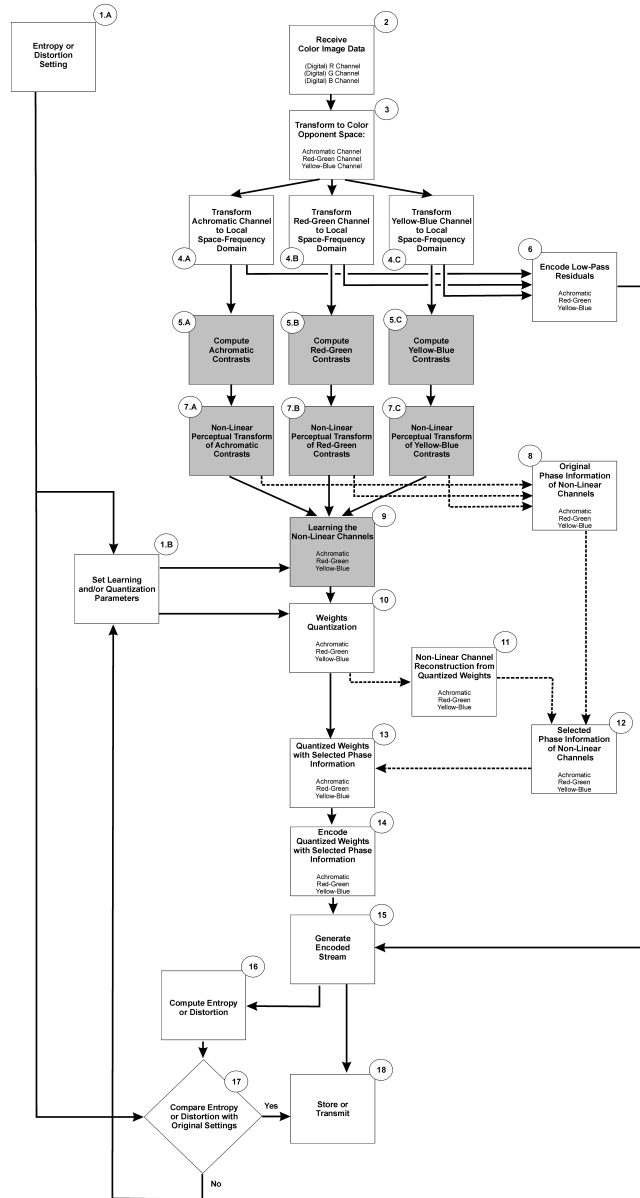


Figure 2: General scheme of the invention. Key blocks have been highlighted.

of the image acquisition device. If such parameters are not available, generic gamma curves and primaries may be used [18]. Second, the X, Y, Z images are combined into the following three channels: Achromatic (A), Red-Green (RG), Yellow-Blue (YB) channels according to some linear or non-linear opponent color transform [19].

Block 4 Transforming each opponent color channel (Achromatic 4.A, Red-Green 4.B, Yellow-Blue 4.C) to a local space-frequency domain. For example, using the block discrete cosine transform (DCT), wavelet transform or any suitable linear or non-linear transformation.

Block 5 Computing each channel contrast with appropriate (perceptually-based) transformations; that is, expressing the amplitudes of the local-frequency signals in 4.A, 4.B and 4.C in the appropriate contrast units to use the perception models in the corresponding steps 7.A, 7.B and 7.C.

In the achromatic case (5.A), this reduces to the standard Michelson contrast [20], but for each of the chromatic channels (5.B, 5.C) both psychophysical experiments with the chromatic basis functions of the transforms 4.B and 4.C, and numerical experiments on a representative color image database, are strictly necessary (see section 5.1).

Block 6 Encode the low-pass residuals of the A, RG and YB, channels from 4.A, 4.B and 4.C, respectively. This information may be DPCM and further entropy-coded.

Block 7 Apply a different non-linear perceptual transformation to each channel. In the case of the achromatic channel (7.A) and assuming a DCT in block 4.A, the non-linear transform introduced in [14, 17] can be applied. In the case of the chromatic channels (7.B, 7.C), different transforms have to be applied, see section 5.2 for their explicit forms.

Block 8 (Optional) Store the original phase information of the non-linearly transformed channels in 7.A, 7.B. and 7.C. This block is only necessary if phase information is independently encoded from amplitude information (related to blocks 11 and 12).

Block 9 Application and training of a machine learning procedure with parameters fixed in 1.B. This procedure selects the most representative and informative coefficients that accurately approximate each signal representation (7.A, 7.B, 7.C). In the case that the optional path leading to block 8 is followed, the learning procedure is applied to the absolute value of the signals. Otherwise, as the phase information is not processed separately, the learning procedure is applied directly to the signals coming from 7.A, 7.B and 7.C.

For example, but not restricted to, training a SVR method, and saving the associated weights and support vectors (see section 5.3).

Block 10 Quantization of the description obtained in block 9 using the quantization parameters set in block 1.B. In the case of using a SVR, a uniform quantization of the weights associated to the selected support vectors may be used.

Block 11 (Optional) Non-linear channel reconstruction from the quantized weights in the previous step 10. This block is only necessary if phase information is independently encoded from amplitude information (block 8). The result from this reconstruction is used to discard the unnecessary phase information coming from block 8 in block 12.

A proper procedure is necessary to invert the non-linear transforms 7.A, 7.B, 7.C. In the case of using divisive normalization transforms [21], a recursive procedure can be used [14].

Invertibility of the non-linearities in 7.A, 7.B and 7.C must be ensured. In [14], inversion of the transform in 7.A was demonstrated. See section 5.2 for a study on the inversion conditions necessary for the proposed transforms in blocks 7.B and 7.C to be valid.

Block 12 (Optional) Discard the phase information coming from optional block 8 according to the discarded coefficients in the reconstruction (optional block 11) after the quantization in block 10.

Block 13 Whenever amplitude and phase information have been separately processed, block 13 gathers the selected amplitude (coming from block 10) and phase (coming from block 12), and forwards the combined information to subsequent block 14. Otherwise, block 13 simply transmits the quantized weights with the corresponding phase information to block 14.

Block 14 Entropy coding of the signals. A variety of entropy coding schemes can be potentially applied [22].

Block 15 Generate encoded stream from the entropy-coded signals in blocks 6 and preceding block 14. The decoder (not shown here) should invert the whole process to reconstruct the image. As noted in block 11, the only non-trivial inversion process is associated to the non-linear transforms in blocks 7.A, 7.B and 7.C. See section 5.2 for a study on the invertibility of the proposed transforms.

Block 16 Compute entropy or distortion.

Block 17 Go to block 1.B unless the pre-specified requirements of either quality (distortion) or size (entropy) are achieved.

Block 18 Store or transmit the encoded color image.

4.2 Further capabilities

The application and use of the invention is not only restricted to the direct application of the previous sequential steps, but it preferably involves the following further considerations:

- Preferably, the method may further include subdividing the input image data received in step 1) into a number of parts, and performing the subsequent steps for each partition of the data.
- Preferably, the method may further consider sorting the input data in a proper way to make the whole process more efficient.
- Preferably, the step of developing the learning machine may consider exclusion of the DC component of the data, or each portion of it, and storing or transmitting the DC components as part of the compressed image.
- Preferably, the method may further include encoding the DC components before storing them as part of the compressed image.
- Preferably, adjustment of contrast and subsequent associated free parameters of the non-linear perceptual transform may involved alternative psychophysical experiments.

- Preferably, the step of training the approximation function may include any of the following (previous) steps: linear or non-linear feature selection or extraction under statistical and/or perceptual facts.
- Preferably, the machine learning procedure is trained to identify a sequence of basis functions and their corresponding weights, guided with proper criteria to optimize the compactness of the resulting encoded data sequence and/or the quality of the encoded image.
- Preferably, the method may further discard weights calculated by the learning method by eliminating those weights without associated basis function, disregard of the proper application of specific filters to the processed approximated signals. Also, the method may preferably further include heuristic rules to discard weights. Preferably, the method may further include discarding the information for any coefficient overpassing any meaningful perceptual or statistical magnitude.
- Preferably, the method may include a system to correct the bias in the approximation so the approximated values are all non-negative. These information may be eventually combined with the associated phase information for signal value coefficients (its sign) in order to properly handle the weights and/or associated basis functions.
- Preferably, the method may include encoding the data stream using entropy coding under any variant.
- Preferably, the step of linear transformation of the data may include the DCT or any form of Gabor-based or wavelet discrete transform.
- Preferably, different learning machines may be applied, either neural networks, spline-based techniques, kernel-based methods, Bayesian networks, Gaussian Processes, or fuzzy logic.
- Preferably, the chromatic transformations can be adapted *per* image, and the contrasts given by the transforms ultimately scaled (or further transformed) to properly fit the subsequent non-linear perceptual transform ranges.
- Preferably, the non-linear perceptual transform may be adapted to user specifications, or its information complemented by other non-linear perceptual or statistical transforms.
- Preferably, the quantization step may include further complementary adaptive coding schemes.

Any of the previous options can be combined to ultimately lead to different functional modes, disregard of the completeness of the general proposed scheme.

5 Detailed description to carry out the invention

In this section, we describe in detail possible implementations of the the key blocks in the general scheme (namely 5.X, 7.X and 9) to carry out the invention.

5.1 Contrast transforms

The first key element is an appropriate contrast transform of the local frequency representation of the A, RG, and YB channels. This is a necessary preprocessing before the non-linear perceptual transforms are applied: in order to apply psychophysically inspired parameters in the non-linear transforms in blocks 7.A, 7.B and 7.C, the input signal has to be expressed in the appropriate units. This section (and the blocks 5.A, 5.B, 5.C) are devoted to express the signal in the appropriate *contrast* units.

The psychophysical literature [23, 24, 25] describes the non-linear response of human visual sensors to oscillating input signals expressed in contrast units. Moreover, it has been found that the shape of the non-linear response for the local-frequency chromatic sensors is similar to the achromatic case [25], when expressing the inputs in the appropriate chromatic contrast units. However, there is no general procedure to define such units for particular spatial basis functions (such as block-DCT or particular wavelets) with chromatic modulation. Therefore, expressing the amplitude of coefficients in contrast units is the key to (1) applying the relative scaling of the parameters of the different achromatic and chromatic channels (as for instance the relative frequency sensitivity [26, 27]), and (2) designing the parameters that control the shape of the chromatic non-linearities from the achromatic one. The contrast transforms in blocks 5.A, 5.B and 5.C are the key to simplify the formulation of the blocks 7.A, 7.B and 7.C.

The definition of contrast in the psychophysical literature is closely related to the particular procedure to construct the stimuli used in the experiments. For instance, chromatic contrast is defined by the luminance modulation of two basis functions of extreme colors, $\mathbf{e}_1^{(i)}$ and $\mathbf{e}_2^{(i)}$, added in counter phase to generate gratings of perceptually uniform brightness [26, 27]. The maximum luminance and color variations (the maximum amplitude of the stimuli, i.e. defining what unit contrast is) is limited by the available color gamut of the particular color reproduction device used in the experiment.

In order to apply the same kind of contrast definition, we have to (1) obtain the equations that relate the amplitudes of local-frequency coefficients and the maximum modulation in color space with the luminance and extreme colors of the equivalent gratings with uniform brightness, and (2) ensure that the extreme colors computed according to the above equations in a big enough color image database are inside the gamut of reproducible colors in the conventional color reproduction devices. The idea is defining the unit of contrast by using the extreme luminance and color variations found in natural images and being consistent with the definition of color contrast in the psychophysical literature.

According to this, the achromatic and chromatic contrast definition in blocks 5.A, 5.B and 5.C is based in two steps:

- First, we obtain the equations that relate achromatic and chromatic contrasts to the amplitudes of local-frequency transform of the achromatic and chromatic channels (next subsection 5.1.1). These equations (for the chromatic channel, i) depend on some extreme colors, $\mathbf{e}_1^{(i)}$ and $\mathbf{e}_2^{(i)}$, obtained from the maximum color deviation in that chromatic direction, $\Delta\mathbf{T}_{\max}^{(i)}$. Besides, when using linear color models, the average luminance of the colors to be mixed to generate uniform brightness gratings is not the same because colors of different chromaticity have different brightness for the same luminance [19]. This implies that a psychophysical experiment is needed to find the right luminance proportion in the mixture. Subsection 5.1.2 describes the procedure and the results of such an experiment in the particular case of the DCT basis functions used in a possible implementation of the invention.

- Second, we study the limits of the maximum deviations, $\Delta \mathbf{T}_{\max}^{(i)}$, in a representative color image database. From the empirical study in subsection 5.1.3, we estimate the maximum deviations that give rise to unit contrast gratings consistent with the assumptions made in the psychophysical literature. These results are examples for the particular DCT choice and the linear YUV choice.

Finally subsection 5.1.4 summarizes the resulting transforms for the appropriate contrast definition used in blocks 5.A, 5.B and 5.C, when DCT and linear YUV are the selected image and color representations.

5.1.1 Formulation: Contrast of achromatic and chromatic gratings from the transform coefficients of images in opponent color spaces

The gratings (particular spatial basis functions) of frequency f used in psychophysical experiments have a different color, described by a 3D vector \mathbf{T} , at every spatial position, x :

$$\mathbf{T}(x) = \mathbf{T}_0 + \Delta \mathbf{T}(f) \cdot B(f, x) \quad (1)$$

where \mathbf{T}_0 is the average color, $\Delta \mathbf{T}(f)$ represents the maximum (peak) chromatic oscillations of that frequency with respect the average color, and $B(f, x)$ is the selected basis function of frequency f , such as those in the DCT. The index, f , identifying the basis function may be more general than the frequency: in the wavelet case, f , would have frequency and space meaning. In eq.1, we added the index, f , to the color deviation, $\Delta \mathbf{T}(f)$, because in images containing more than one basis functions, the amplitude of the color modulation is different for each basis function, $B(f, x)$.

In a linear ‘opponent color space’, such as for instance, the linear YUV color space [28], the three components of vectors $\mathbf{T} = [T_1, T_2, T_3]^\top$, stand for the luminance, $T_1 = Y$, the Yellow-Blue component, $T_2 = U$, and the Red-Green component, $T_3 = V$. Psychophysical experiments use the above gratings to stimulate particular achromatic or chromatic sensors by using particular modulations in the components of $\Delta \mathbf{T}(f)$. The modulations, $\Delta T_i(f)$, are easily related to the coefficients of the discrete spatial transform of corresponding image (color channel) in eq. 1 (see eqs. 20 below). The problem is that the chromatic contrasts in the Yellow-Blue and the Red-Green directions, as they are defined in the psychophysical literature [26, 27], are not trivially related to $\Delta T_2(f)$ and $\Delta T_3(f)$. For the shake of clarity, in the following discussion we will temporarily omit the index f from the color modulations $\Delta T_i(f)$.

Lets start from the easy (achromatic) case. In this case, no modulation is used in the chromatic channels ($\Delta T_2 = \Delta T_3 = 0$), so the grating that isolates the achromatic channel is:

$$\mathbf{T}^{(1)}(x) = \begin{pmatrix} T_{01} \\ T_{02} \\ T_{03} \end{pmatrix} + \begin{pmatrix} \Delta T_1 \\ 0 \\ 0 \end{pmatrix} \cdot B(f, x) \quad (2)$$

and the achromatic contrast is simply defined by the Michelson’s contrast of the achromatic (only luminance) grating:

$$C_{achrom} = \frac{\Delta T_1}{T_{01}} \quad (3)$$

where T_{01} is the average luminance (or the luminance of the average color), and the amplitude ΔT_1 can be obtained from the corresponding local-frequency coefficients of the luminance channel using the equations 20. This definition can be extended to more complex wavelet basis by

dividing each coefficient in each subband by the corresponding low-pass residual at that resolution [20]. As a result, the achromatic contrast is in the range $C_{achrom} \in [0, 1]$ since $\Delta T_1 \in [0, T_{01}]$.

However, in [26, 27], the definition of the chromatic contrast takes into account the way in which chromatic gratings are experimentally designed. In order to generate a purely chromatic modulation around an average color, \mathbf{T}_0 , in a given direction, T_i , with $i = 2, 3$, two luminance gratings with extreme chromaticity from the unit-luminance colors $\mathbf{e}_1^{(i)}$ and $\mathbf{e}_2^{(i)}$, are added in counter phase (see the top row of fig. 3):

$$\mathbf{T}^{(i)}(x) = \mathbf{e}_1^{(i)} \left(\frac{T_{01}}{2}(1 + \eta_i) + \Delta Y_1^{(i)} \cdot B(f, x) \right) + \mathbf{e}_2^{(i)} \left(\frac{T_{01}}{2}(1 - \eta_i) - \Delta Y_2^{(i)} \cdot B(f, x) \right) \quad (4)$$

where the luminance modulations $\Delta Y_j^{(i)}$ are in the ranges: $\Delta Y_1^{(i)} \in [0, \frac{T_{01}}{2}(1 + \eta_i)]$ and $\Delta Y_2^{(i)} \in [0, \frac{T_{01}}{2}(1 - \eta_i)]$ with the following constraints:

- The average color is obtained from the sum of the extreme colors $\mathbf{e}_j^{(i)}$ with their luminance scaled in a particular proportion (given by the factor η_i):

$$\mathbf{T}_0 = \mathbf{e}_1^{(i)} \frac{T_{01}}{2}(1 + \eta_i) + \mathbf{e}_2^{(i)} \frac{T_{01}}{2}(1 - \eta_i) \quad (5)$$

- The luminance modulations, $\Delta Y_j^{(i)}$, are coupled:

$$\Delta Y_2^{(i)} = \frac{(1 - \eta_i)}{(1 + \eta_i)} \Delta Y_1^{(i)} \quad (6)$$

- The overall deviation in the chromatic direction i , $\Delta \mathbf{T}^{(i)}$ induces no modulation in the orthogonal chromatic channel, i.e.:

$$\Delta \mathbf{T}^{(2)} = \begin{pmatrix} \eta_2 T_{01} \\ \Delta T_2 \\ 0 \end{pmatrix} \quad (7)$$

or

$$\Delta \mathbf{T}^{(3)} = \begin{pmatrix} \eta_3 T_{01} \\ 0 \\ \Delta T_3 \end{pmatrix} \quad (8)$$

Note that given the difference between *brightness* and *luminance* [19], in order to obtain iso-brightness gratings, the luminance of the two gratings of different color have different average luminance (thus inducing a residual luminance modulation $\eta_i T_{01}$). The deviation from the average luminance (the factor η_i) has to be experimentally determined for each chromatic channel. Section 5.1.2 is devoted to estimate this factor for DCT basis functions of specific frequency and average color.

In the psychophysics literature, the chromatic contrast for channel i is defined as the Michelson's contrast of the two counter-phase luminance gratings giving rise to the purely chromatic grating:

$$C_{chrom}^{(i)} = \frac{\Delta Y_1^{(i)}}{\frac{T_{01}}{2}(1 + \eta_i)} = \frac{\Delta Y_2^{(i)}}{\frac{T_{01}}{2}(1 - \eta_i)} \quad (9)$$

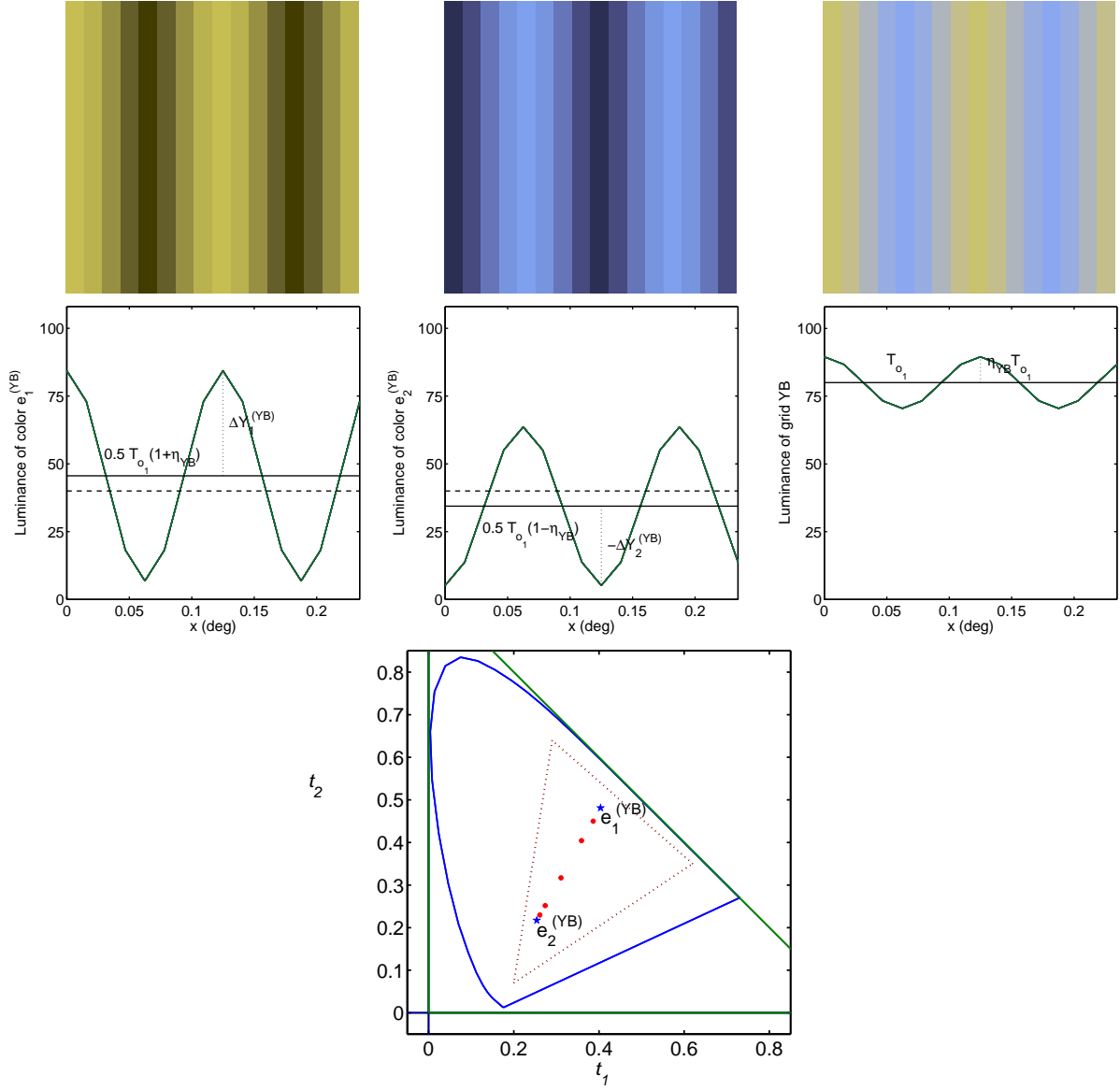


Figure 3: Example of a purely chromatic grating in the Yellow-Blue chromatic direction. In this case, a 16×16 DCT basis function of frequency $f_x = 8$ and $f_y = 0$ in cycl/deg is considered. The average color is $\mathbf{T}_0 = [80 \ 80 \ 80]^T$ in RGB NTSC or $\mathbf{T}_0 = [80 \ 0 \ 0]^T$ in linear YUV, the chromatic contrast is $C_{chrom}^{(2)} = 0.85$ assuming a maximum modulation in the U channel $\Delta T_{2max} = 45$. The top row shows the two counter-phase gratings of extreme chromaticity $\mathbf{e}_1^{(2)}$ (yellow at the left), and $\mathbf{e}_2^{(2)}$ (blue at the center), to generate the purely chromatic grating (at the right). Average luminance values of the yellow and blue gratings have been chosen to obtain constant brightness in the final grating ($\eta_{YB} = 0.14$, see section 5.1.2), and it isolates the response in the Yellow-Blue channel ($\Delta T_2 \neq 0$ and $\Delta T_3 = 0$). The center row shows the luminance of the gratings in the top row. The bottom plot shows the CIE xy chromatic coordinates of the extreme (unit luminance) colors $\mathbf{e}_j^{(2)}$ with $j = 1, 2$ (stars) and the actual colors in the grating (circles). Increasing the chromatic contrast means simultaneously increasing the luminance amplitudes $\Delta Y_j^{(2)}$, thus increasing the chromaticity range covered by the circles in the bottom plot. In the extreme (unit contrast) case, the chromatic range would be exactly determined by the extreme colors $\mathbf{e}_j^{(2)}$. Reducing the chromatic contrast means reducing the luminance amplitudes ΔY_j and reducing the chromatic range. In the zero contrast limit the only color of the grid would be the average color \mathbf{T}_0 (the central circle in the bottom diagram).

From Eqs. (1), (4) and (9), a relation between the color modulation in the i -th chromatic direction, $\Delta\mathbf{T}^{(i)}$, and the corresponding chromatic contrast can be derived:

$$\Delta\mathbf{T}^{(i)} = \left(\mathbf{e}_1^{(i)} - \frac{1 - \eta_i}{1 + \eta_i} \mathbf{e}_2^{(i)} \right) \frac{T_{01}}{2} (1 + \eta_i) C_{chrom}^{(i)} \quad (10)$$

The unit contrast case (taking $C_{chrom}^{(i)} = 1$) gives a relation between what is considered the maximum modulation in this channel, $\Delta\mathbf{T}_{\max}^{(i)}$, and the extreme colors, $\mathbf{e}_j^{(i)}$:

$$\mathbf{e}_1^{(i)} = \frac{1}{T_{01}(1 + \eta_i)} \left(\Delta\mathbf{T}_{\max}^{(i)} + \mathbf{T}_0 \right) \quad (11)$$

$$\mathbf{e}_2^{(i)} = \frac{2}{T_{01}(1 - \eta_i)} \left(\mathbf{T}_0 - \frac{T_{01}}{2} (1 + \eta_i) \mathbf{e}_1^{(i)} \right) \quad (12)$$

Psychophysical experimentation assumes maximum modulation in such a way that the extreme colors are inside the available gamut of the display at hand. According to such criterion, in section 5.1.3 we will explore a representative color image data base to estimate the maximum possible modulation around the average color in each image in such a way that the extreme colors are always inside the gamut of standard displays (e.g. the dashed triangle of the chromaticity diagram in fig. 3).

From eq. 10, it can be obtained the chromatic contrast $C_{chrom}^{(i)}$ from the modulation ΔT_i

$$C_{chrom}^{(i)} = \frac{\Delta T_i}{\left(\mathbf{e}_1^{(i)} - \frac{1 - \eta_i}{1 + \eta_i} \mathbf{e}_2^{(i)} \right)_i \frac{T_{01}}{2} (1 + \eta_i)} = \frac{\Delta T_i}{(\Delta\mathbf{T}_{\max}^{(i)})_i} \quad (13)$$

Given the achromatic and chromatic contrasts as a function of the modulation in each chromatic opponent channel (eqs. 3 and 13), the last piece of information necessary to express the amplitudes of the local-frequency transforms into contrasts is the relationship between the coefficients of the discrete transform and the corresponding modulations.

In the case of the DCT, this relation can be found elsewhere [29]. Here, the indices x , f used before are explicitly declared as $x = (m, n)$ and $f = (p, q)$. With this notation, the 2D-DCT coefficients, $a(p, q)$, of the spatial domain representation of an $M \times N$ image block, $A(m, n)$, are [29]:

$$a(p, q) = \beta_p \beta_q \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} A(m, n) \cos\left(\frac{\pi(2m+1)p}{2M}\right) \cos\left(\frac{\pi(2n+1)q}{2N}\right) \quad (14)$$

where $\beta_p = 1/\sqrt{M}$ for $p = 0$ and $\beta_p = \sqrt{2}/\sqrt{M}$ for $p > 0$, $\beta_q = 1/\sqrt{N}$ for $q = 0$ and $\beta_q = \sqrt{2}/\sqrt{N}$ for $q > 0$.

The inverse of the DCT is written as:

$$A(m, n) = \sum_{p=0}^{M-1} \sum_{q=0}^{N-1} \beta_p \beta_q a(p, q) \cos\left(\frac{\pi(2m+1)p}{2M}\right) \cos\left(\frac{\pi(2n+1)q}{2N}\right) \quad (15)$$

If the acquired image is a simple basis function (in the chromatic channel i) defined as:

$$A_i(m, n) = T_{0_i} + \Delta T_i(p, q) \cos\left(\frac{\pi(2m+1)p}{2M}\right) \cos\left(\frac{\pi(2n+1)q}{2N}\right), \quad (16)$$

one can write the appropriate relation between the coefficients of the image channel i , $a_i(p, q)$ and the actual chromatic modulation in the corresponding frequency $\Delta T_i(p, q)$:

$$T_{0_i} = \beta_o \beta_o a_i(o, o) = \frac{1}{\sqrt{MN}} a_i(o, o) \quad (17)$$

$$\Delta T_i(o, q) = \beta_o \beta_q a_i(o, q) = \frac{\sqrt{2}}{\sqrt{MN}} a_i(o, q) \quad (18)$$

$$\Delta T_i(p, o) = \beta_p \beta_o a_i(p, o) = \frac{\sqrt{2}}{\sqrt{MN}} a_i(p, o) \quad (19)$$

$$\Delta T_i(p, q) = \beta_p \beta_q a_i(p, q) = \frac{2}{\sqrt{MN}} a_i(p, q) \quad (20)$$

5.1.2 Experiment 1: Blue-Yellow and Red-Green luminance proportions in uniform brightness gratings

In this section, we describe the psychophysical experiments carried out to estimate the right proportions of average luminance of the chromatic gratings used to generate purely chromatic gratings in the U and V directions with uniform brightness, i.e. the factors η_i in the above equations.

Stimuli. Stimuli were generated with computer controlled CRT monitors and 8-bit resolution graphic cards. The systems were colorimetrically calibrated and controlled using Matlab[®], with the library COLORLAB [18].

Colors $\mathbf{e}_j^{(i)}$ were chosen along the directions of color space nulling one chromatic channel in YUV space, in such a way that:

- They caused modulations differing in sign only along the i direction from the average color \mathbf{T}_0 . The selected average color was taken from the average color of a public color image data base [30]. In the NTSC RGB color space, we have: $T_0 = [132 \ 116 \ 90]^\top$, with a luminance of 117.8 cd/m^2 .
- The modulation induced was as large as the color gamut of the monitor could allow. Two luminance DCT gratings, of 8 cpd, each with the chromaticity of one of these colors, and initially with the same mean luminance, were generated in counter phase, and added. The resulting image, subtending 0.25 degrees, was presented against a grey background of 50 cd/m^2 in an otherwise dark room.

Prior to the actual measurements and by mean of successive adjustments, the mean luminance of the achromatic point, \mathbf{T}_0 was modified in such a way that:

- Luminance was inside the range of possible values for the monitor used for all pixels in the resulting image.
- A sufficient luminance range remained for the observer to increase the luminance of one of the gratings constituting the final image to reach the constant brightness condition.

With these restrains, the chosen luminance value for the average color was 80 cd/m^2 .

Measurement. Observers, adapted for one minute to the grey background, were shown the images described above and asked to fixate them foveally. The task of the observer was to adjust the variable η_i , in such a way that the two hemiperiods of the resulting grating, appeared to have the same brightness, For this, the method of adjustment (MOA) was used [31]. For each observer, the result of the experiment was the mean of five trials. The final result was obtained as the average of a set of 3 naïf observers, in the [25, 38] age range.

Results. The proportion factors of Yellow-Blue and Red-Green proportions found in the above experiment were:

$$\eta_{YB} = \eta_2 = \eta_U = 0.14 \pm 0.06 \quad (21)$$

$$\eta_{RG} = \eta_3 = \eta_V = 0.02 \pm 0.07 \quad (22)$$

Since the values η_i depend on (1) the chromaticity of the extreme colors $\mathbf{e}_j^{(i)}$, (2) the frequency of the grating, and (3) the spatial nature of the basis function, the above values are strictly valid only for the described conditions. Using different basis functions (e.g. using wavelets instead of block DCT) or different opponent color spaces would require specific experimentation using the above procedure. In the current implementation of the invention, we assume that the proportions found in the above conditions are approximately valid for different average colors and spatial basis functions.

It has to be noted that using non-linear color models with better isolation of the achromatic perception in the corresponding achromatic channel may alleviate the need of this experimental step in order to obtain purely chromatic gratings.

5.1.3 Experiment 2: Maximum modulation in opponent spaces in natural images

In this section, we describe the numerical experiment used to estimate reasonable values for the maximum color modulation for natural images in the chromatic directions of the linear YUV space. The values $(\Delta \mathbf{T}_{\max}^{(i)})_i$ are necessary for the color contrast definition in eq. 13. In order to be consistent with the experimental constraints applied in the psychophysical literature, a particular choice of maximum modulation values has to fulfill the following condition: for a typical natural image with average color \mathbf{T}_0 , the extreme colors in eqs. 12 have to be inside the color gamut of a typical color display.

Experiment. In order to find these values, we started from an arbitrarily large initial guess, e.g. $(\Delta \mathbf{T}_{\max}^{(i)})_i = 256$, and checked the above condition for the $\mathbf{e}_j^{(YB)}$ and $\mathbf{e}_j^{(RG)}$ colors (computed by using eq. 12 with the experimental proportions 22) over a representative set of natural color images (see next section for the details on the database). If, for some image, one of the computed extreme colors was outside the color gamut of a typical CRT monitor, the corresponding maximum modulation, ΔU_{max} or ΔV_{max} , was reduced by a 5%. The search procedure stopped when the condition is satisfied for every image in the database.

Due to the limitations of the database (e.g. limited range of average luminance), the outcome of the above procedure is just a set of recommended values. The recommended values do not guarantee a maximum bound for the chromatic contrast: particular images may exist giving rise to chromatic contrasts larger than 1 specially in low luminance and/or high saturation conditions. For such images larger maximum modulation values may give rise to more convenient contrast values.

Database. We ran the numerical experiment on a database consisting of 100 images from the McGill University color image database [30] (mainly wildlife images) and 25 raw images taken in our lab (mainly human faces). In the latter case, the Canon EOS 20 D camera used in taking the pictures was calibrated with the COLORLAB library [18] and a PhotoResearch SpectraScan PR650 spectroradiometer.

Results. Figure 4 shows two sample images limiting the maximum modulation in the V and U directions respectively. The CIE xy chromatic diagrams below show the colors of the corresponding image (in gray), the average color (in black), and the extreme colors (in red and blue). The dotted triangle in the chromatic diagrams represents the gamut of available colors in a typical CRT monitor (similar to NTSC or PAL RGB primaries). As seen in the chromatic diagrams, the left image limits the modulation in the V (Red-Green) direction while the right image limits the modulation in the U (Yellow-Blue) direction.

The maximum modulations compatible with the psychophysical assumptions in the analyzed database were:

$$\Delta U_{max} = 30 \pm 2 \quad (23)$$

$$\Delta V_{max} = 43 \pm 2 \quad (24)$$

The associated error, ± 2 , corresponds to the 5% step in the search procedure. Besides, due to the above considerations on the limitations of the database, these values are just a convenient recommendation that may be improved by introducing some sort of dependence on the average luminance in more refined implementations of the invention.

5.1.4 Summary: achromatic and chromatic contrast definition

According to the general scheme of the invention (see fig. 2), the set of signal transforms carried out by blocks 3, 4.X and 5.X can be summarized as follows:

$$\mathbf{A}'(x) \xrightarrow{\text{Chrom.Transform.}} \mathbf{A}(x) \xrightarrow{\text{Space-Freq.Transform.}} \mathbf{a}(f) \xrightarrow{\text{Contrast.Transform.}} \mathbf{c}(f) \quad (25)$$

where the components of vectors, $\mathbf{A}'(x)$, are the tristimulus values of the color in the pixel x in the initial (non-opponent) color space (e.g. NTSC RGB); the vectors, $\mathbf{A}(x)$, obtained from the chromatic transform in block 3, are the corresponding colors in an opponent representation (such as linear YUV space); the elements of the vectors, $\mathbf{a}(f)$, are the coefficients of the local space-frequency transforms applied on the images $A_i(x)$, i.e., $a_1(f)$ stands for the coefficients of the transform of the achromatic image, $a_2(f)$ stands for the coefficients of the transform of the Blue-Yellow image, and $a_3(f)$ stands for the coefficients of the transform of the Red-Green image (obtained using the blocks 4.A, 4.B and 4.C respectively). Finally, the contrast transform expresses each amplitude, $a_i(f)$, in achromatic or chromatic contrast units, giving the vectors $\mathbf{c}(f)$. The components of the contrast vector are: $c_1(f) = C_{achrom}(f)$, $c_2(f) = C_{chrom}^{(YB)}(f)$, and $c_3(f) = C_{chrom}^{(RG)}(f)$.

According to the formulation and experimental results presented in this section, if the selected chromatic representation is linear YUV, and the selected spatial representation is block DCT with $M \times N$ block size, the recommended contrast transformations are listed below:

- For the luminance (Y) channel (block 5.A),

$$C_{achrom}(p, o) = \sqrt{2} \frac{a_1(p, o)}{a_1(o, o)} \quad (26)$$

$$C_{achrom}(o, q) = \sqrt{2} \frac{a_1(o, q)}{a_1(o, o)} \quad (27)$$

$$C_{achrom}(p, q) = 2 \frac{a_1(p, q)}{a_1(o, o)} \quad (28)$$

- For the Yellow-Blue (U) channel (block 5.C),

$$C_{chrom}^{(2)}(p, o) = \frac{\sqrt{2}}{\sqrt{MN}} \frac{a_2(p, o)}{30} \quad (29)$$

$$C_{chrom}^{(2)}(o, q) = \frac{\sqrt{2}}{\sqrt{MN}} \frac{a_2(o, q)}{30} \quad (30)$$

$$C_{chrom}^{(2)}(p, q) = \frac{2}{\sqrt{MN}} \frac{a_2(p, q)}{30} \quad (31)$$

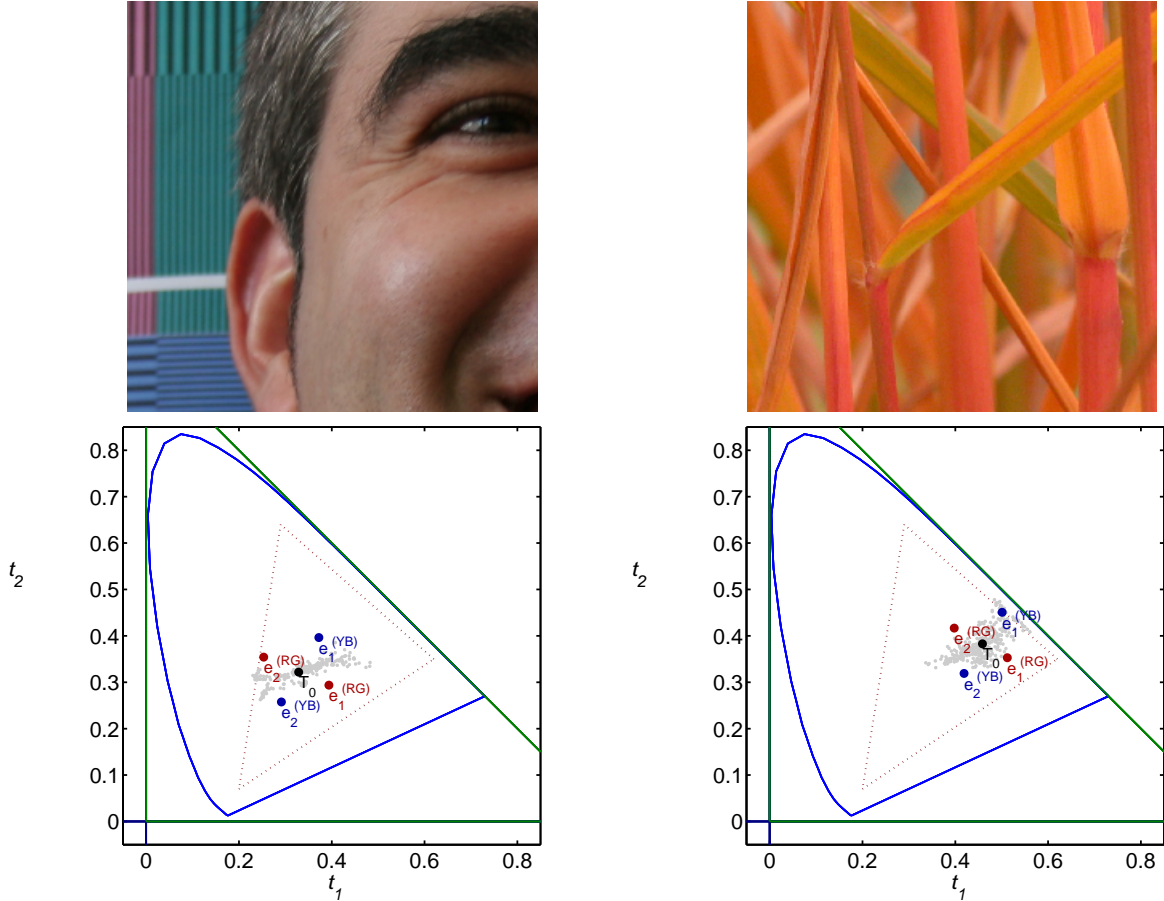


Figure 4: Sample images limiting the maximum modulation in the V and U directions (top) and the corresponding colors (bottom).

- For the Red-Green (V) channel (block 5.B),

$$C_{chrom}^{(3)}(p, o) = \frac{\sqrt{2}}{\sqrt{MN}} \frac{a_3(p, o)}{43} \quad (32)$$

$$C_{chrom}^{(3)}(o, q) = \frac{\sqrt{2}}{\sqrt{MN}} \frac{a_3(o, q)}{43} \quad (33)$$

$$C_{chrom}^{(3)}(p, q) = \frac{2}{\sqrt{MN}} \frac{a_3(p, q)}{43} \quad (34)$$

5.2 Non-linear perceptual transforms

This section is devoted to the description of the non-linear perceptual transforms applied to the local-frequency transforms expressed in achromatic and chromatic contrast units.

5.2.1 Divisive normalization transforms (Achromatic, Yellow-Blue and Red-Green)

The second key process in the present invention deals with the application of a non-linear perceptual transformation to the contrast of the achromatic and chromatic channels. The roots of this transform are motivated by perceptual (contrast gain control and masking experiments) and statistical (efficient coding) facts of image perception.

The whole perceptual transformation is modeled as a three step process: first a set of linear filters are applied to each channel (this step corresponds to blocks 4A, 4B and 4C in Fig 2); next, the contrast of each coefficient is obtained using the procedure described in the previous section (blocks 5A, 5B and 5C in Fig 2); and, finally a non-linear transform is applied to the output of the second stage (blocks 7A, 7B and 7C in Fig 2). This last transform can be a point-wise non-linearity or preferably a non-linear transform [21, 24]. In the later case, the energy of each linear coefficient is normalized by a combination of the energies of neighboring coefficients in frequency, thus for each channel the response at a particular frequency is obtained as follows:

$$r_i(f) = \frac{\text{sgn}(c_i(f)) \cdot |s_i(f) \cdot c_i(f)|^\gamma}{\theta(f) + \sum_{f'=1}^{N^2} h(f, f') |s_i(f') \cdot c_i(f')|^\gamma} \quad (35)$$

where subscript $i = 1, 2, 3$ denotes the channel; $c_i(f)$ are the outputs of the local-frequency analyzers in contrast units; $s_i(f)$ are CSF-shaped functions; γ is an exponent; $\theta(f)$ is a regularizing function; and, $h(f, f')$ determines the interaction neighborhood among coefficients in the non-linear normalization of the energy.

$$h(f, f') \propto \exp\left(-\frac{\|f - f'\|^2}{\sigma_{|f|}^2}\right) \quad (36)$$

where $\sigma_{|f|} = \frac{1}{6}|f| + 0.05$, $|f|$ is given in cycles per degree (cpd), and the Gaussian is normalized to have unit volume.

The above interaction kernel is just an example including intra-block frequency relations, but it could be generalized to include spatial and chromatic interactions.

This general divisive normalization model is well-established for the achromatic case based on achromatic contrast incremental thresholds of sinusoidal or Gabor patches [24] and equivalent physiological experiments [21]. This achromatic model has been successfully applied in several

image processing applications including achromatic image coding [14, 17] and achromatic image denoising [16].

In the chromatic case, similar non-linear behavior has been reported with sinusoidal and Gabor patches [25, 32]. This is why in the proposed implementation of the invention, the parameters controlling the non-linearity in the chromatic case (γ , θ and h) are proposed to be equal to those used in the achromatic case. However, the overall frequency sensitivities are adapted for each chromatic channel, based on previous linear-models work [26, 27].

Since all the psychophysical and physiological results are based on measures using sinusoidal or Gabor patches, the experimental parameters have to be adapted to the particular image representation domain in encoding applications (e.g. DCT, Wavelets, ICA, etc.). If the image representation choice is block-DCT, we recommend to use the particular values illustrated in figs. 5, 6 and 7 following the comments of Dr. Uriegas [32], Dr. Watson [33] and Dr. Heeger [34].

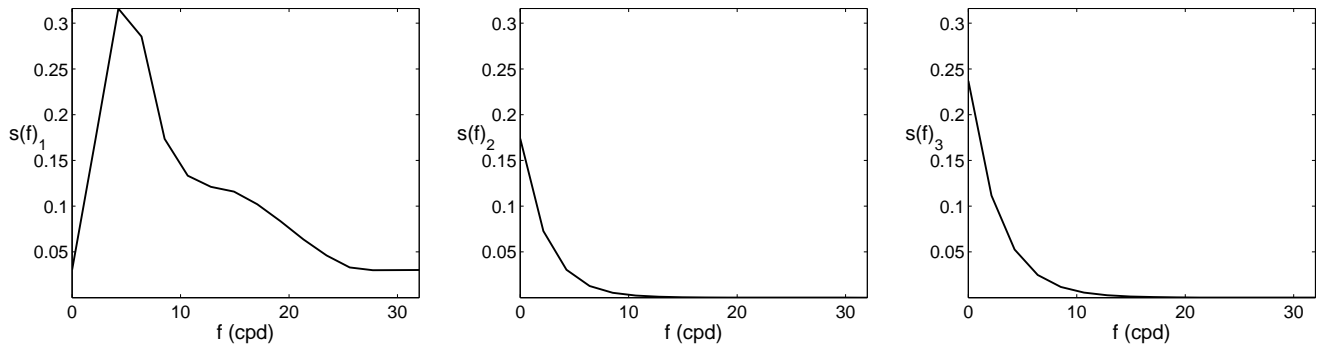


Figure 5: Parameters $s_i(f)$ for each channel in Eq. (35)

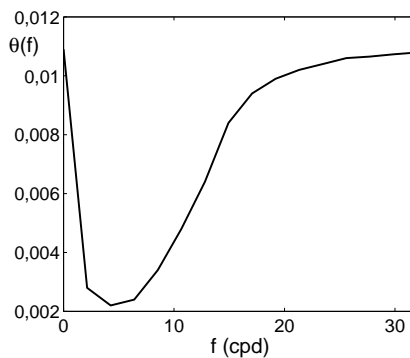


Figure 6: Parameter $\theta(f)$ in Eq. (35)

The response in the achromatic channel has been computed under two input conditions in order to illustrate how the model in eq. 35 accounts for different perceptual facts. First, the response, $r_1(f)$, of a sensor tuned to a frequency f , is computed when the input is $\mathbf{c}_1 = [\mathbf{0}, c_1(f), \mathbf{0}]$ (fig. 8 shows the results for two particular sensors: $r_1(f = 4)$ and $r_1(f = 10)$). Second, the response of a sensor, $r_1(f)$, has been obtained in the presence of an additional pattern with a different spatial frequency, i.e., when the input is $\mathbf{c}_1 = [\mathbf{0}, c_1(f), \mathbf{0}, c_1(f'), \mathbf{0}]$ (fig. 9 shows the results obtained for sensors $r_1(f = 4)$ and $r_1(f = 10)$ when the inputs are $\mathbf{c}_1 = [\mathbf{0}, c_1(4), \mathbf{0}, c_1(6), \mathbf{0}]$ and $\mathbf{c}_1 = [\mathbf{0}, c_1(6), \mathbf{0}, c_1(10), \mathbf{0}]$ respectively). The perceptual facts reproduced are:

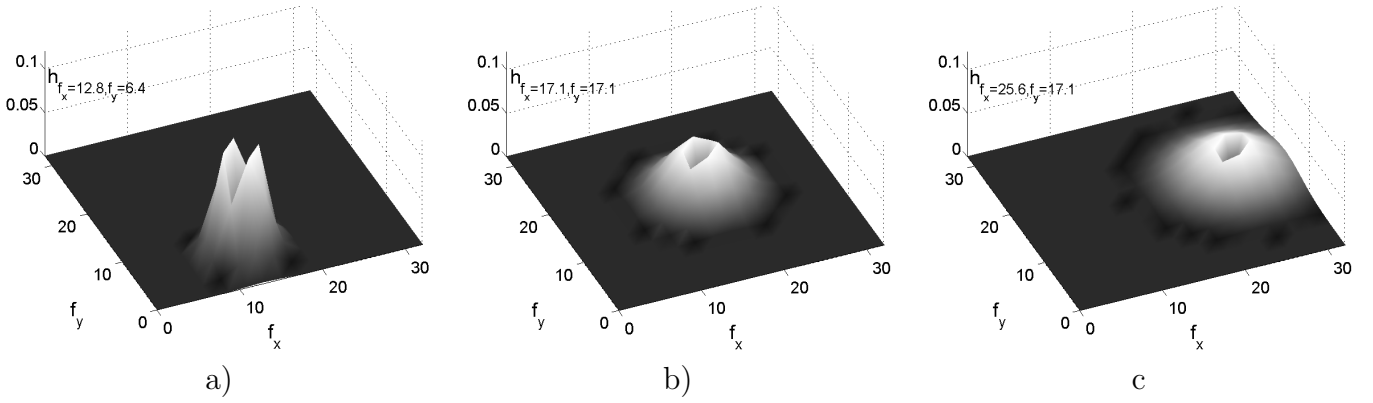


Figure 7: Three illustrative frequency interaction neighborhoods (rows of $h(f, f')$) in Eq. (35). Each surface corresponds to a different frequency.

- *Frequency selectivity*: the visibility of distortions depends on the spatial frequency. Note that in figure 8, the slope of the response curve is larger for 4 cpd than for 10 cpd, thus a larger amount of distortion is required in 10 cpd to obtain the same perceived distortion. In other words, 4 cpd noise is more visible than 10 cpd noise of the same energy. This general behavior is given by the band-pass function $s_1(f)$ (see figure 5).
- *Auto-masking*: the amount of distortion needed to obtain a constant perceptual distortion increases with the contrast of the input. See in figure 8 how Δc increases with the contrast of the stimulus. This is due to the fact that the response is attenuated when increasing the contrast because of the normalization term in the denominator of Eq. (35).
- *Cross-masking*: the attenuation (and the corresponding response saturation and sensitivity decrease) also occurs when other patterns $c_1(f')$ with $f' \neq f$ are present. Note in figure 9 that the required amount of distortion increases as the contrast of the mask of different frequency is increased. Moreover, given the Gaussian shape of the interaction neighborhood, patterns of closer frequencies mask the distortion more effectively than background patterns of very different frequency. Accordingly, the 6 cpd mask induces a larger variation of the acceptable noise in 4 cpd than in 10 cpd.

A similar behavior would be obtained for the chromatic channels.

5.2.2 Experiment 3: checking the inversion of the non-linear transform

Invertibility condition. These non-linear transforms must be invertible in order to reconstruct the image at the decoder from the coded stream. In [14] a closed-form inversion procedure was proposed and the invertibility was studied for the achromatic case. This procedure involves the inversion of the matrix $(I - D_{r_i} \cdot h)$, where I is the identity matrix, D_{r_i} is a diagonal matrix with the absolute value of the elements of $r_i(f)$ and h is the matrix that models the relations between coefficients. The inversion condition states that all the eigenvalues of $D_{r_i} \cdot h$ must be less than one. There is no theoretical guarantee that this condition holds for all images (neither in the already proposed achromatic case, nor in the chromatic cases proposed here). Therefore, the invertibility of the chromatic non-linear transforms (blocks 7.B and 7.C) has to be explicitly checked empirically.

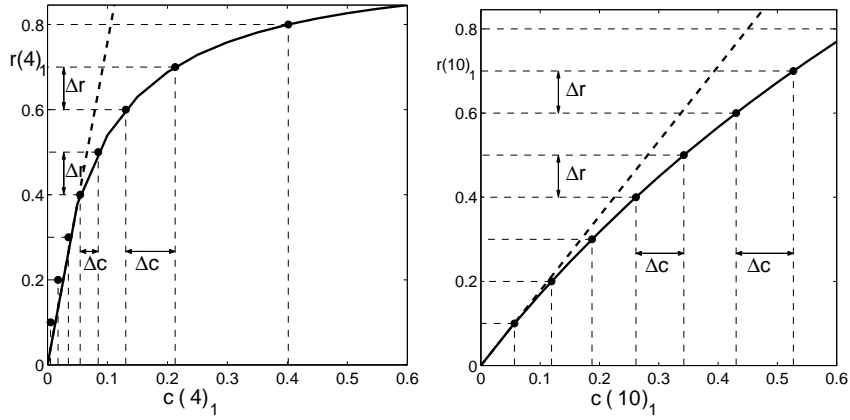


Figure 8: Responses and associated visibility thresholds of the two sensors tuned to frequencies 4 and 10 cpd in auto-masking (zero background) conditions. The required amount of distortion Δc to obtain some specific distortion in the response domain τ is shown for different contrasts of the input pattern.

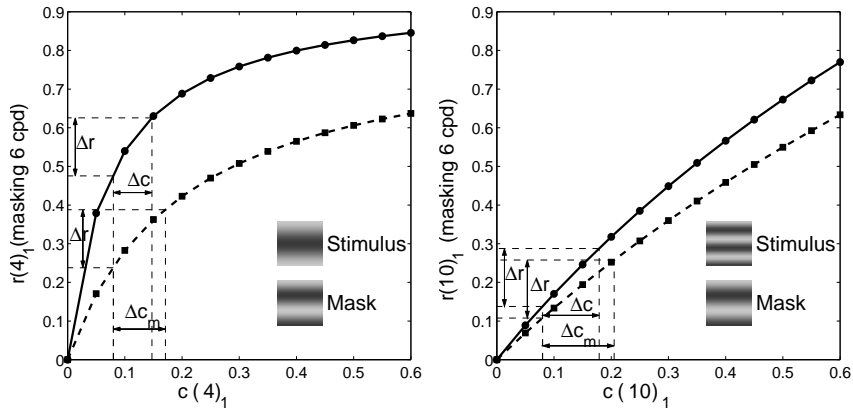


Figure 9: Responses and associated visibility thresholds of the two sensors tuned to frequencies 4 and 10 cpd when masked by a pattern of different frequency (6 cpd) at different contrast: 0 (auto-masking, solid line) and 0.5 (dashed line). In this case, the required amount of distortion Δc to obtain some specific distortion in the response domain τ at a given contrast of the stimulus increases when the contrast of the mask increases.

Testing the invertibility condition in the chromatic cases. The invertibility condition was empirically tested by computing the largest eigenvalue of the matrices $D_{r_2} \cdot h$ and $D_{r_3} \cdot h$ for 25600 16×16 color image blocks from the color image dataset [30]. The invertibility condition was also checked for reconstructed \hat{r}_i from the quantized SVM weights for different compression rates in the range $[0.2, 2.2]$ bits/pix.

In this exploration, the largest eigenvalue for every image block and compression rate was always lower than 1.

5.3 Support Vector Regression (SVR) with Adaptive Profile

An important third basic step in the coding scheme is the introduction of a machine learning method for the selection of the most relevant coefficients in the specific image representation

domain. The method must consider the particular characteristics of the domain, and thus the different perceptual relevance of each coefficient in the working domain. We illustrate this machine learning block with the use of the support vector regression (SVR) method which has demonstrated good capabilities for signal approximation with a small number of parameters (signal samples, or support vectors) [35]. The standard SVR is revised in the following subsection. Noting that the standard approach defines a constant insensitivity independently of the perceptual relevance of each sample, we subsequently propose a novel reformulation of the SVR to work in non-Euclidean domains of data representation. Finally, some remarks are provided to this important fact.

5.3.1 Standard SVR

Three SVR models are applied independently to the coefficients in each non-linear representation domain $\mathbf{r}(f) = [r_1(f), r_2(f), r_3(f)]^\top$. Hereafter, and for the sake of clarity, we will present the notation of the SVR for only one generic channel representation, r . Consequently, we are given a paired set of N coefficients (f, r) , where $f = 1, \dots, N$ represents the spatio-frequency index descriptor of coefficients $r \in \mathbb{R}$ in the particular non-linear perceptually-transformed channel (either $r_1(f)$, $r_2(f)$ or $r_3(f)$).

The standard formulation of the SVR maps the input data f to a higher dimensional space where a linear regression is solved, which is non-linearly related to the input representation space. The regression model is thus defined as

$$\hat{r}_f = \mathcal{F}(f) = \langle \mathbf{w}, \phi(f) \rangle + b \quad (37)$$

where \hat{r} are the r coefficient estimations; $\langle \cdot, \cdot \rangle$ represents the dot product operation; ϕ is a non-linear mapping to a higher dimensional Hilbert space $\phi : \mathbb{R} \rightarrow \mathcal{H}$; \mathbf{w} is a weight vector in the transformed high-dimensional feature space; and b is the bias term in the regression model.

The SVR consists of solving the following regularized problem with linear constraints:

$$\min_{\mathbf{w}, \xi_f, \xi_f^*, b} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \lambda \sum_{f=1}^N (\xi_f + \xi_f^*) \right\} \quad (38)$$

subject to:

$$r_f - (\langle \mathbf{w}, \phi(f) \rangle + b) \leq \varepsilon + \xi_f \quad \forall f = 1, \dots, N \quad (39)$$

$$(\langle \mathbf{w}, \phi(f) \rangle + b) - r_f \leq \varepsilon + \xi_f^* \quad \forall f = 1, \dots, N \quad (40)$$

$$\xi_f, \xi_f^* \geq 0 \quad \forall f = 1, \dots, N \quad (41)$$

Free parameter λ tunes the trade-off between fitting the model to the data (minimizing the errors ξ_f and ξ_f^*) and keeping model weights $\|\mathbf{w}\|$ small (enforcing flatness in the feature space). The method uses the so-called ε -insensitive loss function [35], which penalizes errors larger than ε linearly. It is worth stressing that free parameter ε accounts for the allowed error or distortion.

The usual procedure for solving this optimization problem introduces the linear restrictions (39)-(41) into Eq. (38) by means of Lagrange multipliers $\alpha_f^{(*)}$, computes the Karush-Kuhn-Tucker conditions, and solves the Wolfe's dual problem [36, 37]. This leads to solve a quadratic programming (QP) problem in a space of dual parameters $\alpha_f^{(*)}$ rather than in the space of model parameters \mathbf{w} . A very important result is that, by making zero the derivative of the obtained dual

constrained functional, the weight vector in the feature space is expressed as a linear combination of the mapped samples through the dual variables, that is:

$$\mathbf{w} = \sum_{f=1}^N (\alpha_f - \alpha_f^*) \boldsymbol{\phi}(f). \quad (42)$$

Now, by plugging (42) into (37), one obtains the solution for a particular input f' :

$$\hat{r}_{f'} = \mathcal{F}(f') = \sum_{f=1}^N (\alpha_f - \alpha_f^*) \langle \boldsymbol{\phi}(f), \boldsymbol{\phi}(f') \rangle + b, \quad (43)$$

which explicitly depends of the dot product of mapped samples and the obtained dual weights α_f and α_f^* . This result allows us to work implicitly in a higher dimensional space without even knowing the coordinates of mapped samples explicitly, only the dot products among them. These dot products are called kernel functions, $K(f, f')$, and lead to the final regression function:

$$\hat{r}_{f'} = \mathcal{F}(f') = \sum_{f=1}^N (\alpha_f - \alpha_f^*) K(f, f') + b, \quad (44)$$

where the inner product $\langle \boldsymbol{\phi}(f), \boldsymbol{\phi}(f') \rangle$ is represented with a kernel matrix $K(f, f')$. Note that only samples with non-zero Lagrange multipliers $\alpha_f^{(*)}$ count in the solution and are called *support vectors*. The immediate advantage of the method is that good approximating functions can be obtained with a (relatively) small set of support vectors, leading to the concept of *sparsity* and, in turn, to the idea of inherent compression.

The Gram (or kernel) matrix $K(f, f') = \langle \boldsymbol{\phi}(f), \boldsymbol{\phi}(f') \rangle$ can be seen as a similarity matrix among samples and its proper definition is the key in the learning ability of the SVR method. In all our experiments, the Radial Basis Function (RBF) kernel is used:

$$K(f, f') = \exp \left(- \frac{|f - f'|^2}{2\sigma^2} \right) \quad (45)$$

This kernel mapping introduces a third free parameter to be optimized, the kernel width or lengthscale σ .

5.3.2 SVR with Adaptive Insensitivity Profile

The main problem when considering the previous solution is that we assume that each sample contains *a priori* the same relevance to the function approximation, which in general is not true. This can be easily alleviated by using a different penalization factor for each sample f according to a certain *confidence function* k_f on the samples. This idea can be also extended by using a different insensitivity zone ε for each sample. The proposed SVR with adaptive profile [38] relaxes or tightens the ε -insensitive region depending on each training sample. Now, the objective function becomes [36]:

$$\min_{\mathbf{w}, \xi_f, \xi_f^*, b} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \lambda \sum_{f=1}^N k_f (\xi_f + \xi_f^*) \right\} \quad (46)$$

and restrictions over slack variables become sample-dependent:

$$r_f - (\langle \mathbf{w}, \phi(f) \rangle + b) \leq \frac{\varepsilon}{k_f} + \xi_f \quad \forall f = 1, \dots, N \quad (47)$$

$$(\langle \mathbf{w}, \phi(f) \rangle + b) - r_f \leq \frac{\varepsilon}{k_f} + \xi_f^* \quad \forall f = 1, \dots, N \quad (48)$$

$$\xi_f, \xi_f^* \geq 0 \quad \forall f = 1, \dots, N \quad (49)$$

Therefore, now each sample has its own insensitivity error $\varepsilon_f = \varepsilon/k_f$, which intuitively means that different samples hold different confidence intervals or allowed distortions. By including linear restrictions (47)-(49) in the corresponding functional (46), we can follow as in the standard case, which once again constitutes a QP problem.

5.3.3 Remarks: adaptive insensitivity and the working domain

In SVM-based image coding schemes, such as those in [15, 5, 17] the signal is described by the Lagrange multipliers of the support vectors needed to keep the regression error below the thresholds ε_f . Increasing the thresholds, ε_f , reduces the number of required support vectors, thus reducing the entropy of the encoded image and increasing the distortion. The key point here is choosing ε_f according to a meaningful criterion for the application domain. For example, working in Euclidean domains justifies the use of a constant value of ε_f for all f . This condition is met in the non-linear perceptual transformation included in the invention scheme, but it is not in linear DCT or wavelet domains such as those used in [15], in which each coefficient has a different relevance in representing the signal. See [5, 16, 17] for a more detailed discussion.

6 Experimental results to assess the performance of the method

This section illustrates the performance of an implementation of the method for the compression of a set of color images compared to the JPEG algorithm.

6.1 Experimental setup

The implementation of the proposed method, referred to as C-NL-SVR (Color-Non-Linear-SVR), includes linear YUV space for color representation, 16×16 block-DCTs for image representation, contrast transforms as in section 5.1, non-linear perceptual transform of each color channel with the parameters given in section 5.2 and SVR learning processes of the non-linear representation as described in section 5.3. Finally, the description (weights) selected by the regression algorithm are uniformly quantized and encoded according to the scheme described in section 4. The implementation of JPEG uses the same color representation and 16×16 block-DCT. Quantization matrices based on the Achromatic and Chromatic Mullen CSFs [26] were used.

A set of 25 representative color images was used for comparison purposes (see Fig. 10). All the images were compressed at different bit-rates in the range [0.1, 2.2] bits/pix. It has to be stressed that the bit rate used in consumer electronics (cell phones, low-end digital cameras, semi-professional and high-end digital cameras) ranges between 1 bit/pix for low quality JPEG files to 2.2 bits/pix for high-quality JPEG files.

measures: Structural SIMilarity (SSIM) index [39] and S-CIELab metric [40].

The averaged rate-distortion curves for all images in Fig. 10 are plotted in Fig. 11. In the rate-distortion plots we indicate with bars the standard deviation obtained for each compression rate. These plots show that the proposed C-NL-SVR is better than JPEG (the solid line is under the dashed line for distortions and is over it for similarity, in the interesting bit-rate range). However, a naive interpretation of the standard deviation bars overlapping could lead to question the significance of this eventual gain.

Note however, that in image rate-distortion plots, standard deviation bars overlapping does not necessarily mean equal behavior or statistically negligible gain: the overlapping comes from the fact that different images in the database have intrinsically different complexity giving rise to quite different distortions when encoded at a given bit rate. Figure 12(a) shows an example of the above: in this case, it is shown the rate-distortion behavior of the two methods (solid and dashed curves) for two different images of the database (black and blue curves). Of course, if you average among such class of images the standard deviation is going to be big, but the important thing is that the gain is consistent in every image: note for instance that, if you take some particular entropy for the JPEG result in both images (e.g., red or purple dots in fig. 12(a)), the entropy of the C-NL-SVR result with the same distortion is consistently much smaller.

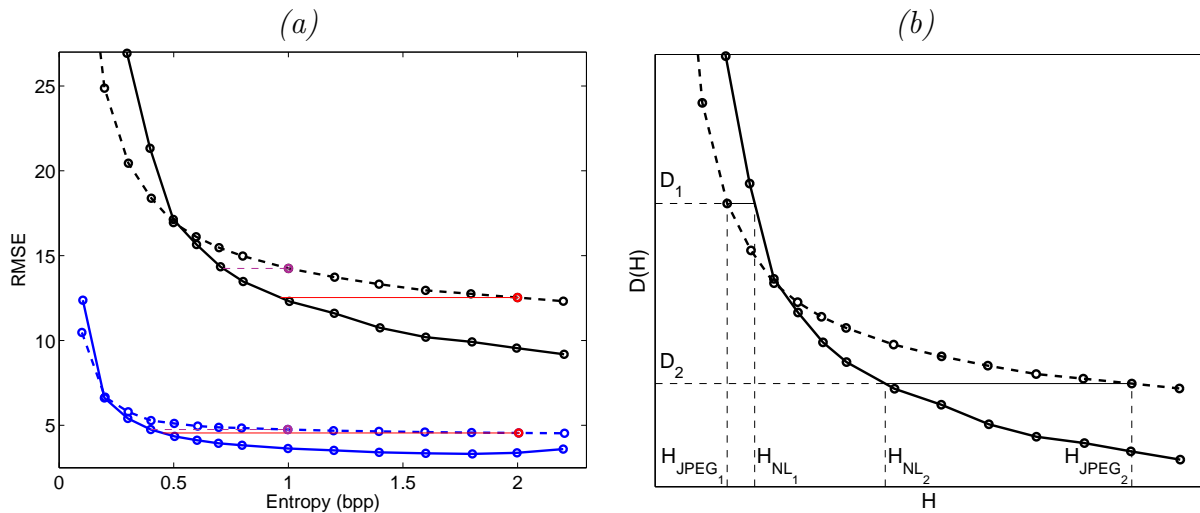


Figure 12: (a) Distortion curves for two example images (in black and blue) using RMSE for the considered 16×16 JPEG (dashed) and the C-NL-SVR approach (solid line). (b) Illustrative scheme for measuring the gain of the C-NL-SVR approach (solid line) versus JPEG (dashed). For a given distortion level, the difference in rate between methods is measured and compared.

Therefore, it makes more sense to define some *compression gain* measure for each image and bit-rate, and average these compression gains over the images in the dataset.

We define the *compression gain* of one method versus the other for a given distortion (or entropy), $D(H)$, in terms of the rate difference for the same distortion level:

$$\text{Compression Gain : } G(H) = \frac{H_{JPEG}(D(H))}{H_{C-NL-SVR}(D(H))} \quad (50)$$

Fig. 12(b) shows an illustrative scheme for measuring the compression gain for a distortion measure on a particular image (zoom of the RMSE rate-distortion for the black curves image of

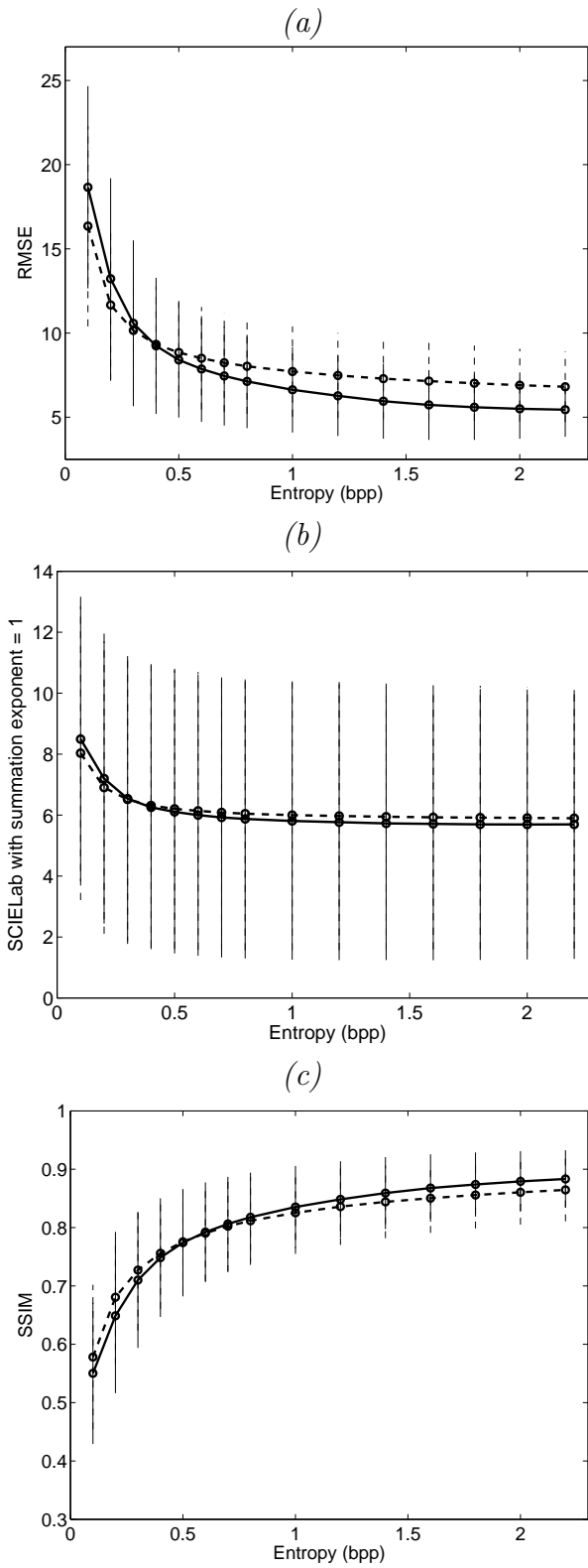


Figure 11: Average rate-distortion curves over the 25 images using (a) RMSE distortion, (b) S-CIELab distortion, and (c) Structural SIMilarity index, for the considered 16×16 JPEG (dashed) and the C-NL-SVR approach (solid line).

the left plot). In the example illustrated in fig. 12(b) we consider the two possible cases (gain smaller than 1 at low entropy and gain bigger than 1 at high entropy). Note that gain values bigger than 1 mean that the proposed method use less bits to represent the image for the same distortion. This compression gain can also be expressed in percentage terms by using:

$$\text{Compression Gain (in \%)} : PG(H) = 100 \cdot (G(H) - 1) \quad (51)$$

In practical terms, for a given entropy, $G = 2$ or $PG = 100\%$ mean that 20 C-NL-SVR images have the same volume in bits than 10 JPEG images (with the same measured quality).

Of course, the particular compression gain values depend on the selected quality measure and the dataset.

Figure 13 show the compression gains obtained by the C-NL-SVR approach versus JPEG for the three considered distortion measurements (RMSE, SSIM and S-CIELab) at every considered entropy in the analyzed bit-rate range. In these plots, we have included the standard deviation of the gains in order to assess the reliability of the approach.

Note that in the commercially meaningful range $[1, 2.2]$ bits/pix, the proposed method largely outperforms JPEG and permits compression improvements in the range $[60, 180]\%$ for RMSE, $[25, 60]\%$ for SSIM, and $[80, 275]\%$ for the S-CIELAB. As the standard deviation bars do not cross below the $G = 1$ line, this means that the gain is consistent over a wide range of images, thus the proposed method clearly outperforms JPEG.

6.3 Visual comparison

Figures 14-18 show representative results of the considered methods on 5 images ('Parrot', 'Lena', 'Roof', 'Flower3', 'Face1') at different bit rates in the range $[1.0, 2.2]$ bits/pix. The visual results confirm that the *numerical* gains shown in Fig. 13 are also *perceptually significant*. In general, JPEG leads to poorer (blocky) results. Also, it is worth noting that high frequency details are smoothed. These effects are highly alleviated by introducing a machine learning procedure like the SVR in the non-linear perceptual domain. See, for instance, Lena's eyes and cheek (Fig.15), her hat's feathers (Fig.15), the better reproduction of the high frequency pattern in the parrot's eye (Fig.14), the vertical roof's stick (Fig.16), and the annoying blocky effect in the flower (Fig.17) and face images (Fig.18).

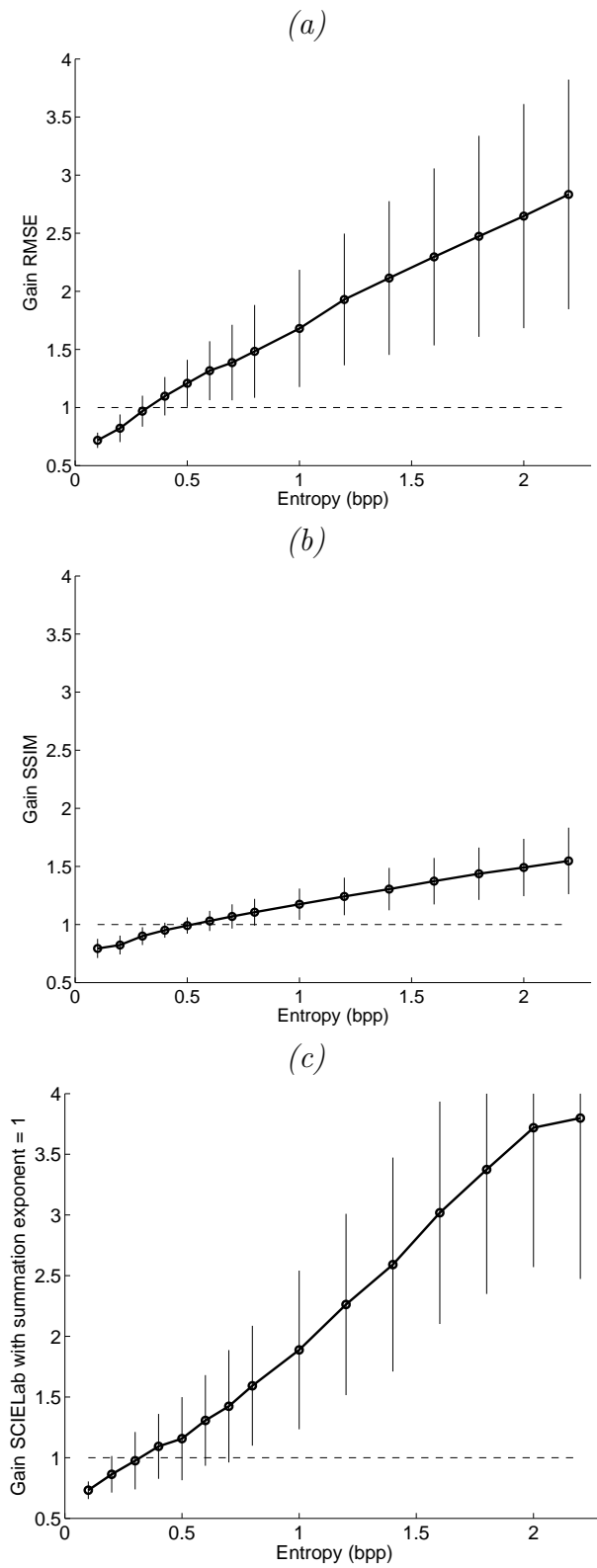


Figure 13: Average gains (with standard deviation bars) of the C-NL-SVR approach with respect to JPEG in terms of (a) RMSE distortion, (b) the Structural SIMilarity index, SSIM [39], and (c) the S-CIElab error [40].

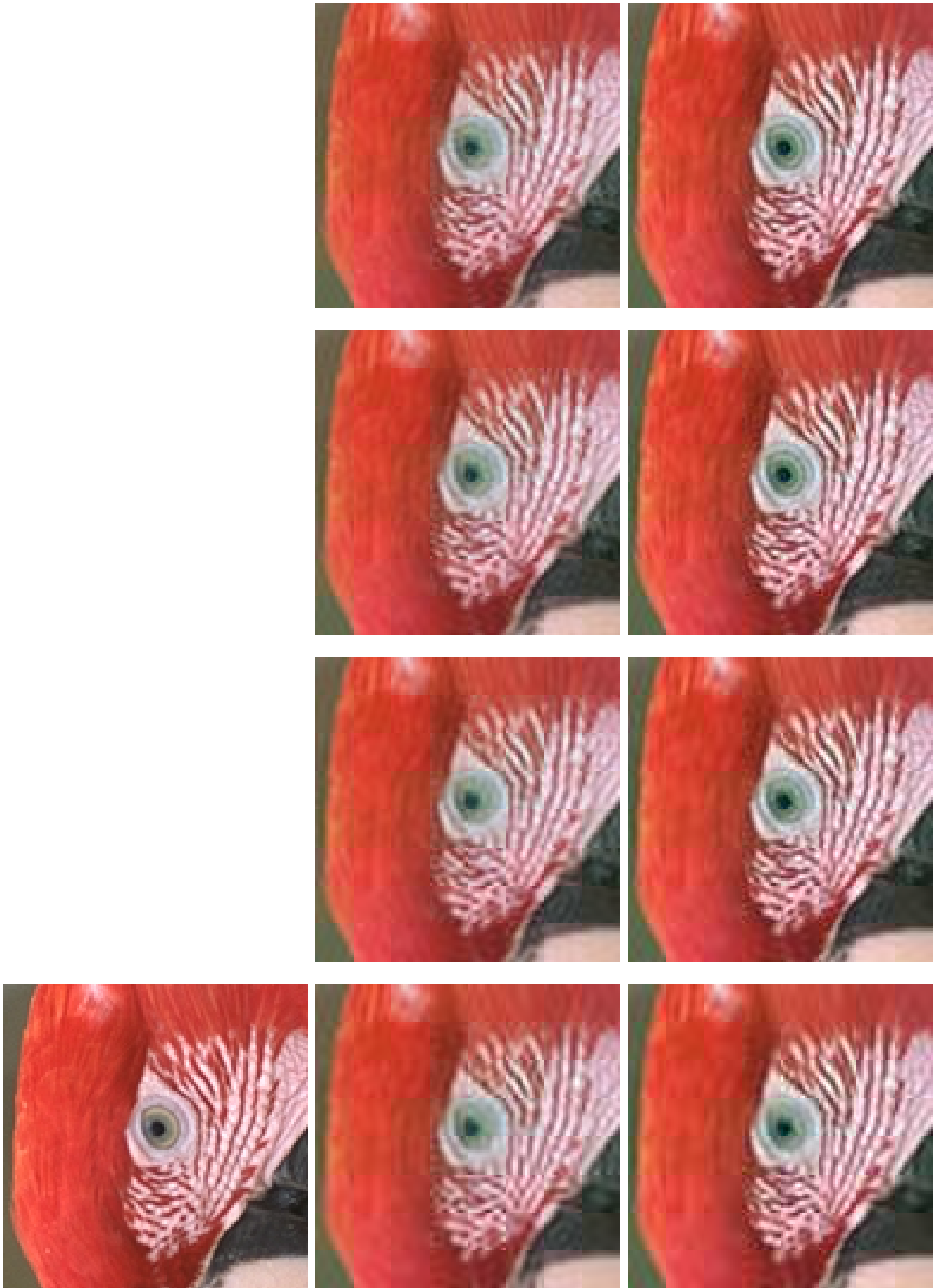


Figure 14: Examples of the decoded 'Parrot (eye)' image with JPEG (top row) and C-NL-SVR (bottom row) at different compression ratios (from left to right: $\{1.0, 1.4, 1.8, 2.2\}$ bpp).

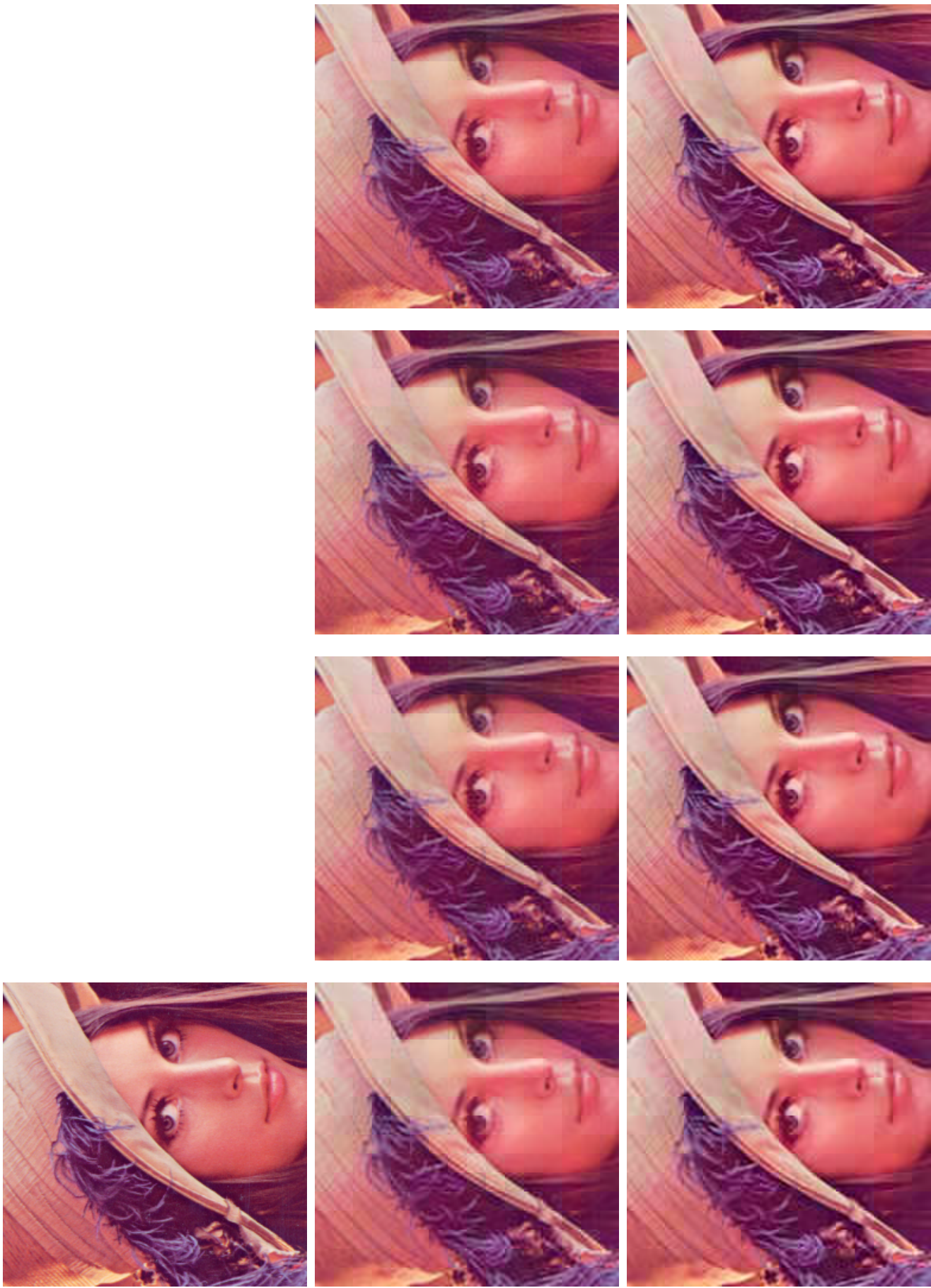


Figure 15: Examples of the decoded 'Lena' image with JPEG (top row) and C-NL-SVR (bottom row) at different compression ratios (from left to right: $\{1.0, 1.4, 1.8, 2.2\}$ bpp).

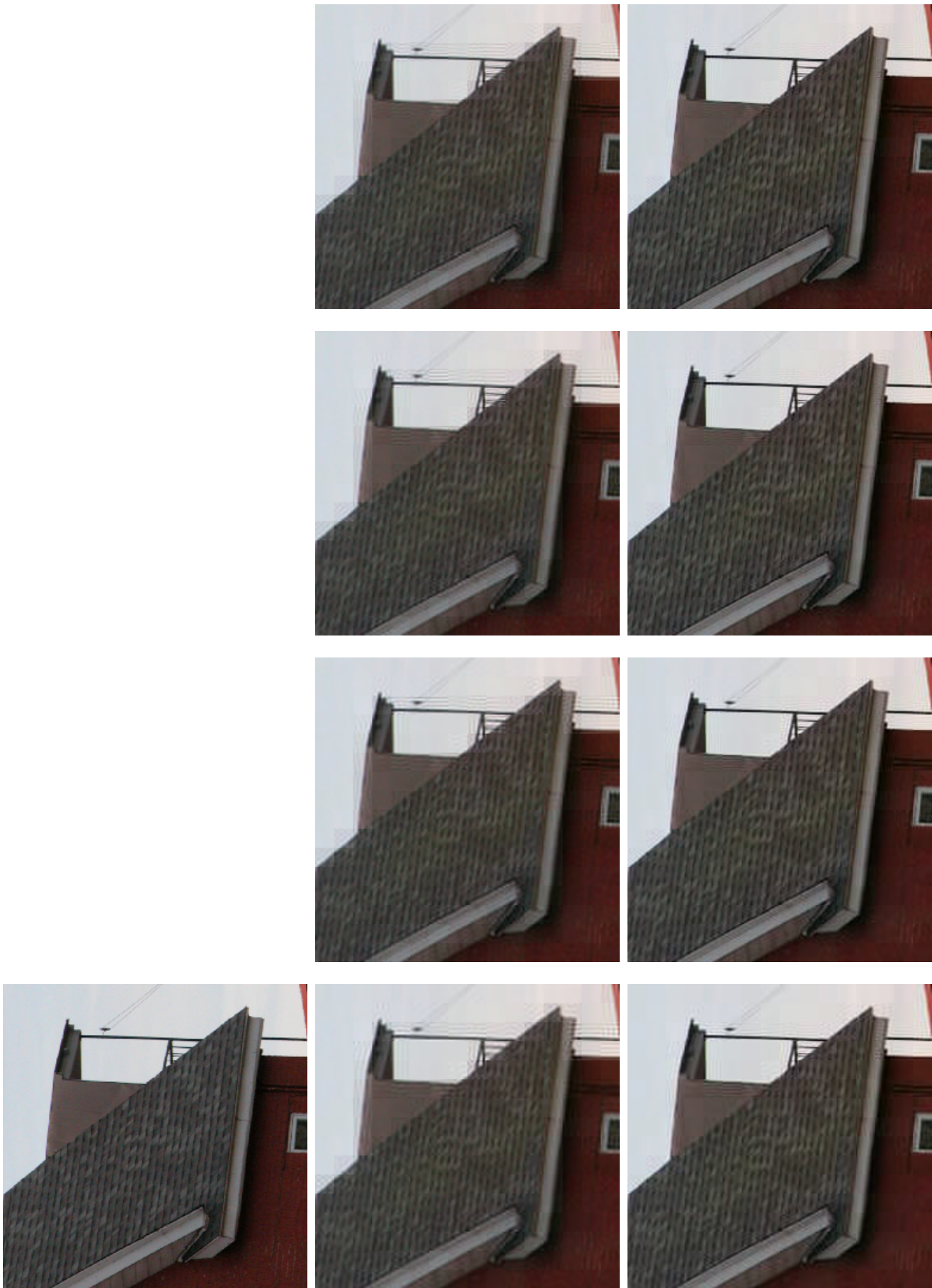


Figure 16: Examples of the decoded 'Roof' image with JPEG (top row) and C-NL-SVR (bottom row) at different compression ratios (from left to right: $\{1.0, 1.4, 1.8, 2.2\}$ bpp).

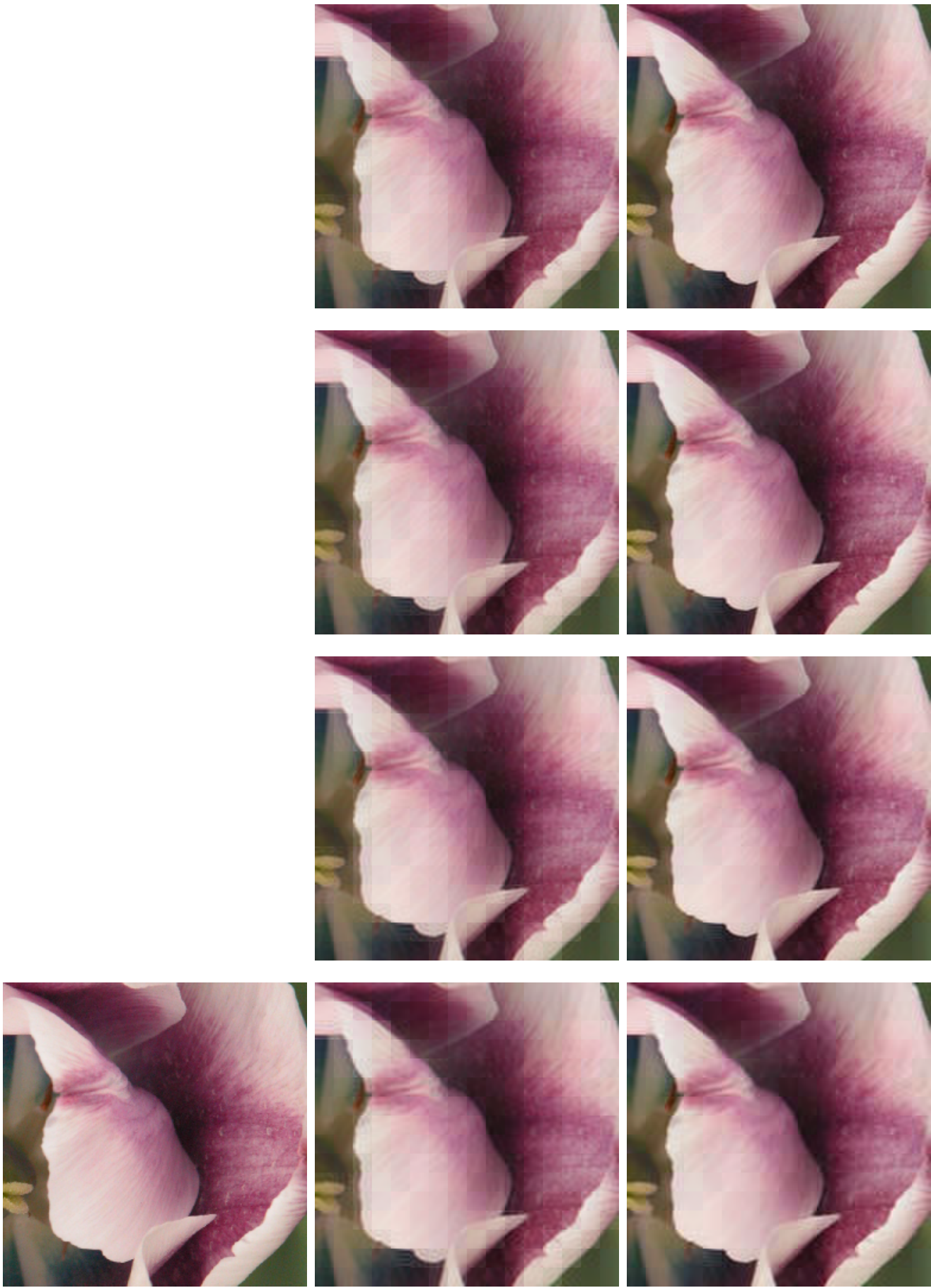


Figure 17: Examples of the decoded 'Flower 3' image with JPEG (top row) and C-NL-SVR (bottom row) at different compression ratios (from left to right: {1.0,1.4,1.8,2.2} bpp).

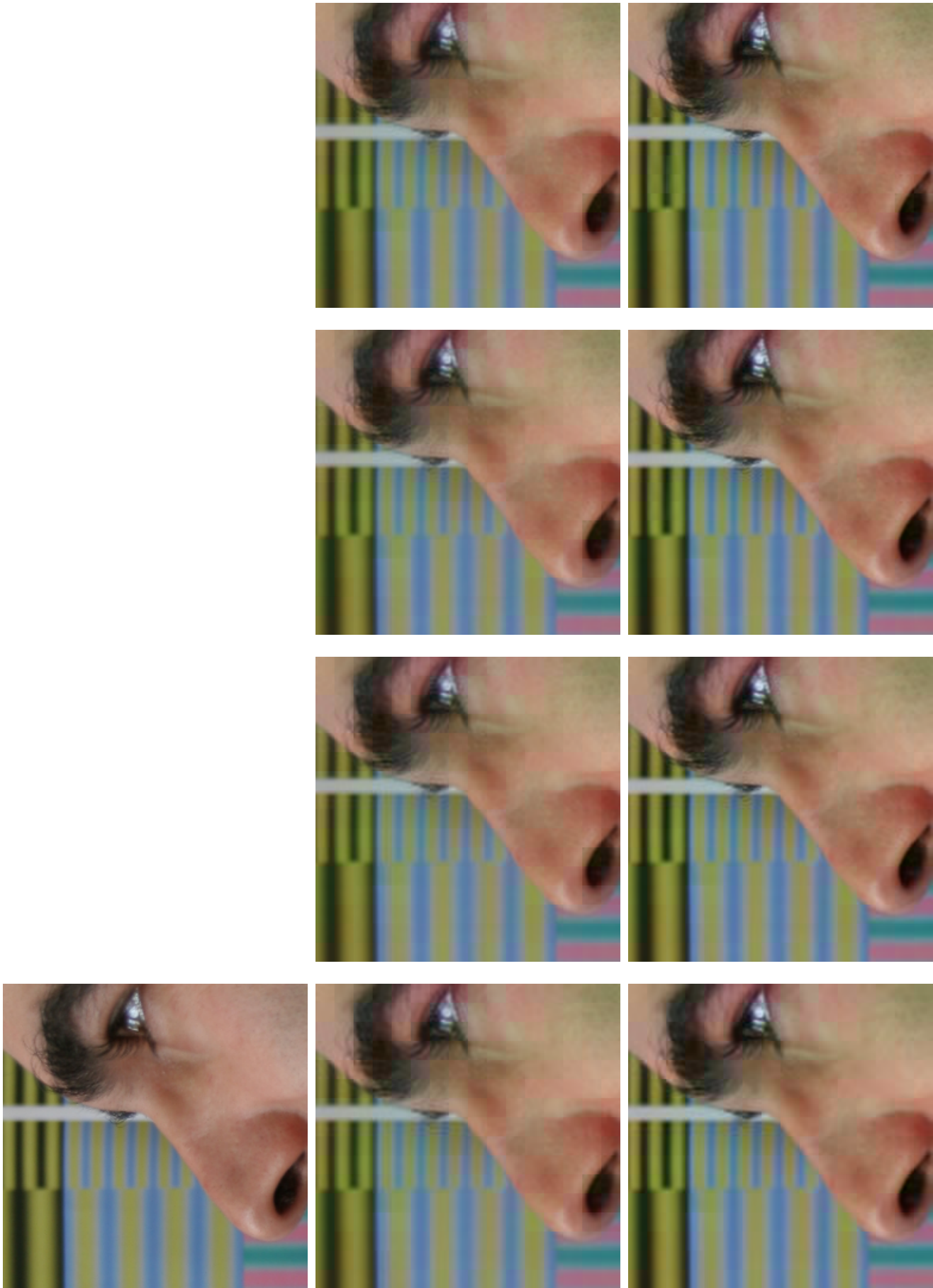


Figure 18: Examples of the decoded 'Face 1' image with JPEG (top row) and C-NL-SVR (bottom row) at different compression ratios (from left to right: $\{1.0, 1.4, 1.8, 2.2\}$ bpp).

References

- [1] G.K. Wallace. The JPEG still picture compression standard. *Communications of the ACM*, 34(4):31–43, 1991.
- [2] M. Naillon and J-B. Theeten. Method of and arrangement for image data compression by means of a neural network. Patent No. 5005206, 1991.
- [3] R. A. Nanni and G. Abraham. Neural-network-based method of image compression. Patent No. 6798914, 2004.
- [4] V. Kecman and J. Robinson. Method, apparatus and software for lossy data compression and function approximation. Patent No. WO/2003/050959, 2003.
- [5] G. Gómez, G. Camps-Valls, J. Gutiérrez, and J. Malo. Perceptual adaptive insensitivity for support vector machine image coding. *IEEE Transactions on Neural Networks*, 16(6):1574–1581, Jun 2005.
- [6] A. B. Watson. Image data compression having minimum perceptual error. Patent No. 5426512, 1995.
- [7] A.B. Watson. DCT quantization matrices visually optimized for individual images. In J. P. Allebach and B. E. Rogowitz, editors, *Human Vision, Visual Processing and Digital Display IV*, volume 1913 of *Proceedings of the SPIE*, pages 202–216. SPIE–The International Society for Optical Engineering, 1993.
- [8] J. Malo, A.M. Pons, and J.M. Artigas. Bit allocation algorithm for codebook design in vector quantization fully based on HVS non-linearities for suprathreshold contrasts. *Electronics Letters*, 31(15):1229–1231, 1995.
- [9] J. Malo, J. Gutiérrez, I. Epifanio, F. Ferri, and J. M. Artigas. Perceptual feed-back in multigrid motion estimation using an improved DCT quantization. *IEEE Transactions on Image Processing*, 10(10):1411–1427, October 2001.
- [10] J. Malo. *Tecnología del color*, chapter Almacenamiento y transmisión de imágenes en color, pages 117–164. Col.lecció Materials. Servei de Publicacions de la Universitat de Valencia, Valencia, 2002.
- [11] D.S. Taubman and M.W. Marcellin. *JPEG2000: Image Compression Fundamentals, Standards and Practice*. Kluwer Academic Publishers, Boston, 2001.
- [12] W. Zeng, S. Daly, and S. Lei. An overview of the visual optimization tools in JPEG2000. *Signal Processing: Image Communication*, 17(1):85–104, 2002.
- [13] Y. Navarro, J. Rovira, J. Gutiérrez, and J.Malo. Gain control for the chromatic channels in JPEG2000. *Proc. of the 10th Intl. Conf. AIC.*, 1:539–542, 2005.
- [14] J. Malo, I. Epifanio, R. Navarro, and R. Simoncelli. Non-linear image representation for efficient perceptual coding. *IEEE Transactions on Image Processing*, 15(1):68–80, 2006.
- [15] J. Robinson and V. Kecman. Combining Support Vector Machine learning with the discrete cosine transform in image compression. *IEEE Transactions on Neural Networks*, 14(4):950–958, July 2003.

- [16] J. Gutiérrez, G. Gómez-Pérez, J. Malo, and G. Camps-Valls. Perceptual image representations for support vector machine image coding. In G. Camps-Valls, J. L. Rojo-Álvarez, and M. Martínez-Ramón, editors, *Kernel Methods in Bioengineering, Signal and Image Processing*. Idea Group Publishing, Hershey, PA (USA), Jan 2007.
- [17] G. Camps-Valls, J. Gutiérrez, G. Gómez, and J. Malo. On the suitable domain for SVM training in image coding. *Journal of Machine Learning Research*, 9(1):49–66, 2008.
- [18] J. Malo and M. J. Luque. *COLORLAB: A Color Processing Toolbox for Matlab*. Dept. d'Òptica, Universitat de València, 2002. Available at: <http://www.uv.es/vista/vistavalencia>.
- [19] M. D. Fairchild. *Color Appearance Models*. Addison-Wesley, New York, 1997.
- [20] E. Peli. Contrast in complex images. *Journal of the Optical Society of America A*, 7:2032–2040, 1990.
- [21] D. J. Heeger. Normalization of cell responses in cat striate cortex. *Visual Neuroscience*, 9:181–198, 1992.
- [22] A. Gersho and R. M. Gray. *Vector Quantization and Signal Compression*. Kluwer Academic Press, Boston, 1992.
- [23] J. M. Foley. Human luminance pattern mechanisms: Masking experiments require a new model. *Journal of the Optical Society of America A*, 11(6):1710–1719, 1994.
- [24] A. B. Watson and J. A. Solomon. A model of visual contrast gain control and pattern masking. *Journal of the Optical Society of America A*, 14(9):2379–2391, September 1997.
- [25] E. Martinez-Uriegas. Color detection and color contrast discrimination thresholds. In *Proceedings of the OSA Annual Meeting ILS–XIII*, page 81, Los Angeles, 1997.
- [26] K. T. Mullen. The contrast sensitivity of human colour vision to red-green and yellow-blue chromatic gratings. *Journal of Physiology*, 359:381–400, 1985.
- [27] D. H. Kelly. Spatiotemporal variation of chromatic and achromatic contrast thresholds. *Journal of the Optical Society of America A*, 73(6):742–749, 1983.
- [28] W. K. Pratt. *Digital image processing (2nd ed.)*. John Wiley & Sons, Inc., New York, NY, USA, 1991.
- [29] J. S. Lim. *Two-dimensional signal and image processing*. Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 1990.
- [30] A. Parraga. *McGill Calibrated Colour Image Database*. Dept. of Vision Science, McGill University, 2003. Available at <http://tabby.vision.mcgill.ca>.
- [31] T. T. Norton, D. A. Corliss, and J. E. Bailey. *The Psychophysical Measurement of Visual Function*. Butterworth Heinemann, 2002.
- [32] E. M. Uriegas. Personal communication on contrast incremental thresholds for chromatic gratings measured at stanford research international. Unpublished results, May 1998.

- [33] A. B. Watson. Personal communication on plausible parameters for divisive normalization in dct basis. Unpublished results, March 2001.
- [34] D. Heeger. Personal communication on plausible parameters for divisive normalization in dct basis. Unpublished results, March 2001.
- [35] A. J. Smola and B. Schölkopf. A tutorial on support vector regression. *Statistics and Computing*, 14:199–222, 2004.
- [36] V. Vapnik. *Statistical Learning Theory*. John Wiley & Sons, New York, 1998.
- [37] B. Schölkopf and A. Smola. *Learning with Kernels – Support Vector Machines, Regularization, Optimization and Beyond*. MIT Press Series, 2002.
- [38] G. Camps-Valls, E. Soria-Olivas, J. Pérez-Ruixo, A. Artés-Rodríguez, F. Pérez-Cruz, and A. Figueiras-Vidal. A profile-dependent kernel-based regression for cyclosporine concentration prediction. In *Neural Information Processing Systems (NIPS) – Workshop on New Directions in Kernel-Based Learning Methods*, Vancouver, Canada, December 2001. Available at <http://www.uv.es/~gcamps>.
- [39] 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*, 13(4):600–612, 2004.
- [40] X. Zhang and B. Wandell. Color image fidelity metrics evaluated using image distortion maps. *Signal Processing*, 70(3):201–214, 1998.