
Quantization-free Lossy Image Compression Using Integer Matrix Factorization

Pooya Ashtari^{1*}†

pooya.ashtari@esat.kuleuven.be

Pourya Behmandpoor^{1†}

pourya.behmandpoor@esat.kuleuven.be

Fateme Nateghi Haredasht²

fnateghi@stanford.edu

Jonathan H. Chen²

jonc101@stanford.edu

Lieven De Lathauwer¹

lieven.delathauwer@kuleuven.be

Sabine Van Huffel¹

sabine.vanhuffel@esat.kuleuven.be

¹Department of Electrical Engineering (ESAT), STADIUS Center, KU Leuven, Belgium

²Department of Medicine, Stanford University, Stanford, CA, USA

Abstract

one paragraph

1 Introduction

2 Related Work

3 Method

3.1 Overall Framework

Figure 1 illustrates an overview of the encoding pipeline for our proposed image compression method using integer matrix factorization (IMF). The encoder accepts an RGB image with dimensions $H \times W$ and a color depth of 8 bits, represented by the tensor $\mathcal{X} \in \{0, \dots, 255\}^{3 \times H \times W}$. Each step of encoding is described in the following.

Color Space Transformation. Analogous to the JPEG standard, the image is initially transformed into the $Y\text{C}_B\text{C}_R$ color space. Let $\mathbf{Y} \in [0, 255]^{H \times W}$ represent the *luma* component, and $\mathbf{C}_B, \mathbf{C}_R \in [0, 255]^{\frac{H}{2} \times \frac{W}{2}}$ represent the blue-difference and red-difference *chroma* components, respectively. Note that as a result of this transformation, the elements of the *luma* (\mathbf{Y}) and *chroma* ($\mathbf{C}_B, \mathbf{C}_R$) matrices are not limited to integers and can take any value within the interval $[0, 255]$.

Chroma Downsampling. After conversion to the $Y\text{C}_B\text{C}_R$ color space, the *chroma* components \mathbf{C}_B and \mathbf{C}_R are downsampled using average-pooling with a kernel size of $(2, 2)$ and a stride of $(2, 2)$, similar to the process used in JPEG. This downsampling exploits the fact that the human visual system perceives far more detail in brightness information (*luma*) than in color saturation (*chroma*).

*Corresponding author. Emails:pooya.ashtari@esat.kuleuven.be, pooya.ash@gmail.com

†Equal contribution

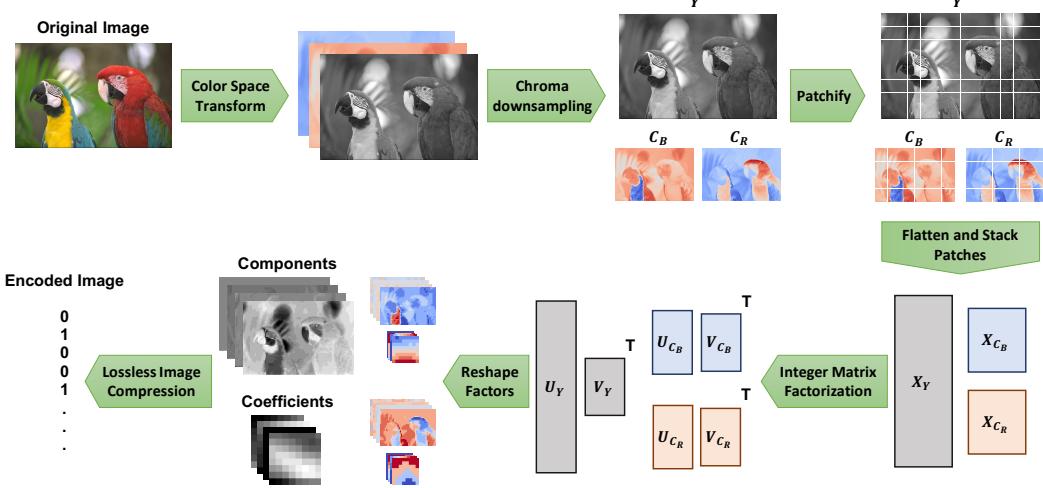


Figure 1 An illustration of the encoder for our image compression method, based on integer matrix factorization.

Patchification. Following *chroma* downsampling, we have three components: the *luma* component $\mathbf{Y} \in [0, 255]^{H \times W}$ and the *chroma* components $\mathbf{C}_B, \mathbf{C}_R \in [0, 255]^{\frac{H}{2} \times \frac{W}{2}}$. Each of the matrices is split into non-overlapping 8×8 patches. If a dimension of a matrix is not divisible by 8, the matrix is first padded to the nearest size divisible by 8 using reflection of the boundary values. These patches are then flattened into row vectors and stacked vertically to form matrices $\mathbf{X}_Y \in [0, 255]^{\frac{HW}{64} \times 64}$, $\mathbf{X}_{C_B} \in [0, 255]^{\frac{HW}{256} \times 64}$, and $\mathbf{X}_{C_R} \in [0, 255]^{\frac{HW}{256} \times 64}$. Later, these matrices will be low-rank approximated using integer matrix factorization (IMF). Note that this patchification technique differs from the block splitting in JPEG, where each block is subject to discrete cosine transform (DCT) and processed independently. This patchification technique not only captures the locality and spatial dependencies of neighboring pixels but also performs better with the matrix decomposition approach to image compression.

Low-rank approximation. We now apply a low-rank approximation to the matrices \mathbf{X}_Y , \mathbf{X}_{C_B} , and \mathbf{X}_{C_R} , which is the core of our compression method that provides a lossy compressed representation of these matrices. The low-rank approximation [2] aims to approximate a given matrix $\mathbf{X} \in \mathbb{R}^{M \times N}$ by

$$\mathbf{X} \approx \mathbf{U}\mathbf{V}^T = \sum_{r=1}^R \mathbf{U}_{:r} \mathbf{V}_{:r}^T, \quad (1)$$

where $\mathbf{U} \in \mathbb{R}^{M \times R}$ and $\mathbf{V} \in \mathbb{R}^{N \times R}$ are *factor matrices* (or simply *factors*), and $R \leq \min(M, N)$ represents the *rank*. We refer to \mathbf{U} as the *basis matrix* and \mathbf{V} as the *coefficient matrix*. By selecting a sufficiently small value for R , the factor matrices \mathbf{U} and \mathbf{V} , with a combined total of $(M + N)R$ elements, offer a compact representation of the original matrix \mathbf{X} , which has MN elements, capturing the most significant patterns in the image. Depending on the loss function used to measure the reconstruction error between \mathbf{X} and the product $\mathbf{U}\mathbf{V}^T$, as well as the constraints on the factor matrices \mathbf{U} and \mathbf{V} , various formulations and variants have been proposed for different purposes. In Section 3.2, we introduce and elaborate on our variant, termed integer matrix factorization (IMF), and argue why it is well-suited and effective for image compression.

Reshape factors. IMF yields integers factor matrices $\mathbf{U}_Y \in \{0, \dots, 255\}^{\frac{HW}{64} \times R}$ and $\mathbf{V}_Y \in \{0, \dots, 255\}^{64 \times R}$; $\mathbf{U}_{C_B} \in \{0, \dots, 255\}^{\frac{HW}{256} \times R}$ and $\mathbf{V}_{C_B} \in \{0, \dots, 255\}^{64 \times R}$; and $\mathbf{U}_{C_R} \in \{0, \dots, 255\}^{\frac{HW}{256} \times R}$ and $\mathbf{V}_{C_R} \in \{0, \dots, 255\}^{64 \times R}$ that correspond to \mathbf{X}_Y , \mathbf{X}_{C_B} , and \mathbf{X}_{C_R} . We reshape these matrices by unfolding their first dimension to obtain R -channel 2D spatial maps,

referred to as *factor maps* and represented by the following tensors:

$$\begin{aligned}\mathcal{U}_Y &\in \{0, \dots, 255\}^{R \times \frac{H}{8} \times \frac{W}{8}}, \\ \mathcal{U}_{C_B}, \mathcal{U}_{C_R} &\in \{0, \dots, 255\}^{R \times \frac{H}{16} \times \frac{W}{16}}, \\ \mathcal{V}_Y, \mathcal{V}_{C_B}, \mathcal{V}_{C_R} &\in \{0, \dots, 255\}^{R \times 8 \times 8}.\end{aligned}\tag{2}$$

Lossless image compression. Since *factor maps* are already integer tensors, we do not need a quantization step, which is commonly present in other lossy image compression methods and adds extra complications. This is the main advantage of our approach. Moreover, each channel of a *factor map* can be treated as a 8-bit grayscale image and therefore can be directly encoded by any standard lossless image compression method, such as PNG or WebP. For images with a resolution of $H, W \gg 64$, which are most common nowadays, the basis maps (\mathcal{U}) are significantly larger than the coefficient maps (\mathcal{V}), accounting for the majority of the storage space. Interestingly, in practice, the IMF basis maps turn out to be meaningful images, each capturing some visual semantic of the image (see Figure 2 for an example). Therefore, our IMF approach can effectively leverage the power of existing lossless image compression algorithms, offering a significant advantage over current methods.

3.2 Integer Matrix Factorization (IMF)

The main building block of our method is integer matrix factorization (IMF), which is responsible for lossy compression of matrices obtained through patchification. IMF can be framed as an optimization problem, aiming to minimize the reconstruction error between the original matrix $\mathbf{X} \in \mathbb{R}^{M \times N}$ and the product $\mathbf{U}\mathbf{V}^\top$, while ensuring that the elements of the factor matrices \mathbf{U} and \mathbf{V} are integers within a specified interval $[\alpha, \beta]$. Formally, the IMF problem can be expressed as:

$$\begin{aligned}\min_{\mathbf{U}, \mathbf{V}} \quad & \|\mathbf{X} - \mathbf{U}\mathbf{V}^\top\|_F^2 \\ \text{s.t.} \quad & \mathbf{U} \in \mathbb{Z}_{[\alpha, \beta]}^{M \times R}, \mathbf{V} \in \mathbb{Z}_{[\alpha, \beta]}^{N \times R},\end{aligned}\tag{3}$$

where $\|\cdot\|_F$ denotes the Frobenius norm; $R \leq \min(M, N)$ represents the *rank*; and $\mathbb{Z}_{[\alpha, \beta]} \triangleq [\alpha, \beta] \cap \mathbb{Z}$ denotes the set of integers within $[\alpha, \beta]$. Without constraints on the factors, the problem would have an analytic solution through the singular value decomposition (SVD), as addressed by the Eckart–Young–Mirsky theorem [2]. If only a nonnegativity constraint were applied (without integrality), variations of nonnegative matrix factorization (NMF) would emerge [4, 3]. The IMF problem (3) poses a challenging integer program, with finding its global minima known to be NP-hard [1, 6]. Only a few iterative algorithms [1, 5] have been proposed to find a “good solution” for some IMF variants in contexts other than image compression. In Section 3.3, we propose an efficient iterative algorithm for the IMF problem (3).

The application of SVD and NMF in image compression is problematic mainly because the resulting factors contain continuous values that must be represented as arrays of floating-point numbers. This necessitates a quantization step that not only adds extra complications but also significantly degrades compression performance due to quantization errors (as demonstrated in Section 4). Conversely, our IMF formulation produces integer factor matrices that can be directly stored and losslessly processed without incurring roundoff errors. The reason for limiting the feasible region to $[\alpha, \beta]$ in our IMF formulation is to enable more compact storage of the factors using standard integral data types, such as `int8` and `int16`, supported by programming languages. Given that the elements of the input matrix \mathbf{X} are in $[0, 255]$, we found the signed `int8` type, which represents integers from -128 to 127, suitable for image compression applications. As a result, our IMF formulation is well-suited for image compression, effectively integrating the factorization and quantization steps into a single, efficient compression process.

3.3 Block Coordinate Descent Scheme for IMF

We propose an efficient algorithm for IMF using the block coordinate descent (BCD) scheme (aka alternating optimization). Starting with some initial parameter values, this approach involves sequentially minimizing the cost function with respect to a single column of a factor at a time, while keeping the other columns of that factor and the entire other factor fixed. This process is repeated until

a stopping criterion is met, such when the change in the cost function value falls below a predefined threshold or the maximum number of iterations is reached. Formally, this involves solving one of the following subproblems at a time:

$$\mathbf{u}_r \leftarrow \arg \min_{\mathbf{u}_r \in \mathbb{Z}_{[\alpha, \beta]}^M} \|\mathbf{E}_r - \mathbf{u}_r \mathbf{v}_r^\top\|_F^2, \quad (4)$$

$$\mathbf{v}_r \leftarrow \arg \min_{\mathbf{v}_r \in \mathbb{Z}_{[\alpha, \beta]}^N} \|\mathbf{E}_r - \mathbf{u}_r \mathbf{v}_r^\top\|_F^2, \quad (5)$$

where $\mathbf{u}_r = U_{:,r}$ and $\mathbf{v}_r = V_{:,r}$ represent the r -th columns of \mathbf{U} and \mathbf{V} , respectively. $\mathbf{E}_r = \mathbf{X} - \sum_{s \neq r}^R \mathbf{u}_s \mathbf{v}_s^\top$ is the residual matrix. We define one iteration of BCD as a complete cycle of updates across all the columns of both factors. In fact, the proposed algorithm is a $2R$ -block coordinate descent procedure, where at each iteration, first the columns of \mathbf{U} and then the columns of \mathbf{V} are updated (see Algorithm 1). Note that subproblem (5) can be transformed into the same form as (4) by simply transposing its error term inside the Frobenius norm. Therefore, we only need to find the best rank-1 approximation with integer elements constrained within a specific interval. Fortunately, this problem has a closed-form solution, as addressed by Theorem 1 below.

Theorem 1 (Integer rank-1 approximation). *Let $\mathbf{E} \in \mathbb{R}^{M \times N}$, $\mathbf{u} \in \mathbb{Z}_{[\alpha, \beta]}^M$, and $\mathbf{v} \in \mathbb{R}^{N \times 1}$. A global solution to the problem*

$$\mathbf{u}^* \in \arg \min_{\mathbf{u} \in \mathbb{Z}_{[\alpha, \beta]}^M} \|\mathbf{E} - \mathbf{u}\mathbf{v}^\top\|_F^2 \quad (6)$$

is given by

$$\mathbf{u}^* = \text{clamp}_{[\alpha, \beta]} \left(\text{round} \left(\frac{\mathbf{E}\mathbf{v}}{\|\mathbf{v}\|^2} \right) \right), \quad (7)$$

where $\text{round}(\mathbf{Z})$ denotes the operator that rounds each element of \mathbf{Z} to the nearest integer, and $\text{clamp}_{[\alpha, \beta]}(\mathbf{Z}) \triangleq \max(\alpha, \min(\mathbf{Z}, \beta))$ denotes the operator that clamps each element of \mathbf{Z} to the interval $[\alpha, \beta]$.

Proof. See Appendix A for the proof. \square

Using Theorem 1, we obtain the following update formulas for our BCD algorithm:

$$\mathbf{u}_r \leftarrow \text{clamp}_{[\alpha, \beta]} \left(\text{round} \left(\frac{\mathbf{E}_r \mathbf{v}_r}{\|\mathbf{v}_r\|^2} \right) \right), \quad (8)$$

$$\mathbf{v}_r \leftarrow \text{clamp}_{[\alpha, \beta]} \left(\text{round} \left(\frac{\mathbf{E}_r^\top \mathbf{u}_r}{\|\mathbf{u}_r\|^2} \right) \right). \quad (9)$$

It is noteworthy that the combination of $\text{round}(\cdot)$ and $\text{clamp}_{[\alpha, \beta]}(\cdot)$ can be interpreted as the element-wise projector to $\mathbb{Z}_{[\alpha, \beta]}$. Thanks to Theorem 1, the proposed algorithm offers an encouraging property that guarantees the cost function is monotonically non-increasing with each update.

Initialization. The initial values of factors can significantly impact the convergence performance of the BCD algorithm. We found that the convergence with naive random initialization can be too slow. To address this issue, we propose an initialization method using SVD. The procedure is straightforward. First, the truncated SVD of the input matrix $\mathbf{X} \in \mathbb{R}^{M \times N}$ is computed as $\tilde{\mathbf{U}}\Sigma\tilde{\mathbf{V}}^\top$, where $\Sigma \in \mathbb{R}^{R \times R}$ is a diagonal matrix corresponding to the R largest singular values. $\tilde{\mathbf{U}} \in \mathbb{R}^{M \times R}$ contains the corresponding left-singular vectors in its columns, and $\tilde{\mathbf{V}} \in \mathbb{R}^{N \times R}$ contains the corresponding right-singular vectors in its columns. The initial factors are then calculated as follows:

$$\mathbf{U}^0 = \text{clamp}_{[\alpha, \beta]}(\text{round}(\tilde{\mathbf{U}}\Sigma^{\frac{1}{2}})), \quad (10)$$

$$\mathbf{V}^0 = \text{clamp}_{[\alpha, \beta]}(\text{round}(\Sigma^{\frac{1}{2}}\tilde{\mathbf{V}})). \quad (11)$$

Essentially, this means we first low-rank approximate \mathbf{X} and then project the elements of the resulting factor matrices into $\mathbb{Z}_{[\alpha, \beta]}$. Algorithm 1 provides the pseudocode for the proposed BCD algorithm for IMF.

Algorithm 1: The proposed block coordinate descent (BCD) algorithm for IMF.

Input: $X \in \mathbb{R}^{M \times N}$, factorization rank R
Output: Factor matrices $U \in \mathbb{Z}_{[\alpha, \beta]}^{M \times R}$ and $V \in \mathbb{Z}_{[\alpha, \beta]}^{N \times R}$

- 1 Initialize U, V using the truncated SVD method, provided by (10) and (11)
- 2 **while** stopping criterion not satisfied **do**
- 3 $A \leftarrow XV, B \leftarrow V^T V$
- 4 **for** $r = 1, \dots, R$ **do**
- 5 $U_{:,r} \leftarrow \text{clamp}_{[\alpha, \beta]} \left(\text{round} \left(\frac{A_{:,r} - \sum_{s \neq r} B_{sr} U_{:,s}}{\|V_{:,r}\|^2} \right) \right)$
- 6 **end**
- 7 $A \leftarrow X^T U; B \leftarrow U^T U$
- 8 **for** $r = 1, \dots, R$ **do**
- 9 $V_{:,r} \leftarrow \text{clamp}_{[\alpha, \beta]} \left(\text{round} \left(\frac{A_{:,r} - \sum_{s \neq r} B_{sr} V_{:,s}}{\|U_{:,r}\|^2} \right) \right)$
- 10 **end**
- 11 **end**
- 12 **return** (U, V)

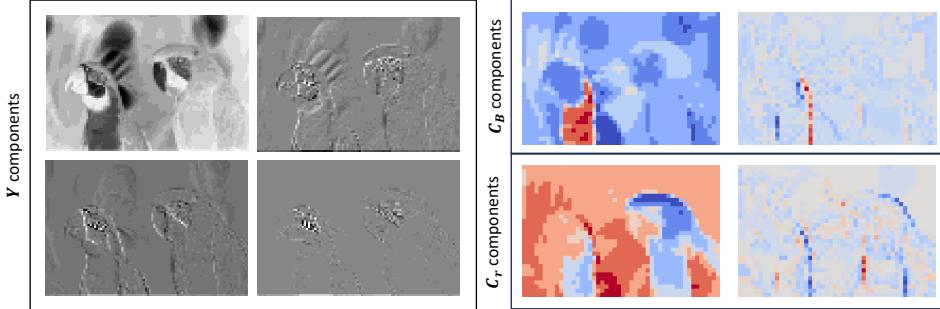


Figure 2 IMF components of the kodim23 image from the Kodak dataset, corresponding to luma (Y), blue-difference (C_b), and red-difference (C_r) chroma.

3.4 Implementation Details

4 Experiments

In this section, the proposed IMF-based compression method is assessed against SVD-based and JPEG [] methods. The performance is reported qualitatively and also based on two criteria, namely rate-distortion performance and image classification performance. Moreover, ablation studies are presented to investigate the effect of different hyperparameters in IMF.

4.1 Qualitative Performance

Qualitative performance is shown on an image selected from the Kodak [] dataset. Fig.2 depicts the first IMF components U_Y, U_{Cb} , and U_{Cr} which are extracted following the procedure elaborated in Section 3.4. It is evident in the figure that the IMF components with higher energy maintain the overall texture of the image in each channel, while components with lower energy focus more on subtle changes.

The qualitative comparison is made in Fig.?? on an image selected from the Kodak [] dataset. The images compressed by the considered compression methods are shown in different bits per pixel (bpp) values, a standard compression measure in the literature []. As can be seen, the IMF-based compression method is capable of maintaining quality compressed images in bpp values as low as ?, outperforming JPEG and SVD-based methods which suffer from maintaining balanced colors as



Figure 3 Qualitative performance comparison between various compression methods. The top row shows performance on kodim21 from the Kodak dataset in the bpp value of 0.3, and the bottom row shows one example of the Clic dataset in the bpp values of 0.14, 0.12, and 0.1 for JPEG, SVD, and IMF, respectively.

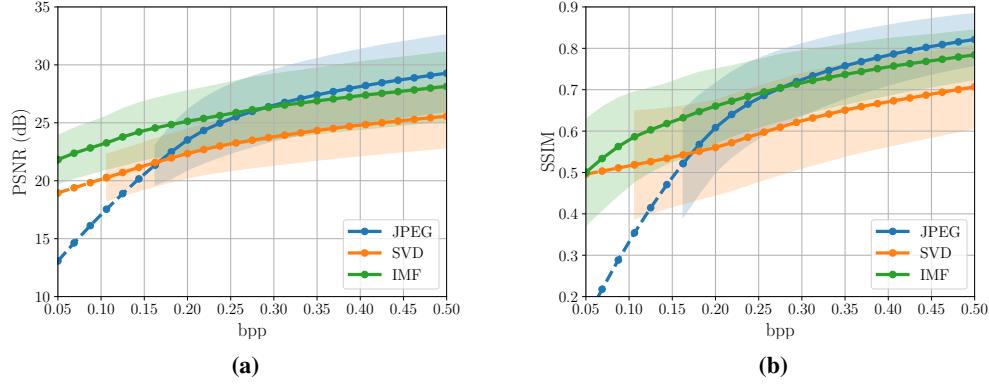


Figure 4 Rate-distortion performance on the Kodak dataset. In panels (a) and (b), the average PSNR and SSIM are plotted against bpp, respectively.

soon as bpp drops below ? and ?, respectively. The artifacts in color are visible, e.g. by JPEG in bpp values starting ?.

4.2 Rate-Distortion Performance

Peak signal-to-noise ratio (PSNR), as well as structural similarity index measure (SSIM) [], are reported versus bpp for the considered methods. The considered datasets are Clic [] and Kodak, consisting of respectively 32 and 24 colored images of various sizes. For each image compressed by any of the considered methods, bpp, PSNR, and SSIM are calculated. Then, for each compression method, PSNR and SSIM values are interpolated linearly in the fixed bpp values in the range (0.05, 0.5). In the following plots, the average performance over all images is reported, along with the standard deviation presented as shadows. For the missing bpp values, the average is extrapolated quadratically and is shown by dashed lines.

In Fig.4, the compression performance on the Kodak dataset is reported. In this figure, the proposed IMF-based method outperforms the SVD-based method, which can be attributed to the quantization errors that SVD is prone to during encoding and decoding, deteriorating its performance in both criteria. In this view, the quantization-free property of IMF effectively guarantees higher performance in different bpp values. It is also evident that the IMF-based method outperforms JPEG in low bpp values. The same performance regime can also be concluded for all the mentioned compression methods on the Clic dataset in Fig.5.

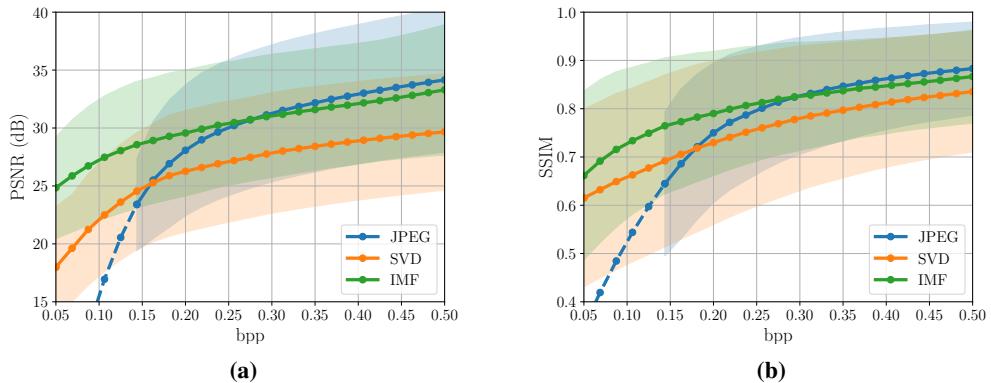


Figure 5 Rate-distortion performance on the CLIC dataset. In panels (a) and (b), the average PSNR and SSIM are plotted against bpp, respectively.

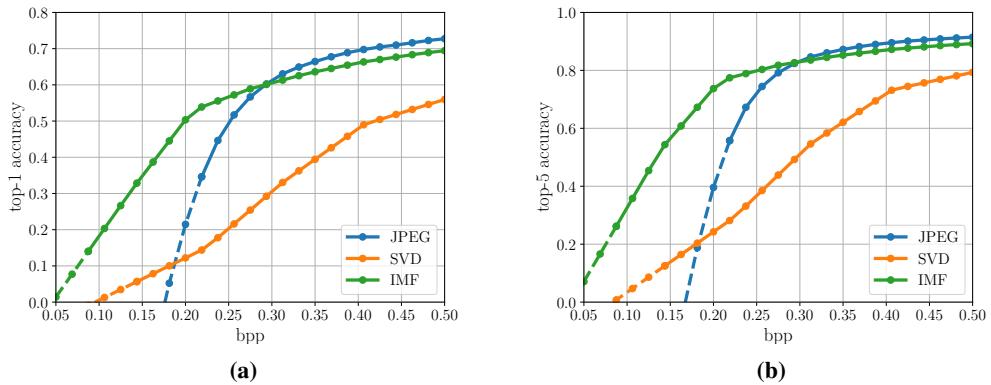


Figure 6 Impact of different compression methods on ImageNet classification accuracy. Panels (a) and (b) show the validation top-1 and top-5 accuracy plotted against bits per pixel (bpp), respectively. A ResNet-50 model pre-trained on the original ImageNet images is evaluated using validation images compressed by different methods.

4.3 ImageNet Classification Performance

As another criterion, classification performance is investigated on the images compressed by the considered compression methods. This criterion focuses on the higher-level information required for object recognition and classification embedded in each image. Furthermore, it highlights the importance of image compression where various vision tasks such as classification are the main objective—rather than maintaining the perceived image quality—while keeping the requirement of resources such as memory, communication bandwidth, computation power, latency budget, etc. as limited as possible. ImageNet [] validation set, consisting of 50000 224×224 colored images in 1000 classes, is considered for this classification task done by a ResNet-50 classifier [], pre-trained on the original ImageNet dataset. The classification performance comparison is made in Fig.6, showing the higher compression performance of IMF for classification, reaching top-1 accuracy of 70% in bpp values as low as 0.2.

4.4 Ablation Studies

In this section, ablation studies are performed, focusing on various hyper-parameters in the IMF-based compression method and their effect on the compression performance.

Patchification. In Fig.7a, patchification effect with different patch sizes is investigated. First, it can be concluded that patchification has a positive effect on performance. This performance boost is

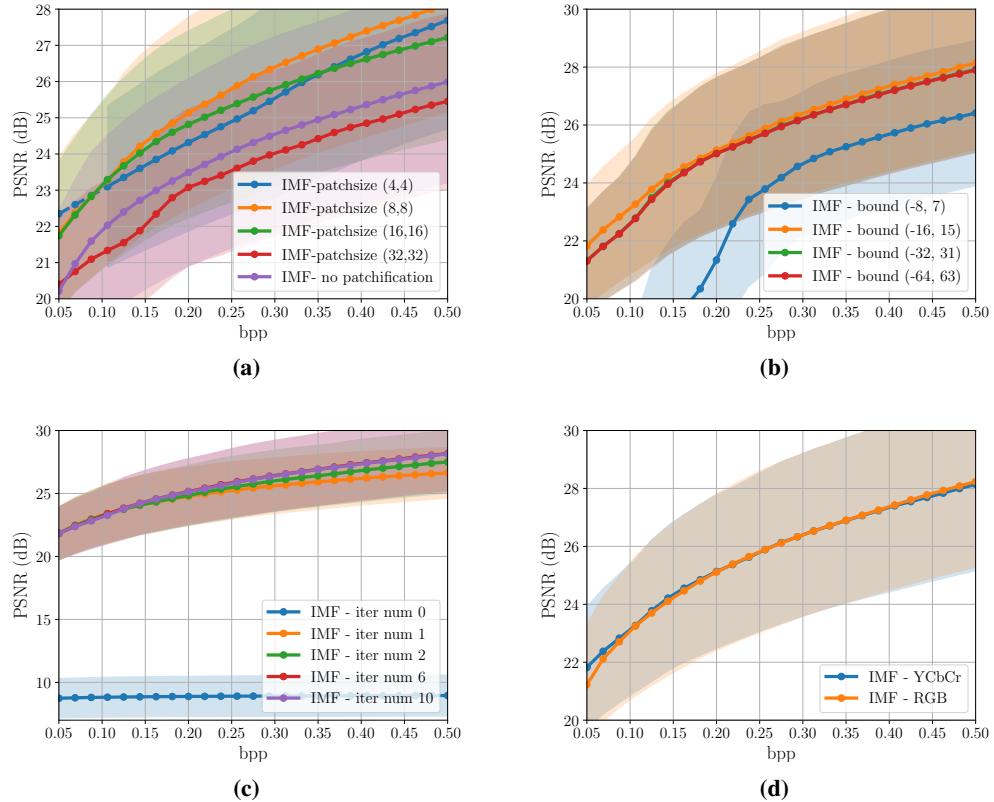


Figure 7 Ablation experiments for the IMF compression method. In all cases, we plot PSNR as a function of bits per pixel (bpp) on the Kodak dataset. (a) Compares IMF compression performance without patchification and different patch sizes. (b) Compares IMF compression performance for different bound values of factor matrices. (c) Compares IMF compression performance for different numbers of BCD iterations. (d) Compares IMF compression performance between RGB and YCbCr color space transform.

mainly due to the fact that in each patch (local) pattern variation is limited and hence IMF has more representation power to reconstruct the original pattern in each patch with fewer components. The performance on various datasets has shown that patches of size 8×8 lead to the best performance. The same conclusion is evident for the Kodak dataset example presented in Fig. 7a.

Factor bounds. As elaborated in Section 3.4, during the IMF optimization, IMF components are constrained into a bound $(-\alpha, \alpha - 1)$. Fig. 7b studies the compression performance versus different bounds. According to the results, the bound $(-16, 15)$ leads to the best performance since the value distribution of IMF components lies mostly in this bound. Hence, dedicating fewer bits to represent this narrower bound compared to the other bounds results in higher compression rates without sacrificing performance.

BCD iteration. The next parameter to study is the inner iteration number required in IMF BCD updates. According to the numerical results on various datasets, the objective value in the IMF optimization drops drastically after 2 iterations, while more iteration numbers have marginal improvement. This observation is shown for the Kodak dataset in Fig. 7c. This feature makes the IMF computationally efficient since with a limited number of iterations a high compression performance can be achieved.

Color space. The compression performance of IMF is studied with two color spaces, namely RGB and YCbCr, in Fig. 7d. Although the compression performance remains unchanged in terms of PSNR, qualitative results reported in Fig. 7 indicate that YCbCr color space can maintain natural colors and brightness more effectively. Consistently, the JPEG method employs this color space as well.

5 Conclusion and Future Work

Acknowledgments and Disclosure of Funding

References

- [1] Bo Dong, Matthew M Lin, and Haesun Park. Integer matrix approximation and data mining. *Journal of scientific computing*, 75:198–224, 2018.
- [2] Carl Eckart and Gale Young. The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3):211–218, 1936.
- [3] Nicholas Gillis. *Nonnegative Matrix Factorization*. SIAM, 2020.
- [4] Daniel Lee and H. Sebastian Seung. Algorithms for Non-negative Matrix Factorization. In T. Leen, T. Dietterich, and V. Tresp, editors, *Advances in Neural Information Processing Systems*, volume 13. MIT Press, 2000. URL https://proceedings.neurips.cc/paper_files/paper/2000/file/f9d1152547c0bde01830b7e8bd60024c-Paper.pdf.
- [5] Matthew M Lin, Bo Dong, and Moody T Chu. Integer matrix factorization and its application.
- [6] Peter van Emde Boas. Another np-complete problem and the complexity of computing short vectors in a lattice. *Tecnical Report, Department of Mathmatics, University of Amsterdam*, 1981.

A Proof of Theorem ??

Theorem 2 (Global convergence). *Let $(\mathbf{U}^k)_{k \in \mathbb{N}}$ and $(\mathbf{V}^k)_{k \in \mathbb{N}}$ be sequences generated by the proposed Algorithm ???. Then both sequences are convergent to a locally optimal point of the optimization problem (3).*

Proof. To study the convergence of the proposed Algorithm ???, we recast the optimization problem (3) to the following equivalent problem:

$$\begin{aligned} & \underset{\mathbf{U} \in \mathbb{R}^{M \times R}, \mathbf{V} \in \mathbb{R}^{N \times R}}{\text{minimize}} \quad \Psi(\mathbf{U}, \mathbf{V}) := H(\mathbf{U}, \mathbf{V}) + f(\mathbf{U}) + g(\mathbf{V}), \\ & \text{where} \quad H(\mathbf{U}, \mathbf{V}) := \|\mathbf{X} - \mathbf{U}\mathbf{V}^T\|_F^2, \\ & \quad f(\mathbf{U}) := \delta_{[a,b]}(\mathbf{U}) + \delta_{\mathbb{Z}}(\mathbf{U}), \\ & \quad g(\mathbf{V}) := \delta_{[a,b]}(\mathbf{V}) + \delta_{\mathbb{Z}}(\mathbf{V}), \end{aligned} \quad (12)$$

with $\delta_{\mathcal{B}}(\cdot)$ as the indicator function of the nonempty set \mathcal{B} where $\delta_{\mathcal{B}}(\mathbf{x}) = 0$ if $\mathbf{x} \in \mathcal{B}$ and $\delta_{\mathcal{B}}(\mathbf{x}) = +\infty$, otherwise. By the definition of functions above, it is easy to confirm that the problems (3) and (12) are equivalent.

The unconstrained optimization problem (12) consists of the sum of a differentiable (smooth) convex function H with nonsmooth nonconvex functions f and g . This problem has been extensively studied in the literature under the class of unconstrained nonconvex nonsmooth minimization problems. One of the common algorithms applied to such a problem class is the well-known forward-backward-splitting (FBS) []. In Algorithm ??, the variables \mathbf{U} and \mathbf{V} are updated sequentially following block coordinate (BC) descent minimization algorithms, also often called Gauss-Seidel updates or alternating minimization. Hence, in this convergence study we are interested in the algorithms that allow BC-type updates for the nonconvex nonsmooth problem of (12) []. Specifically, we focus on the proximal alternating linearized minimization (PALM) algorithm, relating its convergence behavior to that of Algorithm ???. To that end, we show that the updates of Algorithm ?? are equivalent to the updates of PALM on the recast problem of (12), and all the assumptions necessary for the convergence of PALM is satisfied by our problem setting.

The PALM algorithm [] can be summarized as follows:

1. Initialize $\mathbf{U}^0 \in \mathbb{R}^{M \times R}, \mathbf{V}^0 \in \mathbb{R}^{N \times R}$
2. For each iteration $k = 0, 1, \dots$

$$\begin{aligned} (a) \quad & \mathbf{U}^{k+1} \in \text{prox}_{c_k}^f \left(\mathbf{U}^k - \frac{1}{c_k} \nabla_{\mathbf{U}} H(\mathbf{U}^k, \mathbf{V}^k) \right), \text{ with } c_k > L_1(\mathbf{V}^k) \\ (b) \quad & \mathbf{V}^{k+1} \in \text{prox}_{d_k}^g \left(\mathbf{V}^k - \frac{1}{d_k} \nabla_{\mathbf{V}} H(\mathbf{U}^{k+1}, \mathbf{V}^k) \right), \text{ with } d_k > L_2(\mathbf{U}^{k+1}) \end{aligned} \quad (13)$$

where the proximal map for an extended proper lower semicontinuous (nonsmooth) function $\varphi : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ and $\gamma > 0$ is defined as $\text{prox}_{\gamma}^{\varphi}(\mathbf{x}) := \arg \min_{\mathbf{w} \in \mathbb{R}^n} \{ \varphi(\mathbf{w}) + \frac{\gamma}{2} \|\mathbf{w} - \mathbf{x}\|_2^2 \}$, ∇_x is the partial derivative with respect to x , and $L_1 > 0, L_2 > 0$ are local Lipschitz moduli, defined in the following assumption II.

The following proposition investigates the necessary assumptions for convergence of PALM and Algorithm ??.

Proposition 1 (Meeting required assumptions). *The assumptions necessary for the convergence of iterations in (13) are satisfied by the functions involved in the problem (12), specifically:*

1. *The indicator functions $\delta_{[a,b]}$ and $\delta_{\mathbb{Z}}$ are proper and lower semicontinuous functions, so do the functions f and g ;*
2. *For any fixed \mathbf{V} , the partial gradient $\nabla_{\mathbf{U}} H(\mathbf{U}, \mathbf{V})$ is globally Lipschitz continuous with modulus $L_1(\mathbf{V}) = \|\mathbf{V}^T \mathbf{V}\|_F$ defined by*

$$\|\nabla_{\mathbf{U}} H(\mathbf{U}_1, \mathbf{V}) - \nabla_{\mathbf{U}} H(\mathbf{U}_2, \mathbf{V})\| \leq L_1(\mathbf{V}) \|\mathbf{U}_1 - \mathbf{U}_2\|, \quad \forall \mathbf{U}_1, \mathbf{U}_2 \in \mathbb{R}^{M \times R},$$

where $\|\cdot\|$ in this section denotes the ℓ_2 -norm of the vectorized input with the proper dimension (here, with the input in $\mathbb{R}^{M \times R}$). The similar Lipschitz continuity is evident for $\nabla_{\mathbf{V}} H(\mathbf{U}, \mathbf{V})$ as well with modulus $L_2(\mathbf{U}) = \|\mathbf{U}\mathbf{U}^T\|_F$.

3. The sequences \mathbf{U}^k and \mathbf{V}^k are bounded due to the indicator functions $\delta_{[a,b]}$ with bounded a and b . Hence the moduli $L_1(\mathbf{V}^k)$ and $L_2(\mathbf{U}^k)$ are bounded from below and from above for all $k \in \mathbb{N}$.
4. The function H is twice differentiable, hence, its full gradient $\nabla H(\mathbf{U}, \mathbf{V})$ is Lipschitz continuous on the bounded set $\mathbf{U} \in [a, b]^{M \times R}$, $\mathbf{V} \in [a, b]^{N \times R}$. Namely, with $M > 0$:

$$\begin{aligned} \|(\nabla_{\mathbf{U}} H(\mathbf{U}_1, \mathbf{V}_1) - \nabla_{\mathbf{U}} H(\mathbf{U}_2, \mathbf{V}_2), \nabla_{\mathbf{V}} H(\mathbf{U}_1, \mathbf{V}_1) - \nabla_{\mathbf{V}} H(\mathbf{U}_2, \mathbf{V}_2))\| \\ \leq M \|(\mathbf{U}_1 - \mathbf{U}_2, \mathbf{V}_1 - \mathbf{V}_2)\|, \end{aligned} \quad (14)$$

where (\cdot, \cdot) denotes the concatenation of the two arguments.

5. The sets $[a, b]$ and integer numbers are semi-algebraic; so are their indicator functions. The function H is also polynomial, hence it is semi-algebraic. The sum of these functions results in a semi-algebraic function Ψ in (12), hence Ψ is a KL function.

By Proposition 1, the optimization problem (12)—and equivalently the original problem (3)—can be solved by PALM, which is convergence to a locally optimal point, specifically a critical point $(\mathbf{U}^*, \mathbf{V}^*)$ where $0 \in \partial\Psi(\mathbf{U}^*, \mathbf{V}^*)$, with ∂ as the subdifferential of Ψ .

In the following, we highlight that the updates in () can be implemented more simply and more efficiently by Algorithm ?? for the problem of image compression. It is noted that the so-called *forward* steps $\mathbf{U}^k - \frac{1}{c_k} \nabla_{\mathbf{U}} H(\mathbf{U}^k, \mathbf{V}^k)$ and $\mathbf{V}^k - \frac{1}{d_k} \nabla_{\mathbf{V}} H(\mathbf{U}^{k+1}, \mathbf{V}^k)$ in the prox operators can be replaced by the simple closed-form updates of Algorithm ?? in steps () and (), respectively. This is thanks to the special form of the functions $H(\cdot, \mathbf{V}^k)$ and $H(\mathbf{U}^{k+1}, \cdot)$ being quadratic functions, each having a global optimal point which ensures a descent in each forward step. Furthermore, the proximal operators $\text{prox}_{c_k}^f$ and $\text{prox}_{d_k}^g$ can efficiently be implemented by the steps () and () in Algorithm ???. The equivalence of these steps is proven in the following lemma.

Lemma 1 (prox implementation). Define the following (elementwise) operators $P_{[a,b]} : \mathbb{R} \rightarrow \mathbb{R}$, and $T_{\mathbb{Z}} : \mathbb{R} \rightarrow \mathbb{R}$:

$$P_{[a,b]}(x) := \min\{\max\{a, x\}, b\}, \quad (15)$$

$$T_{\mathbb{Z}}(x) := \max\{\lfloor x \rfloor, \lfloor x + 0.5 \rfloor\}. \quad (16)$$

Then $\text{prox}_{c_k}^f(\mathbf{W}) = T_{\mathbb{Z}}(P_{[a,b]}(\mathbf{W}))$ and $\text{prox}_{d_k}^g(\mathbf{Z}) = T_{\mathbb{Z}}(P_{[a,b]}(\mathbf{Z}))$ for any $\mathbf{W} \in \mathbb{R}^{M \times R}$ and $\mathbf{Z} \in \mathbb{R}^{N \times R}$ and $T_{\mathbb{Z}}(P_{[a,b]}(\cdot))$ being an elementwise operator on the input matrices.

Proof. Define the following norms for a given matrix $\mathbf{W} \in \mathbb{R}^{M \times R}$:

$$\|\mathbf{W}\|_{[a,b]} := \sum_{i,j | a \leq \mathbf{W}_{ij} \leq b} \mathbf{W}_{ij}^2, \quad \|\mathbf{W}\|_a := \sum_{i,j | \mathbf{W}_{ij} < a} \mathbf{W}_{ij}^2, \quad \|\mathbf{W}\|_b := \sum_{i,j | \mathbf{W}_{ij} > b} \mathbf{W}_{ij}^2.$$

Moreover, note that the rounding operator in (16) can be equivalently driven by the following proximal operator:

$$T_{\mathbb{Z}}(\mathbf{W}) = \arg \min_{\mathbf{U}} \{\|\mathbf{U} - \mathbf{W}\|_F^2 \mid \mathbf{U} \in \mathbb{Z}^{M \times R}\}. \quad (17)$$

The proximal operator can $\text{prox}_{c_k}^f(\mathbf{W})$ can be rewritten as

$$\begin{aligned}
\text{prox}_{c_k}^f(\mathbf{W}) &= \arg \min_{\mathbf{U}} \{\delta_{[a,b]}(\mathbf{U}) + \delta_{\mathbb{Z}}(\mathbf{U}) + \frac{c_k}{2} \|\mathbf{U} - \mathbf{W}\|_F^2\} \\
&= \arg \min_{\mathbf{U}} \{\|\mathbf{U} - \mathbf{W}\|_F^2 \mid \mathbf{U} \in \mathbb{Z}_{[a,b]}^{M \times R}\} \\
&= \arg \min_{\mathbf{U}} \{\|\mathbf{U} - \mathbf{W}\|_{[a,b]}^2 + \|\mathbf{U} - \mathbf{A}\|_a^2 + \|\mathbf{U} - \mathbf{B}\|_b^2 \mid \mathbf{U} \in \mathbb{Z}_{[a,b]}^{M \times R}\} \\
&= \arg \min_{\mathbf{U}} \{\|\mathbf{U} - \mathbf{W}\|_{[a,b]}^2 + \|\mathbf{U} - \mathbf{A}\|_a^2 + \|\mathbf{U} - \mathbf{B}\|_b^2 \mid \mathbf{U} \in \mathbb{Z}^{M \times R}\} \\
&= \arg \min_{\mathbf{U}} \{\|\mathbf{U} - P_{[a,b]}(\mathbf{W})\|_F^2 \mid \mathbf{U} \in \mathbb{Z}^{M \times R}\} \\
&= T_{\mathbb{Z}}(P_{[a,b]}(\mathbf{W})).
\end{aligned}$$

The first equality is due to the definition of prox which is equivalent to the second equality. In the third equality the matrices $A \in \mathbb{R}^{M \times R}$ and $B \in \mathbb{R}^{M \times R}$ have elements all equal to a and b , respectively. The third equality is due to the fact that replacing $\|\mathbf{U} - \mathbf{W}\|_a^2 + \|\mathbf{U} - \mathbf{W}\|_b^2$ with $\|\mathbf{U} - \mathbf{A}\|_a^2 + \|\mathbf{U} - \mathbf{B}\|_b^2$ has no effect on the solution of the minimization. The fourth equality is also trivial due to the involved norms in the third equality. The fifth equality can be easily confirmed by the definition of the projection $P_{[a,b]}$. Finally, in the last equality (17) is invoked. A similar proof can be trivially followed for $\text{prox}_{d_k}^f(\mathbf{Z}) = T_{\mathbb{Z}}(P_{[a,b]}(\mathbf{Z}))$ as well. \square

Now that the equivalence of updates () with the simple and closed-form steps in Algorithm ?? is fully established, and the assumptions required for the convergence are verified to be met by the problem (12) equivalent to the original problem (3), the following proposition can be invoked to prove the convergence of Algorithm ??.

Proposition 2 (Global convergence []). *Let $((\mathbf{U}^k, \mathbf{V}^k))_{k \in \mathbb{N}}$ be a sequence generated by Algorithm ??, then the sequence converges to a locally optimal point of problem (3) which is a critical point $(\mathbf{U}^*, \mathbf{V}^*)$ of the problem (12).*

\square