

Quantization-aware Matrix Factorization for Low Bit Rate Image Compression

Pooya Ashtari^{a,1,*}, Pourya Behmandpoor^{a,1}, Fateme Nateghi Haredasht^b, Jonathan H. Chen^b, Panagiotis Patrinos^a, Sabine Van Huffel^a

^a*Department of Electrical Engineering (ESAT), STADIUS Center, KU Leuven, Belgium*

^b*Stanford Center for Biomedical Informatics Research, Stanford University, CA, USA*

Abstract

Lossy image compression is essential for efficient transmission and storage. Traditional compression methods mainly rely on discrete cosine transform (DCT) or singular value decomposition (SVD), both of which represent image data in continuous domains and, therefore, necessitate carefully designed quantizers. Notably, these methods consider quantization as a separate step, which prevents quantization errors from being incorporated into the compression process and degrades the reconstruction quality, particularly in SVD-based methods. To address this issue, we introduce a quantization-aware matrix factorization (QMF) to develop a novel lossy image compression method. QMF provides a low-rank representation of the image data as a product of two smaller matrices, with elements constrained to bounded integer values, thereby effectively integrating quantization with low-rank approximation. We propose an efficient, provably convergent iterative algorithm for QMF using a block coordinate descent scheme, with subproblems having closed-form solutions. Our experiments demonstrate that our method consistently outperforms JPEG at low bit rates below 0.25 bits per pixel and remains comparable at higher bit rates. We also demonstrated that our method has an improved capability to preserve visual semantics compared to JPEG at low bit rates by evaluating an ImageNet pre-trained classifier on compressed images. The project is available at <https://github.com/pashtari/lrf>.

Keywords: Matrix Factorization, Low-rank Approximation, Quantization, Image Compression.

1. Introduction

Lossy image compression involves reducing the storage size of digital images by discarding some image data that are redundant or less perceptible to the human eye. This is crucial for efficiently storing and transmitting images, particularly in applications where bandwidth or storage resources are limited, such as web browsing, streaming, and mobile platforms. Lossy image compression methods enable adjusting the degree of compression, providing a selectable tradeoff between storage size and image quality. Widely used methods such as JPEG [1] and JPEG 2000

*Corresponding author: Pooya Ashtari (pooya.ashtari@esat.kuleuven.be).

¹Pooya Ashtari and Pourya Behmandpoor contributed equally to this work.

[2] follow the *transform coding* paradigm [3]. They use orthogonal linear transformations, such as discrete cosine transform (DCT) [4] and discrete wavelet transform (DWT) [5], to decorrelate small image blocks. Since these transforms map image data into a continuous domain, quantization is necessary before coding into bytes. Unfortunately, as quantization errors can significantly degrade compression performance, the quantizers must be carefully crafted to minimize this impact, which further complicates codec design.

Another promising paradigm relies on low-rank approximation techniques, with singular value decomposition (SVD) being a notable example. SVD is recognized as the deterministically optimal transform for energy compaction [6]. In practice, current SVD-based methods [6, 7, 8] can represent image data only with factors that contain floating-point elements, necessitating a quantization step prior to any byte-level processing. The quantization step often introduces errors, which result in suboptimal compression performance.

Motivated by this, we introduce quantization-aware matrix factorization (QMF) and, based on it, develop an effective lossy image compression method. Unlike traditional compression methods, the proposed approach integrates quantization into the optimization process rather than treating it as a separate step before byte-level processing. Our QMF formulation provides a low-rank representation of image data as the product of two smaller factor matrices. The quantization is integrated via introducing constraints in the optimization process, where the elements of the factor matrices are constrained to *bounded integer* values. These elements, with discrete values represented as bounded integers, can be directly stored using standard integral data types—such as `int8` and `int16` supported by programming languages—and losslessly processed, making QMF arguably better suited than SVD for image compression. Another advantage of QMF is that the reshaped factor matrices can be treated as 8-bit grayscale images, allowing any lossless image compression standard to be seamlessly integrated into the proposed framework. We propose an efficient iterative algorithm for QMF using a block coordinate descent (BCD) scheme, where each column of a factor matrix is taken as a block and updated one at a time using a closed-form solution.

Our contributions are summarized as follows. We propose a novel optimization framework that enables the integration of quantization and low-rank approximation for image compression. Moreover, we introduce an efficient algorithm for the QMF problem and prove its convergence. Finally, to the best of our knowledge, this work is the first effort to explore QMF for image compression, presenting the first algorithm based on a low-rank approach that significantly outperforms SVD and competes favorably with JPEG, particularly at low bit rates. Our method narrows the gap between factorization and quantization by integrating them into a single layer and optimizing the compression system.

2. Related Work

Transform Coding. Transform coding is a widely used approach in lossy image compression, leveraging mathematical transforms to decorrelate pixel values and represent image data more compactly. One of the earliest and most influential methods is the discrete cosine transform (DCT) [4], used in JPEG [1], which converts image data into the frequency domain, prioritizing lower frequencies to retain perceptually significant information. The discrete wavelet transform (DWT) [5], used in JPEG 2000 [2], offers improved performance by capturing both frequency and location information, leading to better handling of edges and textures [9]. More recently,

the WebP [10] and HEIF [11, 12] formats combine DCT and intra-frame prediction to achieve superior compression and quality compared to JPEG.

Learned Image Compression (LIC). Recently, learned image compression (LIC) has gained attention for potentially outperforming traditional methods by leveraging deep neural networks. Ballé et al. [13] pioneered this area with an end-to-end trainable convolutional neural network based on variational autoencoders. Cheng et al. [14] incorporated a simplified attention module and discretized Gaussian mixture likelihoods for achieving a more accurate and flexible entropy model. Liu et al. [15] combined transformers and CNNs to exploit the local modeling ability of convolutions and the global modeling ability of the attention mechanism. Yang and Mandt [16] introduced diffusion models into LIC, using a denoising decoder to iteratively reconstruct a compressed image. Despite these advancements, the high computational complexity of LIC methods remains a significant limitation, particularly for real-time applications and resource-constrained environments.

Low-rank Techniques. Low-rank approximation can provide a compact representation by decomposing image data into smaller components. Notably, truncated singular value decomposition (tSVD) is a classical technique that decomposes images into singular values and vectors, retaining only the most significant components to achieve compression [6, 7]. Hou et al. [8] proposed sparse low-rank matrix approximation (SLRMA) for data compression, which is able to explore both the intra- and inter-coherence of data samples simultaneously from the perspective of optimization and transformation. More recently, Yuan and Haimi-Cohen [17] introduced a graph-based low-rank regularization to reduce compression artifacts near block boundaries at low bit rates.

Integer Matrix Factorization. There are applications where meaningful representation of data as discrete factor matrices is crucial. While typical low-rank techniques like SVD and nonnegative matrix factorization (NMF) are inappropriate for such applications, integer matrix factorization (IMF) ensures the integrality of factors to achieve this goal. Lin et al. [18] investigates IMF to effectively handle discrete data matrices for cluster analysis and pattern discovery. Dong et al. [19] introduce an alternative least squares method for IMF, verifying its effectiveness with some data mining applications. However, the application of IMF in image compression remains unexplored.

While existing IMF methods generally constrain factor elements to the entire set of integer values, we propose quantized matrix factorization (QMF), which minimizes the objective function over a bounded interval of integers, thereby modeling a uniformly quantized domain. Furthermore, we introduce a block coordinate descent (BCD)-based algorithm to solve the QMF problem, which is both computationally efficient and provably convergent. This work investigates the potential of QMF for image compression, arguing that it can serve as a powerful tool for this purpose.

3. Method

3.1. Overall Encoding Framework

The proposed compression method follows a *transform coding* paradigm, but it does not involve a separate quantization step. Figure 1 illustrates an overview of our encoding pipeline based on quantization-aware matrix factorization (QMF). The encoder accepts an RGB image with

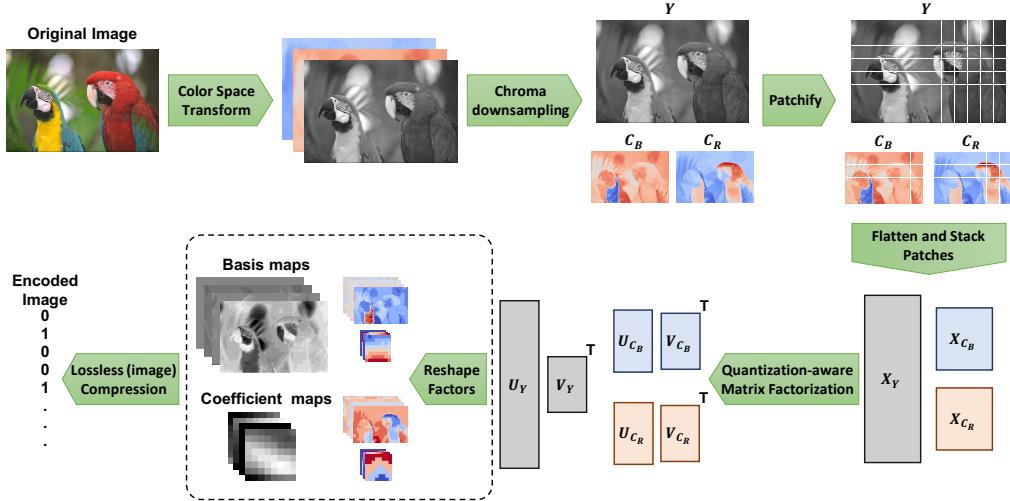


Figure 1: An illustration of the encoder for our image compression method.

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 YC_BC_R color space. Let $Y \in [0, 255]^{H \times W}$ represent the *luma* component, and $C_B, C_R \in [0, 255]^{\frac{H}{2} \times \frac{W}{2}}$ represent the blue-difference and red-difference *chroma* components, respectively. Note that as a result of this transformation, the elements of the *luma* (Y) and *chroma* (C_B, C_R) matrices are not limited to integers and can take any value within the interval $[0, 255]$.

Chroma Downsampling. After conversion to the YC_BC_R color space, the *chroma* components C_B and 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*).

Patchification. After *chroma* downsampling, we have three components: the *luma* component $Y \in [0, 255]^{H \times W}$ and the *chroma* components $C_B, C_R \in [0, 255]^{\frac{H}{2} \times \frac{W}{2}}$. 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 $X_Y \in [0, 255]^{\frac{HW}{64} \times 64}$, $X_{C_B} \in [0, 255]^{\frac{HW}{256} \times 64}$, and $X_{C_R} \in [0, 255]^{\frac{HW}{256} \times 64}$. Later, these matrices will be low-rank approximated using QMF. Note that this patchification technique differs from the block splitting in JPEG, where each block is subject to DCT individually and processed independently. This patchification technique not only captures the locality and spatial dependencies of neighboring pixels but also performs better when combined with the matrix decomposition approach for image compression.

Low-rank Approximation. We now apply a low-rank approximation to the matrices X_Y , X_{C_B} , and X_{C_R} , which is the core of our compression method that provides a lossy compressed representation.

tation of these matrices. The low-rank approximation [22] aims to approximate a given matrix $X \in \mathbb{R}^{M \times N}$ by

$$X \approx UV^\top = \sum_{r=1}^R U_{:,r} V_{:,r}^\top, \quad (1)$$

where $U \in \mathbb{R}^{M \times R}$ and $V \in \mathbb{R}^{N \times R}$ are *factor matrices* (or simply *factors*), $R \leq \min(M, N)$ represents the *rank*, $U_{:,r}$ and $V_{:,r}$ represent the r -th columns of U and V , respectively. We refer to U as the *basis matrix* and V as the *coefficient matrix*. By selecting a sufficiently small value for R , the factor matrices U and V , with a combined total of $(M + N)R$ elements, offer a compact representation of the original matrix 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 X and the product UV^\top , as well as the constraints on the factor matrices U and V , various formulations and variants have been proposed for different purposes [23, 24, 18]. In Section 3.3, we introduce and elaborate on our variant, termed quantization-aware matrix factorization (QMF), and argue why it is well-suited and effective for image compression.

Lossless Compression. QMF yields factor matrices $U_Y \in \{0, \dots, 255\}^{\frac{HW}{64} \times R}$ and $V_Y \in \{0, \dots, 255\}^{64 \times R}$; $U_{C_B} \in \{0, \dots, 255\}^{\frac{HW}{256} \times R}$ and $V_{C_B} \in \{0, \dots, 255\}^{64 \times R}$; and $U_{C_R} \in \{0, \dots, 255\}^{\frac{HW}{256} \times R}$ and $V_{C_R} \in \{0, \dots, 255\}^{64 \times R}$ that correspond to X_Y , X_{C_B} , and X_{C_R} . Since these matrices have elements constrained to integer values (allowing seamless integration of quantization with their optimization process), they can be directly encoded using any standard lossless data compression method, such as zlib [25]. In contrast, other lossy image compression methods typically require a separate quantization step, introducing errors that cannot be incorporated or considered during the compression process.

Alternatively, we can first reshape the factor matrices by unfolding their first dimension to obtain R -channel 2D spatial maps, referred to as *factor maps* and represented by the following tensors:

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

As each channel of a *factor map* can be treated as an 8-bit grayscale image, we can encode it by any standard lossless image compression method such as PNG. 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 QMF *basis maps* turn out to be meaningful images, each capturing some visual semantic of the image (see Figure 2 for an example). Therefore, our QMF approach can effectively leverage the power of existing lossless image compression algorithms, offering a significant advantage over current methods. However, in this work, we take the first approach and use the zlib library [25] to encode factor matrices, creating a stand-alone codec that is independent from other image compression methods.

3.2. Decoding

The decoder receives an encoded image and reconstructs the RGB image by applying the inverse of the operations used by the encoder, starting from the last layer and moving to the first. Initially, the factor matrices are produced by losslessly decompressing the encoded image. The matrices

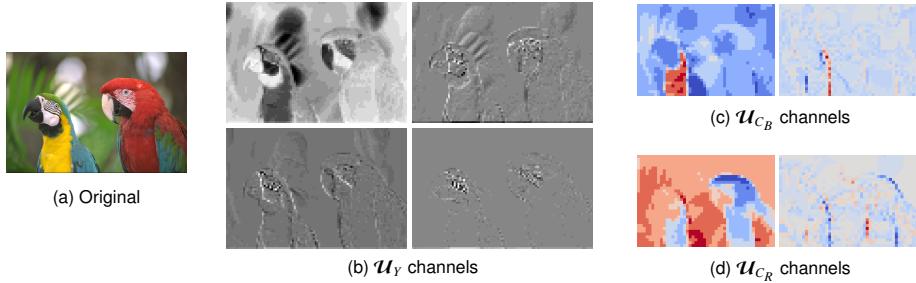


Figure 2: The channels of QMF basis maps for the `kodim23` image from Kodak. (a) shows the original image. The QMF basis maps corresponding to luma (b), blue-difference (c), and red-difference chroma (d) are shown. The channels of basis map with higher energy maintain the overall texture of the original image, where channels with lower energy focus more on subtle changes.

X_Y , X_{C_B} , and X_{C_R} are calculated through the product of the corresponding factor matrices, according to (1). The *luma* and downsampled *chroma* components are then obtained by reshaping X_Y , X_{C_B} , and X_{C_R} back into their spatial forms, following the inverse of the patchification step. Subsequently, the downsampled *chroma* components are upsampled to their original size using nearest-neighbor interpolation. Finally, the YC_BC_R image is converted back into an RGB image.

3.3. Quantization-aware Matrix Factorization (QMF)

The main building block of our method is quantization-aware matrix factorization (QMF), which is responsible for the lossy compression of matrices obtained through patchification. QMF 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{UV}^\top , while ensuring, as an integrated quantization step, that the elements of the factor matrices \mathbf{U} and \mathbf{V} are integers within a specified interval $[\alpha, \beta]$ with integer endpoints, i.e., $\alpha, \beta \in \mathbb{Z}$. Formally, the QMF problem can be expressed as:

$$\begin{aligned} & \underset{\mathbf{U}, \mathbf{V}}{\text{minimize}} && \|\mathbf{X} - \mathbf{UV}^\top\|_F^2 \\ & \text{subject to} && \mathbf{U} \in \mathbb{Z}_{[\alpha, \beta]}^{M \times R}, \mathbf{V} \in \mathbb{Z}_{[\alpha, \beta]}^{N \times R} \end{aligned} \quad (3)$$

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 singular value decomposition (SVD), as addressed by the Eckart–Young–Mirsky theorem [22]. If only a nonnegativity constraint were applied (without integrality), variations of nonnegative matrix factorization (NMF) would emerge [23, 26]. The QMF problem (3) poses a challenging integer program, with finding its global minima known to be NP-hard [19, 27]. Only a few iterative algorithms [19, 18] have been proposed to find a “good solution” for some QMF variants in contexts other than image compression. In Section 3.4, we propose an efficient iterative algorithm for the QMF problem (3).

The existing lossy image compression methods based on SVD and NMF approach the problem as an optimization task, followed by a separate quantization step. The optimization focuses on finding factors that minimize the reconstruction error. Before any byte-level processing, a quantization step is applied to project the floating-point elements of the resulting factors onto a set of discrete values. However, because the quantization is performed separately from the

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

Input: $X \in \mathbb{R}^{M \times N}$, factorization rank R , factor bounds $[\alpha, \beta]$, # iterations K

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

```

1: Initialize  $U^{\text{init}}$ ,  $V^{\text{init}}$  using the truncated SVD method, provided by (8) and (9), and set  $k = 0$ 
2: while  $k < K$  do
3:    $k \leftarrow k + 1$ 
4:    $A \leftarrow X V^k$ 
5:    $B \leftarrow V^{k^\top} V^k$ 
6:   for  $r = 1, \dots, R$  do
7:      $U_{:,r}^{k+1/2} \leftarrow \frac{A_{:,r} - \sum_{s=1}^{r-1} B_{sr} U_{:,s}^{k+1} - \sum_{s=r+1}^M B_{sr} U_{:,s}^k}{\|V_{:,r}^k\|^2}$ 
8:      $U_{:,r}^{k+1} \leftarrow \text{clamp}_{[\alpha, \beta]}(\text{round}(U_{:,r}^{k+1/2}))$ 
9:   end for
10:   $A \leftarrow X^\top U^{k+1}$ 
11:   $B \leftarrow U^{k+1^\top} U^{k+1}$ 
12:  for  $r = 1, \dots, R$  do
13:     $V_{:,r}^{k+1/2} \leftarrow \frac{A_{:,r} - \sum_{s=1}^{r-1} B_{sr} V_{:,s}^{k+1} - \sum_{s=r+1}^N B_{sr} V_{:,s}^k}{\|U_{:,r}^{k+1}\|^2}$ 
14:     $V_{:,r}^{k+1} \leftarrow \text{clamp}_{[\alpha, \beta]}(\text{round}(V_{:,r}^{k+1/2}))$ 
15:  end for
16: end while
17: return  $(U^K, V^K)$ 

```

optimization, the quantization errors cannot be incorporated into the compression process. This separation leads to suboptimal compression performance (as demonstrated in Section 4) and additional complications in designing quantizers. In contrast, our QMF formulation, through the *constrained* optimization in (3), produces integer factor matrices that minimize the reconstruction error while ensuring that all elements are discrete. These integer factor matrices can be directly stored and processed losslessly without introducing roundoff errors. The reason for limiting the feasible region to $[\alpha, \beta]$ in our QMF 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 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 QMF formulation is well-suited for image compression, effectively integrating the factorization and quantization steps into a single, efficient compression process.

3.4. Block Coordinate Descent Scheme for QMF

We propose an efficient algorithm for QMF using the block coordinate descent (BCD) scheme (aka alternating optimization). The pseudocode is provided in Algorithm 1. 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 as 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:

$$\mathbf{u}_r \leftarrow \arg \min_{\mathbf{u}_r \in \mathbb{Z}_{[\alpha, \beta]}^{M \times 1}} \|\mathbf{E}_r - \mathbf{u}_r \mathbf{v}_r^\top\|_F^2, \quad (4)$$

$$\mathbf{v}_r \leftarrow \arg \min_{\mathbf{v}_r \in \mathbb{Z}_{[\alpha, \beta]}^{N \times 1}} \|\mathbf{E}_r - \mathbf{u}_r \mathbf{v}_r^\top\|_F^2, \quad (5)$$

where $\mathbf{u}_r \triangleq U_{:,r}$ and $\mathbf{v}_r \triangleq V_{:,r}$ represent the r -th columns of \mathbf{U} and \mathbf{V} , respectively. $\mathbf{E}_r \triangleq \mathbf{X} - \sum_{s \neq r}^R \mathbf{u}_s \mathbf{v}_s^\top$ 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 (Monotonicity). *The global optima of subproblems (4) and (5) can be represented by closed-form solutions as follows:*

$$\mathbf{u}_r \leftarrow \text{clamp}_{[\alpha, \beta]} \left(\text{round} \left(\frac{\mathbf{E}_r \mathbf{v}_r}{\|\mathbf{v}_r\|^2} \right) \right), \quad (6)$$

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

where $\text{round}(\mathbf{Z})$ denotes an element-wise 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 an element-wise operator that clamps each element of \mathbf{Z} to the interval $[\alpha, \beta]$. Moreover, the cost function in (3) is monotonically nonincreasing over BCD iterations of Algorithm 1, which involve sequential updates of (6) and (7) over columns of \mathbf{U} and \mathbf{V} .

Proof. See Appendix A for the proof. \square

It is noteworthy that the combination of $\text{round}(\cdot)$ and $\text{clamp}_{[\alpha, \beta]}(\cdot)$ in (6) and (7) can be interpreted as the element-wise projector to $\mathbb{Z}_{[\alpha, \beta]}$. In addition, updates (6) and (7) are presented in Algorithm 1 at steps 7 and 8, and steps 13 and 14, respectively. In Theorem 2, the convergence of Algorithm 1 employing these closed-form solutions is established.

Theorem 2 (Convergence). *Let $(\mathbf{U}_{:,r}^k)_{k \in \mathbb{N}}$ and $(\mathbf{V}_{:,r}^k)_{k \in \mathbb{N}}$ for $r \in \{1, \dots, R\}$ be sequences generated by the proposed Algorithm 1. Then all sequences are bounded and convergent to a locally optimal point of optimization problem (3).*

Proof. See Appendix B for the proof. \square

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. $\tilde{\mathbf{U}} \in \mathbb{R}^{M \times R}$ and $\tilde{\mathbf{V}} \in \mathbb{R}^{N \times R}$ contain the corresponding left-singular vectors and right-singular vectors in their

columns, respectively. The initial factors are then calculated as follows:

$$\mathbf{U}^{\text{init}} = \text{clamp}_{[\alpha, \beta]}(\text{round}(\tilde{\mathbf{U}}\Sigma^{\frac{1}{2}})), \quad (8)$$

$$\mathbf{V}^{\text{init}} = \text{clamp}_{[\alpha, \beta]}(\text{round}(\Sigma^{\frac{1}{2}}\tilde{\mathbf{V}})). \quad (9)$$

Essentially, this means that instead of performing a constrained optimization, we first low-rank approximate \mathbf{X} and then satisfy the constraints by projecting the elements of the resulting factor matrices onto $\mathbb{Z}_{[\alpha, \beta]}$.

4. Experiments

4.1. Setup

QMF Configuration. In our QMF implementation, we used a default patch size of 8×8 . The default factor bounds were set to $[-16, 15]$. Unless otherwise specified, the number of BCD iterations was set to 10 although our ablation studies in Section 4.6 suggest that even 2 iterations may suffice in practice (see Figure 6b). For lossless compression of factors, we encoded and decoded each column of a factor separately using the zlib library [25]. We also tested other lossless compression methods, such as zstd [28] and Huffman coding [29], which demonstrated a comparable performance. However, as zlib is well-established, simple, and offers fast performance in Python, the experimental results are reported for this compression method.

Baseline Codecs. We compared our QMF method against JPEG and SVD baselines. For JPEG compression, we used the Pillow library [30]. Our SVD baseline follows the same framework as the proposed method (described in 3.1) but substitutes truncated SVD for QMF. This is followed by uniform quantization of the SVD factor matrices before lossless compression using zlib [25]. This differs from QMF compression, which benefits from the integrality of factors by directly encoding them with zlib, eliminating the need for a separate quantization step.

Datasets. To validate the effectiveness of our method, we conducted experiments using the widely-used **Kodak** dataset [31], consisting of 24 lossless images with a resolution of 768×512 . To evaluate the robustness of our method in a higher-resolution setting, we also experimented with the **CLIC 2024** validation dataset [32], which contains 30 high-resolution, high-quality images. Additionally, we assessed the compression methods by their ability to retain visual semantics. This was achieved by evaluating a pre-trained ImageNet classifier on compressed images from the **ImageNet** validation set [33], consisting of 50,000 images with a resolution of 224×224 across 1,000 classes.

Metrics. To evaluate the rate-distortion performance of methods on the Kodak and CLIC 2024 datasets, we measured the bit rate in bits per pixel (bpp) and assessed the quality of the reconstructed images using peak signal-to-noise ratio (PSNR) and structural similarity index measure (SSIM). Then, these metrics were plotted as functions of bit rate for each method to illustrate their rate-distortion performance. To control the quality of the reconstructed images in QMF and SVD, similar to JPEG, we defined a quality factor $Q \in [0, 1]$, where 0 represents the highest compression and 1 represents no compression. To determine the factorization rank R in Algorithm 1, we used $R = \max \{\text{round}(Q \times \min\{M, N\}), 1\}$.

More precisely, to construct a rate-distortion curve for each method on each dataset, we evaluated various qualities Q for each image. For each quality, we first measured the PSNR/SSIM values

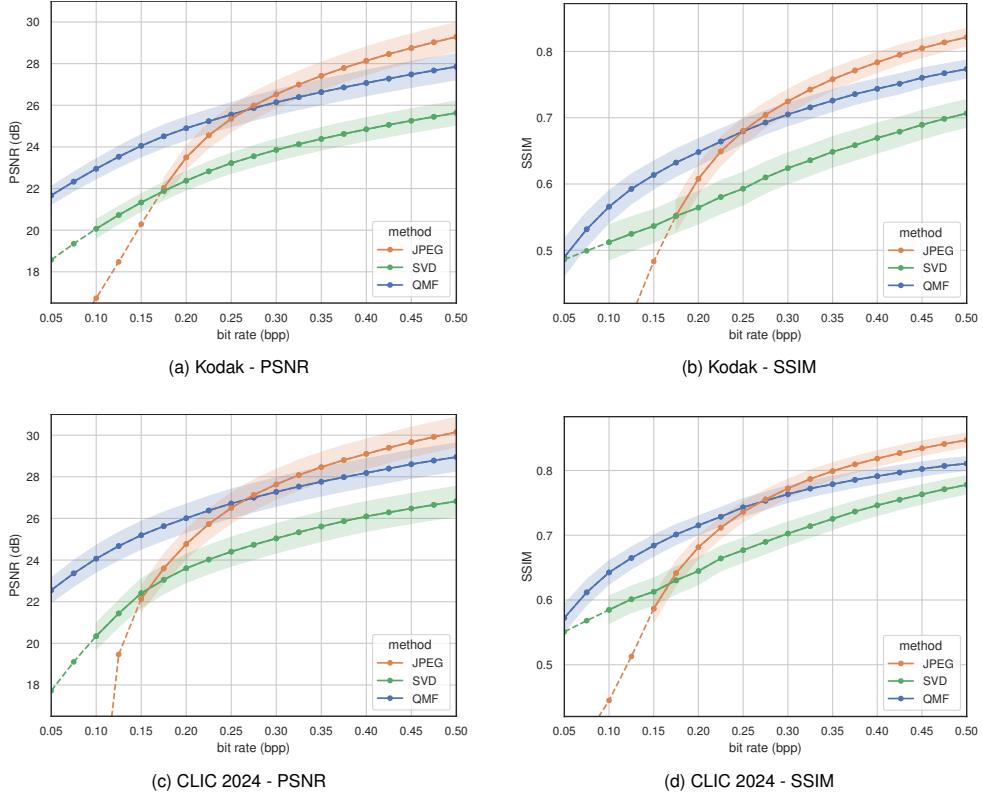


Figure 3: Rate-distortion performance on the Kodak (top panels) and CLIC 2024 (bottom panels) datasets. The average PSNR (left panels) and average SSIM (right panels) for each method are plotted as functions of bit rate. Shaded areas represent standard errors. Dashed lines indicate extrapolated values predicted using LOESS [34] for extremely low bit rates that are otherwise unattainable.

at the corresponding bit rate. Next, PSNR/SSIM values were interpolated at evenly spaced bit rates ranging from 0.05 bpp to 0.5 bpp using LOESS (locally estimated scatterplot smoothing) [34]. Finally, the interpolated values were averaged over all images at each of these bit rates.

4.2. Rate-Distortion Performance

Figure 3 illustrates the rate-distortion curves comparing the performance of QMF, SVD, and JPEG compression methods.

Kodak. On the Kodak dataset, as shown in Figures 3a and 3b, our QMF method consistently outperforms JPEG at low bit rates below 0.25 bpp and remains comparable at higher bit rates in terms of both PSNR and SSIM. Furthermore, QMF significantly surpasses the SVD-based baseline across all bit rates.

CLIC 2024. A similar trend is observed with the CLIC 2024 dataset, as shown in Figures 3c and 3d. Here, the PSNR (Figure 3c) and SSIM (Figure 3d) results further confirm the competitive performance of QMF across all bit rates, with a particularly notable margin at bit rates lower

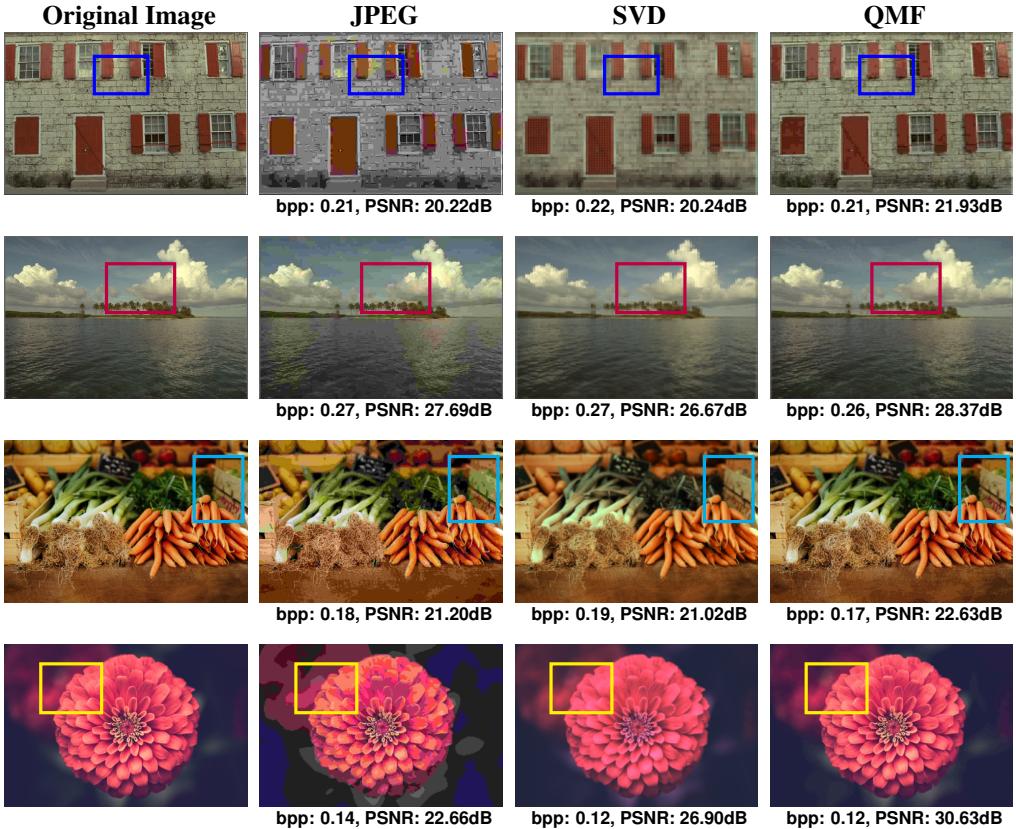


Figure 4: Qualitative performance comparison on example images from the Kodak (top two rows) and the CLIC 2024 (bottom two rows) datasets. Each column shows the original image, JPEG, SVD, and QMF compression results respectively. The bit rate and PSNR values for each compressed image is reported. The colored bounding boxes highlight artifacts produced by JPEG and SVD compression.

than 0.25 bpp. Specifically, at a bit rate of 0.15 bpp, QMF achieves an PSNR of over 25 dB, compared to approximately 22 dB for both JPEG and SVD. This supports the robustness of QMF in preserving visual quality across different datasets.

4.3. Qualitative Performance

Figure 4 compares various compression methods using images from the Kodak (top two rows) and CLIC 2024 (bottom two rows) datasets, compressed at similar bit rates.

In the building image (first row), JPEG compression, with a PSNR of 20.22 dB at a bit rate of 0.21 bpp, introduces *blocking artifacts* and changes the facade color, as visible in the blue boxes. SVD compression reduces these artifacts but causes blurriness. Our QMF compression, with a similar bit rate but a higher PSNR (21.93 dB), maintains both texture and sharpness with minimal artifacts.

In the seascape image (second row), JPEG causes blocking and significant *color bleeding artifacts*, such as the redness in the cloud area marked by the red boxes and also on the water surface

Table 1: Mean decoding CPU times for different compression methods at bit rates of 0.15 bpp and 0.25 bpp, measured on the Kodak and CLIC 2024 datasets.

| Method | Kodak | | CLIC 2024 | |
|--------|---------------------|---------------------|---------------------|---------------------|
| | Bit rate = 0.15 bpp | Bit rate = 0.25 bpp | Bit rate = 0.15 bpp | Bit rate = 0.25 bpp |
| JPEG | 4.54 ms | 4.23 ms | 26.76 ms | 25.75 ms |
| SVD | 1.33 ms | 1.23 ms | 5.29 ms | 4.82 ms |
| QMF | 2.82 ms | 2.66 ms | 9.91 ms | 9.06 ms |

(outside the red box). SVD reduces color distortion but still has blockiness and blurriness. QMF preserves the color and texture of clouds and water more effectively, resulting in a more visually pleasing image.

In the vegetables image (third row), JPEG yields visible *color distortion* (marked by the cyan boxes), while SVD introduces significant blurriness. QMF, however, effectively preserves the color fidelity and detail.

In the flower image (fourth row), JPEG compression, with a PSNR of 20.22 dB at a bit rate of 0.14 bpp, exhibits severe *color banding artifacts* around the flower boundary. SVD compression offers smoother gradients but remains blurry. Our QMF compression maintains the gradient fidelity and intricate petal distinctions, achieving a significantly higher PSNR of 30.63 dB at a lower bit rate of 0.12 bpp.

4.4. Run Time

The decoding times at bit rates of 0.15 bpp and 0.25 bpp for each method on Kodak and CLIC 2024 are reported in Table 1. All experiments in this section were conducted on 2 Xeon Gold 6140 CPUs @ 2.3 GHz (Skylake), each with 18 cores, and with 192 GiB RAM.

QMF and SVD have a significant advantage in decoding speed over JPEG, with SVD being the fastest. Specifically, QMF decodes more than twice as fast as JPEG on the CLIC 2024 dataset across all bit rates. This is due to the heavier FFT operation in the JPEG decoder compared to the lighter matrix multiplication in the QMF decoder. Overall, QMF is preferable for applications requiring high image quality at low bit rates, especially in scenarios where compressed images are frequently accessed or displayed while encoding occurs less often. Examples include web browsing, image hosting, mobile applications, satellite imagery for maps, and interactive gaming applications. These use cases often involve large numbers of images or thumbnails, where faster decoding times are essential for seamless user experiences, even if encoding takes more time or demands higher computational resources.

4.5. ImageNet Classification Performance

It is relevant to assess the ability of different compression methods in preserving the visual semantic information in images. To this end, we investigate the performance of an image classifier on images compressed using various compression methods. This is particularly crucial in scenarios where the ultimate goal is a vision task such as image classification, rather than maintaining perceived image quality, and we compress images before classification to minimize resource requirements, such as memory and communication bandwidth.

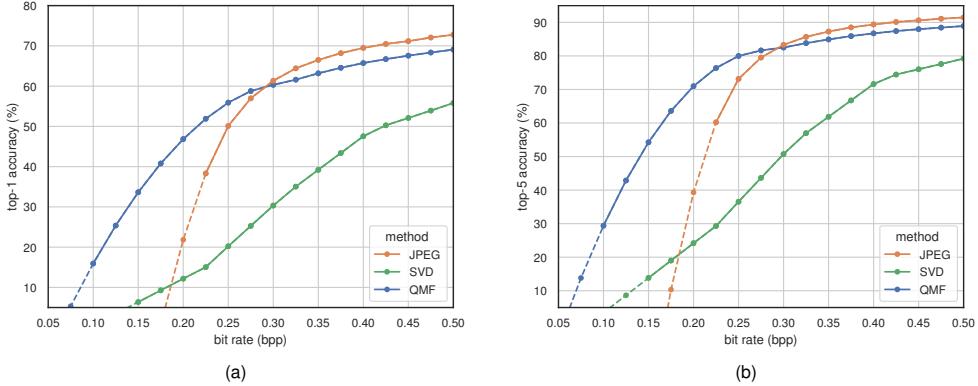


Figure 5: Impact of different compression methods on ImageNet classification accuracy. A ResNet-50 classifier pre-trained on the original ImageNet images is evaluated using validation images compressed by different methods. Panels (a) and (b) show top-1 and top-5 accuracy plotted as a function of bit rate, respectively. Dashed lines indicate extrapolated values predicted using LOESS [34] for extremely low bit rates that are otherwise unattainable.

In this experiment, we employed a ResNet-50 classifier [35], pre-trained on the original ImageNet [36] dataset, to classify compressed images from the ImageNet validation set using different compression methods. The classification performance comparison is presented in Figure 5. Notably, the results indicate that QMF compression achieves over a 5% improvement in top-1 accuracy compared to JPEG at bit rates under 0.25 bpp and reaches a top-5 accuracy exceeding 70% at a bit rate of 0.2 bpp. QMF compression leads to higher classification accuracies than JPEG at bit rates up to approximately 0.30 bpp.

4.6. Ablation Studies

We conducted ablation studies to investigate the impact of factor bounds, the number of BCD iterations, and patch size on the compression performance of our QMF method. All experiments in this section were performed using the Kodak dataset. We followed the QMF configuration described in Section 4.1 and varied only the parameters under ablation one at a time.

Factor Bounds. Figure 6a shows the average PSNR as a function of bit rate for QMF using various factor bounds $[\alpha, \beta]$ in Algorithm 1. The results indicate that the interval $[-16, 15]$ yields the optimal performance, showing moderate improvement over both $[-32, 31]$ and $[-128, 127]$, while significantly outperforming $[-8, 7]$. In fact, constraining the factor elements within a sufficiently narrow range can reduce the bit allocation needed, thereby leading to higher compression ratios. Note that in all these cases, the factor elements are represented as the `int8` data type.

BCD Iterations. The next parameter studied is the number of BCD iterations K in Algorithm 1, where each BCD iteration involves one complete cycle of updates across all the columns of both factors. Figure 6a shows the average PSNR plotted against the bit rate for QMF with different numbers of iterations $K \in \{0, 1, 2, 5, 10\}$. As expected, more iterations consistently resulted in higher PSNR for QMF compression. Without any BCD iterations ($K = 0$) and relying solely on the SVD-based initialization given by (8) and (9), the results became very poor. However, performance improved significantly after a few iterations, with more than $K = 5$ iterations yielding only marginal gains. We found that $K = 10$ iterations are sufficient in practice for image

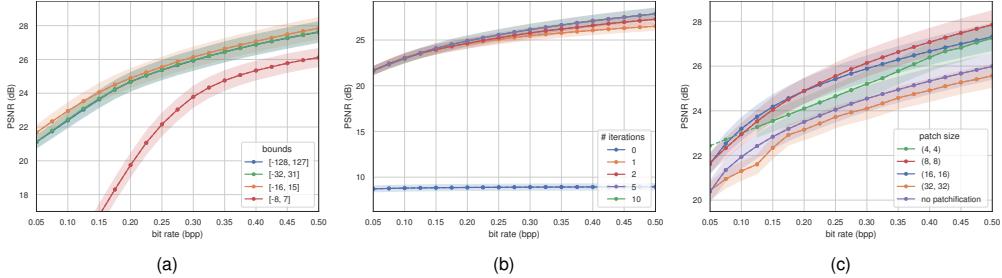


Figure 6: Ablation studies for QMF. The average PSNR on the Kodak dataset is plotted as a function of bit rate under various experimental conditions: (a) varying the bounds $[\alpha, \beta]$ for the elements of the factor matrices, (b) changing the number of BCD iterations, and (c) adjusting the patch size.

compression applications. This makes QMF computationally efficient, as decent compression performance can be achieved even with a limited number of BCD iterations.

Patchification. Figure 6c explores the impact of different patch sizes on QMF performance in terms of PSNR. As observed, a patch size of (8, 8) yields the best performance. A patch size of (16, 16) follows closely, with only marginally lower PSNR at higher bit rates. Conversely, larger patch sizes like (32, 32) or omitting the patchification step altogether significantly degrade compression performance.

5. Discussion

All our comparative results (Figure 3 and Figure 5) consistently show that our QMF method outperforms JPEG in both maintaining image quality and preserving visual semantics at low bit rates and remains comparable at higher bit rates. Moreover, QMF consistently demonstrates superior performance compared to SVD across all bit rates. This superiority can be attributed to the integration of quantization with low-rank approximation in QMF, which enables more accurate reconstruction. In contrast, the high sensitivity of SVD to quantization errors, which arise during a separate quantization step, significantly degrades the reconstruction quality.

As observed in Figure 6a, contracting the QMF factor bounds from $[-128, 127]$ to $[-16, 15]$ consistently improves the rate-distortion performance. Generally, narrowing the factor bounds $[\alpha, \beta]$ can potentially lower the entropy, thereby improving the effectiveness of lossless compression in the final stage of our framework and subsequently reducing the bit rate. However, this reduction in entropy comes at the cost of increased reconstruction error, as the feasible set in (3) becomes more constrained. This trade-off between entropy and reconstruction quality limits the compression performance of QMF. Therefore, it would be beneficial to moderately expand the factor bounds $[\alpha, \beta]$ while simultaneously controlling the entropy of the elements in the factor matrices. We plan to address this in the future by incorporating an entropy-aware regularization term into the current QMF objective function.

Patchification with an appropriate patch size (e.g., (8, 8)) helps capture local spatial dependencies and, as confirmed by our results in Figure 6c, positively impacts the performance of QMF and SVD. However, discontinuities at patch boundaries can introduce *blocking artifacts*, similar to JPEG compression at very low bit rates (see the building image example in Figure 4). Moreover, while JPEG suffers more from *color distortion* (e.g., *color bleeding* and *color banding*) at low bit

rates, QMF and SVD are more affected by *blurriness*, as observed in the seascape image example in Figure 4. As a potential solution for future work, a deep neural network could be trained to remove these artifacts and then integrated as a post-processing module to further enhance the quality of QMF-compressed images.

6. Conclusion

This work presents a novel lossy image compression method based on quantization-aware matrix factorization (QMF). By representing image data as the product of two smaller factor matrices with elements constrained to bounded integer values, the proposed QMF approach effectively integrates quantization with low-rank approximation. In contrast, traditional compression methods such as JPEG and SVD consider quantization as a separate step, where quantization errors cannot be incorporated into the compression process. The reshaped factor matrices in QMF are compatible with existing lossless compression standards, enhancing the overall flexibility and efficiency of our method. Our proposed iterative algorithm, utilizing a block coordinate descent scheme, has proven to be both efficient and convergent. Experimental results demonstrate that the QMF method significantly outperforms JPEG in terms of PSNR and SIMM at low bit rates and maintains better visual semantic information. This advantage underscores the potential of QMF to set a new standard in lossy image compression, bridging the gap between factorization and quantization.

Acknowledgments

This research received funding from Flanders AI Research Program. Sabine Van Huffel and Pooya Ashtari are affiliated to Leuven.AI - KU Leuven institute for AI, B-3000, Leuven, Belgium.

Appendix A. Proof of Theorem 1

We start by proving the closed-form solution (6), noting that the proof for (7) follows the same reasoning. The objective function in the subproblem (4) can be reformulated as follows:

$$\arg \min_{\mathbf{u}_r \in \mathbb{Z}_{[\alpha, \beta]}^M} \|\mathbf{E}_r - \mathbf{u}_r \mathbf{v}_r^\top\|_F = \arg \min_{\mathbf{u}_r \in \mathbb{Z}_{[\alpha, \beta]}^M} \sum_{i=1}^M \sum_{j=1}^N (e_{ij}^r - u_i^r v_j^r)^2, \quad (\text{A.1})$$

where e_{ij}^r denotes the element of matrix \mathbf{E}_r in the i th row and j th column, and u_i^r and v_j^r are the i th and j th elements of vectors \mathbf{u}_r and \mathbf{v}_r , respectively. Since the elements of \mathbf{E}_r and \mathbf{v}_r are fixed in problem (4), the optimization (A.1) can be decoupled into M optimizations as follows:

$$\arg \min_{u_i^r \in \mathbb{Z}_{[\alpha, \beta]}} q_i(u_i^r), \quad \text{where } q_i(u_i^r) \triangleq \sum_{j=1}^N (e_{ij}^r - u_i^r v_j^r)^2, \quad \forall i \in \{1, \dots, M\}. \quad (\text{A.2})$$

The objective functions $q_i(u_i^r)$ in (A.2) are single-variable quadratic problems. Hence, the global optimum in each decoupled optimization problem can be achieved by finding the minimum of each quadratic problem and then projecting it onto the set $\mathbb{Z}_{[\alpha, \beta]}$. The minimum of each quadratic function in (A.2), denoted by \bar{u}_i^r , can be simply found by

$$\nabla_{u_i^r} q_i(u_i^r) = 0 \implies \bar{u}_i^r = \sum_{j=1}^N e_{ij}^r v_j^r / \sum_{j=1}^N v_j^2, \quad (\text{A.3})$$

where ∇_x is the partial derivative with respect to x . Since q_i has a constant curvature (second derivative) and $q_i(\bar{u}_i^r + d)$ is nondecreasing with increasing $|d|$, the value in the set $\mathbb{Z}_{[\alpha, \beta]}$ which is closest to \bar{u}_i^r is the global minimizer of (A.2). This value can be reached by projecting \bar{u}_i^r onto the set $\mathbb{Z}_{[\alpha, \beta]}$, namely $u_i^{r*} = \text{clamp}_{[\alpha, \beta]}(\text{round}(\bar{u}_i^r))$, which is presented for all $i \in \{1, \dots, M\}$ in a compact form in (6).

Since $\mathbf{u}_r^* \triangleq (u_1^{r*}, \dots, u_M^{r*})$ is the global optimum of optimization (A.1), it is evident that

$$\|\mathbf{E}_r - \mathbf{u}_r^* \mathbf{v}_r^\top\|_F \leq \|\mathbf{E}_r - \mathbf{u}_r \mathbf{v}_r^\top\|_F. \quad (\text{A.4})$$

This inequality guarantees a nonincreasing cost function over one update of \mathbf{u}_r . Following the same reasoning for updates of \mathbf{v}_r in (7), it can be concluded that in each update of (6) and (7), the cost function is nonincreasing. Therefore, the sequential updates over the columns of \mathbf{U} and \mathbf{V} in Algorithm 1 result in a monotonically nonincreasing cost function in (3).

Appendix B. Proof of Theorem 2

To study the convergence of the proposed Algorithm 1, we recast the optimization problem (3) into the following equivalent problem:

$$\underset{\mathbf{U}_{:r} \in \mathbb{R}^M, \mathbf{V}_{:r} \in \mathbb{R}^N, \forall r \in \{1, \dots, R\}}{\text{minimize}} \quad \Psi(\mathbf{U}, \mathbf{V}) \quad (\text{B.1})$$

where

$$\begin{aligned} \Psi(\mathbf{U}, \mathbf{V}) &\triangleq f_0(\mathbf{U}, \mathbf{V}) + \sum_{r=1}^R f(\mathbf{U}_{:r}) + \sum_{r=1}^R g(\mathbf{V}_{:r}), \\ f_0(\mathbf{U}, \mathbf{V}) &\triangleq \|\mathbf{X} - \mathbf{UV}^\top\|_F^2, \\ f(\mathbf{U}_{:r}) &\triangleq \delta_{[a,b]}(\mathbf{U}_{:r}) + \delta_{\mathbb{Z}}(\mathbf{U}_{:r}), \\ g(\mathbf{V}_{:r}) &\triangleq \delta_{[a,b]}(\mathbf{V}_{:r}) + \delta_{\mathbb{Z}}(\mathbf{V}_{:r}), \end{aligned}$$

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 (B.1) are equivalent.

The unconstrained optimization problem (B.1) consists of the sum of a differentiable (smooth) function f_0 and nonsmooth functions f and g . This problem has been extensively studied in the literature under the class of nonconvex nonsmooth minimization problems. In Algorithm 1, the blocks $U_{:,r}$ and $V_{:,r}$ are updated sequentially following block coordinate descent (BCD) minimization algorithms, also often called Gauss-Seidel updates or alternating optimization [37, 38]. Hence, in this convergence study, we are interested in algorithms that allow BCD-like updates for the nonconvex nonsmooth problem of (B.1) [39, 40]. Specifically, we focus on the proximal alternating linearized minimization (PALM) algorithm [40], to relate its convergence behavior to that of Algorithm 1. To that end, we show that the updates of Algorithm 1 are related to the updates of PALM on the recast problem of (B.1), and all the assumptions necessary for the convergence of PALM are satisfied by our problem setting. It is noted that, for the sake of presentation and without loss of generality, in this proof, we assume each of the matrices U and V has only one column ($R = 1$); hence, we only have two blocks in the BCD updates. The iterates in PALM and the presented proof can be trivially extended for more than two blocks.

The PALM algorithm can be summarized as follows:

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

$$\begin{aligned} (a) \quad U^{k+1} &\in \text{prox}_{c_k}^f \left(U^k - \frac{1}{c_k} \nabla_U f_0(U^k, V^k) \right), \\ (b) \quad V^{k+1} &\in \text{prox}_{d_k}^g \left(V^k - \frac{1}{d_k} \nabla_V f_0(U^{k+1}, V^k) \right), \end{aligned} \quad (\text{B.2})$$

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 $\text{prox}_\gamma^\varphi(\mathbf{x}) \triangleq \arg \min_{\mathbf{w} \in \mathbb{R}^n} \{\varphi(\mathbf{w}) + \frac{\gamma}{2} \|\mathbf{w} - \mathbf{x}\|_2^2\}$. In (B.2), $c_k > L_1(V^k)$ and $d_k > L_2(U^{k+1})$ where $L_1 > 0$, $L_2 > 0$ are local Lipschitz moduli, defined in the following proposition.

The following proposition investigates the necessary assumptions (cf. [40, Asm. 1 and Asm. 2]) for convergence of iterates in (B.2).

Proposition 1 (Meeting required assumptions). *The assumptions necessary for the convergence of iterates in (B.2) are satisfied by the functions involved in the problem (B.1), 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 V , the partial gradient $\nabla_U f_0(U, V)$ is globally Lipschitz continuous with modulus $L_1(V) = \|V^T V\|_{\text{F}}$. Therefore, for all $U_1, U_2 \in \mathbb{R}^{M \times R}$ the following holds*

$$\|\nabla_U f_0(U_1, V) - \nabla_U f_0(U_2, V)\| \leq L_1(V) \|U_1 - U_2\|,$$

where $\|\cdot\|$ denotes the ℓ_2 -norm of the vectorized input with the proper dimension (here, with the input in $\mathbb{R}^{MR \times 1}$). The similar Lipschitz continuity is evident for $\nabla_V f_0(U, V)$ as well with modulus $L_2(U) = \|UU^T\|_{\text{F}}$.

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 f_0 is twice differentiable, hence, its full gradient $\nabla f_0(\mathbf{U}, \mathbf{V})$ 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 $M > 0$:

$$\|(\nabla_U f_0(\mathbf{U}_1, \mathbf{V}_1) - \nabla_U f_0(\mathbf{U}_2, \mathbf{V}_2), \nabla_V f_0(\mathbf{U}_1, \mathbf{V}_1) - \nabla_V f_0(\mathbf{U}_2, \mathbf{V}_2))\| \leq M \|(\mathbf{U}_1 - \mathbf{U}_2, \mathbf{V}_1 - \mathbf{V}_2)\|,$$

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 f_0 is also polynomial, hence it is semi-algebraic. The sum of these functions results in a semi-algebraic function Ψ in (B.1), hence Ψ is a Kurdyka-Łojasiewicz (KL) function.

By Proposition 1, the optimization problem (B.1) can be solved by the iterates in (B.2), due to the following proposition:

Proposition 2 (Global convergence [40]). *With the assumptions in proposition 1 being met by the problem (B.1), let $((\mathbf{U}^k, \mathbf{V}^k))_{k \in \mathbb{N}}$ be a sequence generated by the iterates in (B.2). Then the sequence converges to a critical point $(\mathbf{U}^*, \mathbf{V}^*)$ of the problem (B.1), where $0 \in \partial\Psi(\mathbf{U}^*, \mathbf{V}^*)$, with ∂ as the subdifferential of Ψ .*

It is noted that the so-called *forward* steps $\mathbf{U}^k - \frac{1}{c_k} \nabla_U f_0(\mathbf{U}^k, \mathbf{V}^k)$ and $\mathbf{V}^k - \frac{1}{d_k} \nabla_V f_0(\mathbf{U}^{k+1}, \mathbf{V}^k)$ in the prox operators in (B.2) are replaced by the simple closed-form solutions $E_r v_r / \|v_r\|^2$ and $E_r^\top u_r / \|u_r\|^2$ in Algorithm 1 at steps 7 and 13 (cf. (6) and (7)), respectively. In the case where the iterates (B.2) are extended to multi-block updates, each block represents one column. This is thanks to the special form of the functions $f_0(\cdot, \mathbf{V}^k)$ and $f_0(\mathbf{U}^{k+1}, \cdot)$ 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 operators *round* and *clamp*_[α, β] in (6) and (7) (and equivalently in Algorithm 1 at steps 8 and 14). The equivalence of these steps is proven in the following lemma.

Lemma 1 (prox implementation). *Consider the operators *round* and *clamp*_[α, β] defined in (6) and (7). Then $\text{prox}_{c_k}^f(\mathbf{W}) = \text{round}(\text{clamp}_{[\alpha, \beta]}(\mathbf{W}))$ and $\text{prox}_{d_k}^g(\mathbf{Z}) = \text{round}(\text{clamp}_{[\alpha, \beta]}(\mathbf{Z}))$ for any $\mathbf{W} \in \mathbb{R}^{M \times R}$, $\mathbf{Z} \in \mathbb{R}^{N \times R}$, and *round*(*clamp*_[α, β](\cdot)) being an elementwise operator on the input matrices.*

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

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

Moreover, note that the *round* operator can be equivalently driven by the following proximal operator:

$$\text{round}(\mathbf{W}) = \arg \min_{\mathbf{U} \in \mathbb{Z}^{M \times R}} \{\|\mathbf{U} - \mathbf{W}\|_F^2\}. \quad (\text{B.3})$$

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

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

The first equality is due to the definition of prox which is equivalent to the second equality. In the third equality the matrices $\mathbf{A} \in \mathbb{R}^{M \times R}$ and $\mathbf{B} \in \mathbb{R}^{M \times R}$ 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 $\text{clamp}_{[\alpha,\beta]}$. Finally, in the last equality, (B.3) is invoked. It is noted that in the implementation, $\text{round}(\text{clamp}_{[\alpha,\beta]}(\cdot)) = \text{clamp}_{[\alpha,\beta]}(\text{round}(\cdot))$ due to the integrality of the bounds $\alpha, \beta \in \mathbb{Z}$. A similar proof can be trivially followed for $\text{prox}_{d_k}^g(\mathbf{Z}) = \text{round}(\text{clamp}_{[\alpha,\beta]}(\mathbf{Z}))$ as well. \square

Now that the equivalence of iterates (B.2) with the simple and closed-form steps in Algorithm 1 is fully established, and the assumptions required for the convergence are verified in proposition 1 to be met by problems (B.1) and (3), proposition 2 can be trivially invoked to establish the convergence of Algorithm 1 to a locally optimal point of problem (3).

References

- [1] G. K. Wallace, The JPEG still picture compression standard, *Communications of the ACM* 34 (1991) 30–44.
- [2] A. Skodras, C. Christopoulos, T. Ebrahimi, The JPEG 2000 still image compression standard, *IEEE Signal processing magazine* 18 (2001) 36–58.
- [3] V. K. Goyal, Theoretical foundations of transform coding, *IEEE Signal Processing Magazine* 18 (2001) 9–21.
- [4] N. Ahmed, T. Natarajan, K. R. Rao, Discrete cosine transform, *IEEE transactions on Computers* 100 (1974) 90–93.
- [5] M. Antonini, M. Barlaud, P. Mathieu, I. Daubechies, Image coding using wavelet transform., *IEEE Trans. Image Processing* 1 (1992) 20–5.
- [6] H. Andrews, C. Patterson, Singular value decomposition (SVD) image coding, *IEEE transactions on Communications* 24 (1976) 425–432.
- [7] H. Prasantha, H. Shashidhara, K. B. Murthy, Image compression using SVD, in: International conference on computational intelligence and multimedia applications (ICCIMA 2007), volume 3, IEEE, 2007, pp. 143–145.
- [8] J. Hou, L.-P. Chau, N. Magnenat-Thalmann, Y. He, Sparse low-rank matrix approximation for data compression, *IEEE Transactions on Circuits and Systems for Video Technology* 27 (2015) 1043–1054.
- [9] J. M. Shapiro, Embedded image coding using zerotrees of wavelet coefficients, *IEEE Transactions on signal processing* 41 (1993) 3445–3462.
- [10] G. Developers, WebP Compression Techniques, <https://developers.google.com/speed/webp>, 2011. Accessed: 2024-05-17.
- [11] J. Lainema, M. M. Hannuksela, V. K. M. Vadakital, E. B. Aksu, HEVC still image coding and high efficiency image file format, in: 2016 IEEE International Conference on Image Processing (ICIP), IEEE, 2016, pp. 71–75.
- [12] M. M. Hannuksela, J. Lainema, V. K. M. Vadakital, The high efficiency image file format standard [standards in a nutshell], *IEEE Signal Processing Magazine* 32 (2015) 150–156.
- [13] J. Ballé, D. Minnen, S. Singh, S. J. Hwang, N. Johnston, Variational image compression with a scale hyperprior, *arXiv preprint arXiv:1802.01436* (2018).
- [14] Z. Cheng, H. Sun, M. Takeuchi, J. Katto, Learned image compression with discretized gaussian mixture likelihoods and attention modules, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2020, pp. 7939–7948.
- [15] J. Liu, H. Sun, J. Katto, Learned image compression with mixed transformer-CNN architectures, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2023, pp. 14388–14397.
- [16] R. Yang, S. Mandt, Lossy image compression with conditional diffusion models, *Advances in Neural Information Processing Systems* 36 (2024).
- [17] X. Yuan, R. Haimi-Cohen, Image compression based on compressive sensing: End-to-end comparison with JPEG, *IEEE Transactions on Multimedia* 22 (2020) 2889–2904.
- [18] M. M. Lin, B. Dong, M. T. Chu, Integer matrix factorization and its application, *Tech. Rep.* (2005).
- [19] B. Dong, M. M. Lin, H. Park, Integer matrix approximation and data mining, *Journal of scientific computing* 75 (2018) 198–224.
- [20] P. Miettinen, T. Mielikäinen, A. Gionis, G. Das, H. Mannila, The discrete basis problem, *IEEE transactions on knowledge and data engineering* 20 (2008) 1348–1362.
- [21] S. Ravanbakhsh, B. Póczos, R. Greiner, Boolean matrix factorization and noisy completion via message passing, in: International Conference on Machine Learning, PMLR, 2016, pp. 945–954.
- [22] C. Eckart, G. Young, The approximation of one matrix by another of lower rank, *Psychometrika* 1 (1936) 211–218.
- [23] D. Lee, H. S. Seung, Algorithms for Non-negative Matrix Factorization, in: T. Leen, T. Dietterich, V. Tresp (Eds.), *Advances in Neural Information Processing Systems*, volume 13, MIT Press, 2000.
- [24] C. H. Ding, T. Li, M. I. Jordan, Convex and semi-nonnegative matrix factorizations, *IEEE transactions on pattern analysis and machine intelligence* 32 (2008) 45–55.
- [25] P. Deutsch, J.-L. Gailly, Zlib compressed data format specification version 3.3, Technical Report, 1996.
- [26] N. Gillis, Nonnegative Matrix Factorization, SIAM, 2020.
- [27] P. van Emde Boas, Another NP-complete problem and the complexity of computing short vectors in a lattice, Technical Report, Department of Mathematics, University of Amsterdam (1981).
- [28] Y. Collet, contributors, Zstandard - fast real-time compression algorithm, 2016. URL: <https://facebook.github.io/zstd/>.
- [29] D. A. Huffman, A method for the construction of minimum-redundancy codes, *Proceedings of the IRE* 40 (1952) 1098–1101.
- [30] A. Clark, Pillow (PIL Fork) Documentation, 2015. URL: <https://buildmedia.readthedocs.org/media/pdf/pillow/latest/pillow.pdf>.
- [31] E. Kodak, Kodak lossless true color image suite, 1993. URL: <https://r0k.us/graphics/kodak/>.
- [32] Organizers, Challenge on Learned Image Compression, 2024. URL: <https://www.compression.cc/>.
- [33] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein,

- et al., Imagenet large scale visual recognition challenge, *International journal of computer vision* 115 (2015) 211–252.
- [34] W. S. Cleveland, S. J. Devlin, Locally weighted regression: an approach to regression analysis by local fitting, *Journal of the American statistical association* 83 (1988) 596–610.
 - [35] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
 - [36] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, L. Fei-Fei, Imagenet: A large-scale hierarchical image database, in: *2009 IEEE conference on computer vision and pattern recognition*, Ieee, 2009, pp. 248–255.
 - [37] Y. Nesterov, Efficiency of coordinate descent methods on huge-scale optimization problems, *SIAM Journal on Optimization* 22 (2012) 341–362.
 - [38] H. Attouch, J. Bolte, B. F. Svaiter, Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forward-backward splitting, and regularized Gauss-Seidel methods, *Mathematical Programming* 137 (2013) 91–129.
 - [39] A. Beck, L. Tetruashvili, On the convergence of block coordinate descent type methods, *SIAM journal on Optimization* 23 (2013) 2037–2060.
 - [40] J. Bolte, S. Sabach, M. Teboulle, Proximal alternating linearized minimization for nonconvex and nonsmooth problems, *Mathematical Programming* 146 (2014) 459–494.