
IMAGE COMPRESSION

CS663 Project Report

Authors

Brian Mackwan (Roll Number : 22b0413)
Ekansh Ravi Shankar (Roll Number : 22b1032)
Ved Danait (Roll Number : 22b1818)
Indian Institute of Technology, Bombay

Contents

1	Introduction	3
2	Datasets	3
3	JPEG Implementation and Analysis	3
3.1	Details of the Scheme	3
3.2	Results	9
4	Comparison with CV2 Implementation	29
5	Experiments with Various Quantization Matrices	34
6	PCA Implementation and Comparison with JPEG	45
6.1	Details of Implementation	45
6.2	Comparison for Original Dataset	47
7	Sparse Orthonormal Transforms and Comparison with Custom JPEG	50
7.1	Details of Implementation	50
7.2	Comparison with Custom JPEG	52
8	Conclusions and Wrapping Up	57
9	Contributions	57

1 Introduction

We implement a lossy Image Compression scheme using the JPEG algorithm and conduct some analysis of the quality of the implementation relative to CV2's implementation on metrics such as "RMSE vs BPP" and "Compression Rate vs Quality Factor".

We also perform an experiment for 5 different Quantization Matrices : JPEG Standard, Highly Detailed, Low Detail, JPEG2000-like and Uniform and conduct an analysis based on the same metrics. We also see how the compressed image behaves for different kinds of images.

We also implement a compression algorithm using PCA and compare it with our custom implementation of JPEG for our dataset of 19 grayscale images. We see that JPEG outperforms PCA in most of these cases.

In addition, we also thoroughly studied, and implemented the paper, Osman Gokhan Sezer, Onur G. Guleryuz and Yucel Altunbasak, "Approximation and Compression With Sparse Orthonormal Transforms", IEEE Transactions on Image Processing, 2015, which gives us an algorithm to obtain transformations for compression by solving optimization problems (called "SOT" from now onwards). We will see how our custom SOT compares with our custom JPEG implementation.

2 Datasets

For all the JPEG analysis we have conducted in this report, we use the 19 grayscale bmp images from the website Hlevkin Test Images which are of appropriate size (dimensions multiple of 8). We could use 20, but there was some issue with Zelda2 so we left it out.

3 JPEG Implementation and Analysis

3.1 Details of the Scheme

This implementation of the JPEG compression algorithm follows the standard steps:

1. Block-based Discrete Cosine Transform (DCT).
2. Quantization using a scaled matrix based on the quality factor.
3. Encoding using Zig-Zag traversal, Run-Length Encoding (RLE), and Huffman coding.
4. Decoding involves reversing the above steps to restore the image.

Quantization Matrix and Scaling

The class JPEG initializes a default quantization matrix, scaled by the input quality factor Q. The default quantization matrix is :

$$M = \begin{bmatrix} 16 & 11 & 10 & 16 & 24 & 40 & 51 & 61 \\ 12 & 12 & 14 & 19 & 26 & 58 & 60 & 55 \\ 14 & 13 & 16 & 24 & 40 & 57 & 69 & 56 \\ 14 & 17 & 22 & 29 & 51 & 87 & 80 & 62 \\ 18 & 22 & 37 & 56 & 68 & 109 & 103 & 77 \\ 24 & 35 & 55 & 64 & 81 & 104 & 113 & 92 \\ 49 & 64 & 78 & 87 & 103 & 121 & 120 & 101 \\ 72 & 92 & 95 & 98 & 112 & 100 & 103 & 99 \end{bmatrix}$$

We scale this using the quality factor Q as :

```
1 self.scaled_quantization_matrix = np.round(self.quantization_matrix * (50.0 / Q))
```

Discrete Cosine Transform (DCT) and Inverse DCT

Two helper functions handle 2D DCT and inverse DCT on 8x8 image blocks:

```
1 def _dct2(self, block):
2     return dct(dct(block.T, norm='ortho').T, norm='ortho')
3
4 def _idct2(self, block):
5     return idct(idct(block.T, norm='ortho').T, norm='ortho')
```

Image Chunking and Flattening

Images are divided into 8x8 blocks for processing. After decoding, the blocks are recombined into the original shape. We do this using the functions `_chunkify(self, image)` and `_dechunkify(self, chunks, original_shape)`. Another pair of helper functions we use is `_zigzag(self, matrix)` and `_reverse_zigzag(self, flattened, rows, cols)`. This is used to flatten the quantized patch into a 1 dimensional array in order to apply RLE and Huffman Encoding. The reverse function is used while decoding when we decode in order to restore the shape of the obtained matrix after retrieving the flattened matrix while decoding the compressed file.

Run-Length Encoding (RLE) and Huffman Coding

The RLE algorithm compresses the image's non-zero coefficients, and Huffman coding is used to store them efficiently.

```
1 def _run_length_encoding(self, matrix):
2     rle = []
3     zero_count = 0
4     for value in matrix:
5         if value == 0:
6             zero_count += 1
7             if zero_count == 16:
8                 rle.append((15, 0))
9                 zero_count = 0
10            else:
11                rle.append((zero_count, value))
12                zero_count = 0
13    return rle
```

We have used a standard implementation for Huffman encoding which takes the values given after the RLE compression and assigns bitstrings to efficiently encode the values provided with respect to their frequencies.

Encoding Process

1. Quantize the DCT-transformed 8x8 blocks.
2. Apply Zig-Zag traversal.
3. Compress using RLE and encode with Huffman coding.
4. Store results in a binary file.

```
1 def compress(self, image, filename):
2     image = image.astype('float32') # cast to float
3     image = image - 128.0 # rescale the values
4     chunks = self._chunkify(image) # break into 8x8 chunks
5
6     flattened_matrices = []
7     for chunk in chunks:
8         quantized_matrix = self._quantize(self._dct2(chunk))
9         flattened_matrix = self._zigzag(quantized_matrix)
10        flattened_matrices.append(flattened_matrix)
11
12    self._encode(flattened_matrices, filename)
```

```

1   def _encode(self, flattened_matrices, filename):
2       with open(filename, 'wb') as f:
3           rle_result = []
4
5           for i in range(0, len(flattened_matrices)):
6               rle = self._run_length_encoding(flattened_matrices[i])
7               rle_result.extend(rle)
8               rle_result.append((0,0))
9
10          values = [value for count, value in rle_result] # Only values for Huffman
11              encoding
12          huffman_codes = self._huffman_encoding(values)
13
14          # Store the Quality Factor in the file
15          f.write(np.array([self.Q], dtype=np.uint8).tobytes())
16          f.flush()
17
18          # Store the Huffman table in the file (value -> huffman code mapping)
19          f.write(np.array([len(huffman_codes)], dtype=np.uint16).tobytes()) # Write the
19              number of unique values in Huffman table
20          f.flush()
21
22
23          # Collect all Huffman codes
24          all_huffman_codes = []
25
26          for value, code in huffman_codes.items():
27              f.write(np.array([value], dtype=np.int16).tobytes())
28              f.flush()
29              size_in_bits = len(code)
30              f.write(np.array([size_in_bits], dtype=np.uint8).tobytes()) # Write the
30                  size of the Huffman code
31              f.flush()
32              all_huffman_codes.append(code)
33
34          all_huffman_codes_str = ''.join(all_huffman_codes)
35          total_size_in_bits = len(all_huffman_codes_str)
36
37          f.write(np.array([total_size_in_bits], dtype=np.uint32).tobytes())
38          f.flush()
39
40          padding_length = 0
41          if len(all_huffman_codes_str) % 8 != 0:
42              padding_length = 8 - (len(all_huffman_codes_str) % 8)
43          all_huffman_codes_str = all_huffman_codes_str + '0' * padding_length
44
45          bitstream = bitstring.BitStream(bin=all_huffman_codes_str)
46          f.write(bitstream.bytes) # Write the Huffman code as bytes
47          f.flush()
48
49          # To store the RLE Result
50          # Store the size of the rle_result
51          f.write(np.array([len(rle_result)], dtype=np.uint32).tobytes())
52          f.flush()
53
54          all_huffman_codes = []
55          all_huffman_codes_str = ''
56
57          for count, value in rle_result:
58              huffman_code = huffman_codes[value] # Find Huffman code for this value
59              size_in_bits = len(huffman_code) # The size in bits of the Huffman code
60
61              f.write(np.array([count], dtype=np.uint8).tobytes())
62              f.write(np.array([size_in_bits], dtype=np.uint8).tobytes())
63
64              f.flush()

```

```

65         all_huffman_codes.append(huffman_code)
66
67
68     all_huffman_codes_str = ''.join(all_huffman_codes)
69     total_size_in_bits = len(all_huffman_codes_str)
70
71     # Calculate padding to ensure byte alignment
72     if total_size_in_bits % 8 != 0:
73         padding_length = 8 - (total_size_in_bits % 8)
74         all_huffman_codes_str += '0' * padding_length # Add padding to the
75             bitstream
76     else:
77         padding_length = 0
78
79     # Write the total size in bits to the file (before padding)
80     f.write(np.array([total_size_in_bits], dtype=np.uint32).tobytes())
81     f.flush()
82
83     # Convert the padded bitstream into a byte array
84     bitstream = bitstring.BitStream(bin=all_huffman_codes_str)
85     f.write(bitstream.bytes) # Write the Huffman codes as bytes
86     f.flush()

```

Decoding Process

1. Decode Huffman codes and RLE data from the file.
2. Reconstruct blocks using inverse Zig-Zag and dequantization.
3. Apply inverse DCT and combine blocks to restore the image.

```

1 def decompress(self, filename):
2     with open(filename, 'rb') as f:
3         # Read the Quality Factor (Q)
4         Q = np.fromfile(f, dtype=np.uint8, count=1)[0]
5
6         # Read the Huffman table (number of entries)
7         num_huffman_codes = np.fromfile(f, dtype=np.uint16, count=1)[0]
8
9         huffman_codes = {}
10
11        table_entries = []
12        # Read each Huffman code (value, size)
13        for _ in range(num_huffman_codes):
14            value = np.fromfile(f, dtype=np.int16, count=1)[0]
15            size_in_bits = np.fromfile(f, dtype=np.uint8, count=1)[0]
16            table_entries.append((value, size_in_bits))
17
18        length_huffman_codes = np.fromfile(f, dtype=np.uint32, count=1)[0]
19        huffman_codes_bitstream = f.read((length_huffman_codes + 7) // 8) # Read the
20            total size in bytes
21        huffman_codes_bitstream = bitstring.BitStream(bytes=huffman_codes_bitstream).bin
22
23        current_pos = 0
24        for value, size_in_bits in table_entries:
25            huffman_codes[value] = huffman_codes_bitstream[current_pos:current_pos+
26                size_in_bits]
27            current_pos+=size_in_bits
28
29        # get the number of rle values
30        num_rle_values = np.fromfile(f, dtype=np.uint32, count=1)[0]
31
32        # 3. Decode the RLE-encoded values and sizes
33        rle_result = []

```

```

33     table_entries = []
34
35     for _ in range(num_rle_values):
36         count = np.fromfile(f, dtype=np.uint8, count=1)[0]
37         size_in_bits = np.fromfile(f, dtype=np.uint8, count=1)[0]
38         table_entries.append((count, size_in_bits))
39
40     # print(table_entries)
41
42     length_huffman_codes = np.fromfile(f, dtype=np.uint32, count=1)[0]
43     huffman_codes_bitstream = f.read((length_huffman_codes + 7) // 8) # Read the
44     # total size in bytes
45     huffman_codes_bitstream = bitstring.BitStream(bytes=huffman_codes_bitstream).bin
46
47     current_pos = 0
48     for count, size_in_bits in table_entries:
49         value = self._decode_huffman_code(huffman_codes_bitstream[current_pos:
50                                         current_pos+size_in_bits], huffman_codes)
51         rle_result.append((count, value))
52         current_pos+=size_in_bits
53
54
55     # Reconstruct the original flattened matrices from RLE
56     flattened_matrices = []
57
58     current_patch = []
59     for count, value in rle_result:
60         if count == 0 and value == 0:
61             current_patch = current_patch + [0]*(64 - len(current_patch))
62             flattened_matrices.append(np.array(current_patch))
63             current_patch = []
64         else:
65             current_patch = current_patch + [0]*count + [value]
66
67
68     # convert it back to the original patch dimensions
69     for i in range(len(flattened_matrices)):
70
71         flattened_matrices[i] = self._reverse_zigzag(flattened_matrices[i], 8, 8)
72
73         # now de-quantize the matrix and rescale
74         flattened_matrices[i] = (flattened_matrices[i] * self.
75             scaled_quantization_matrix)
76         flattened_matrices[i] = self._idct2(flattened_matrices[i]) + 128
77
78     # combine the patches together
79     restored_image = self._dechunkify(flattened_matrices, (8*int(np.sqrt(len(
80         flattened_matrices))), 8*int(np.sqrt(len(flattened_matrices))))))
81
82     return restored_image

```

Format of the Compressed Image

The structure of a JPEG image file is represented as follows:

- **Quality Factor (Q):** A single byte representing the image compression quality.
- **Length of Huffman Table:** Specifies the number of Huffman codes. Each entry contains:
 - *Value* (`int8`): The symbol associated with the Huffman code.
 - *Size of Code* (`uint8`): Length of the Huffman code in bits.
- **Total Length of Concatenated Huffman Codes:** Specifies the combined size of all Huffman codes in bits (`uint32`).
- **Huffman Codes:** The concatenated Huffman bitstream.

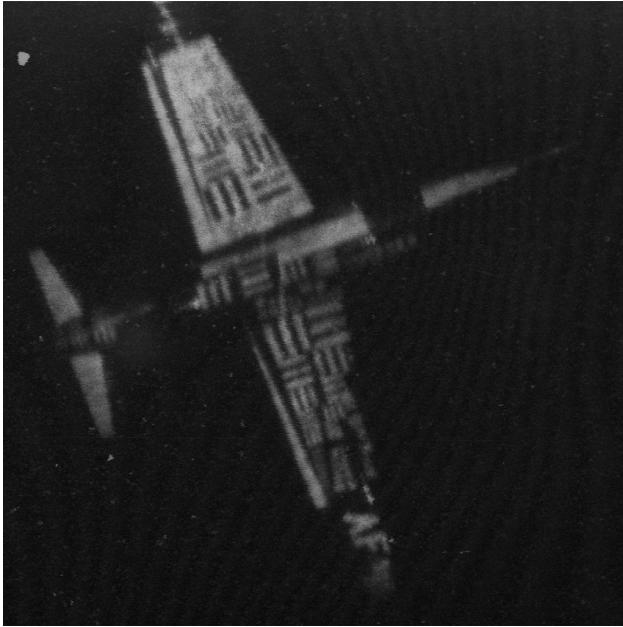
Field	Data Type
Quality Factor (Q)	uint8
Length of Huffman Table	uint16
- Value	int8
- Size of Code	uint8
Total Length of concatenated Huffman Codes	uint32
Huffman Codes	(Variable Length)
Size of Run-Length Encoded (RLE) Result	uint32
- Count	uint8
- Size of Code	uint8
Total Length of concatenated Huffman Codes	uint32
Huffman Codes	(Variable Length)

Structure of JPEG Image File

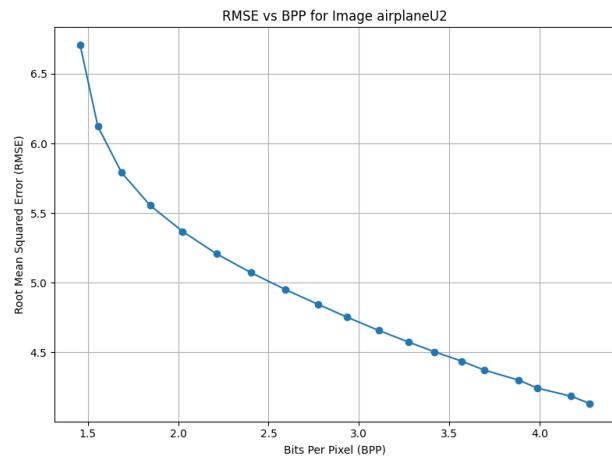
- **Size of Run-Length Encoded (RLE) Result:** Total number of RLE (count, value) pairs (uint32).
 - *Count* (uint8): Number of consecutive zeros.
 - *Size of Code* (uint8): Length of the corresponding Huffman code in bits.

3.2 Results

AirplaneU2

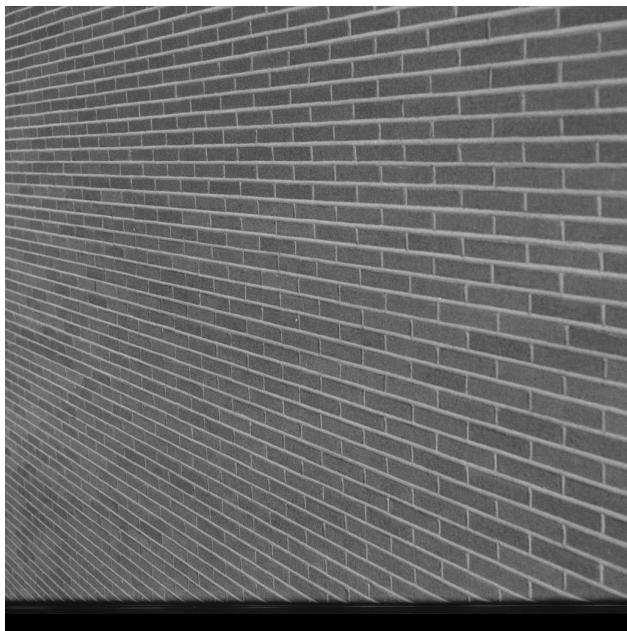


Original Image

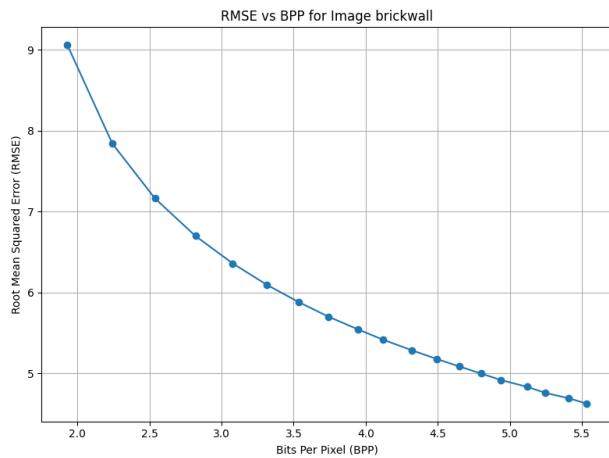


RMSE vs. BPP for varying quality factors (Q).

Brickwall



Original Image

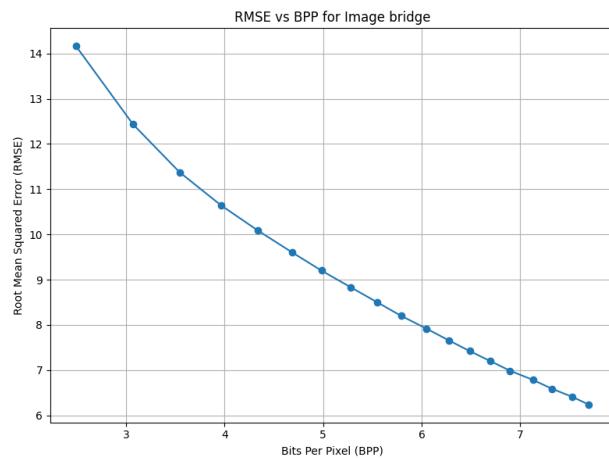


RMSE vs. BPP for varying quality factors (Q).

Bridge



Original Image

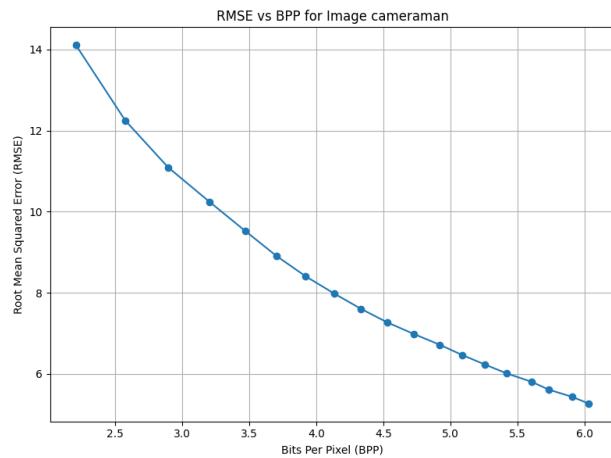


RMSE vs. BPP for varying quality factors (Q).

Cameraman

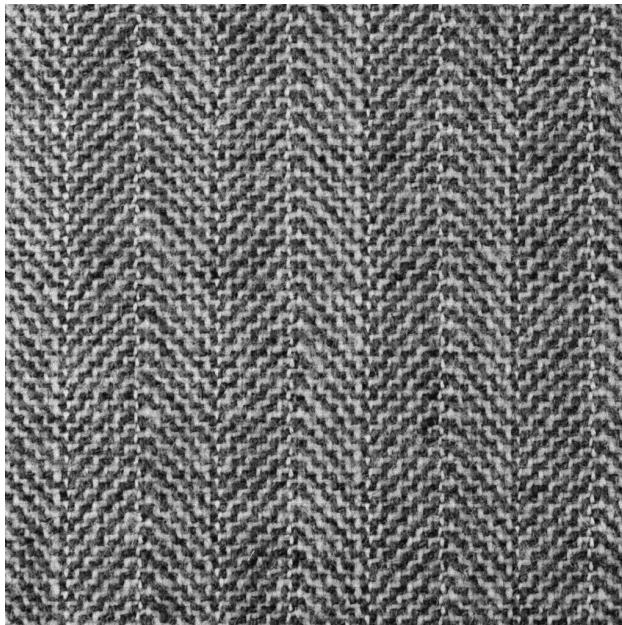


Original Image

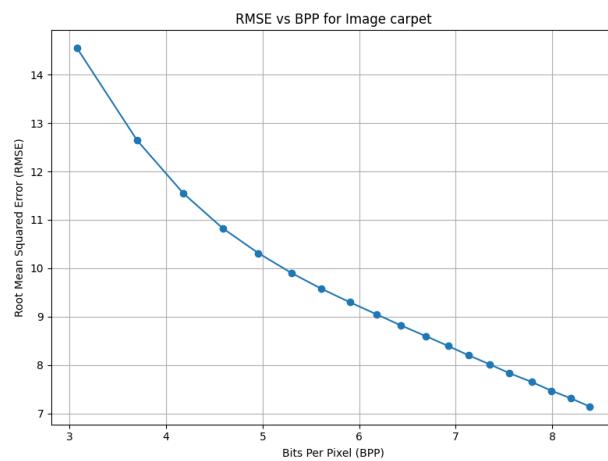


RMSE vs. BPP for varying quality factors (Q).

Carpet



Original Image

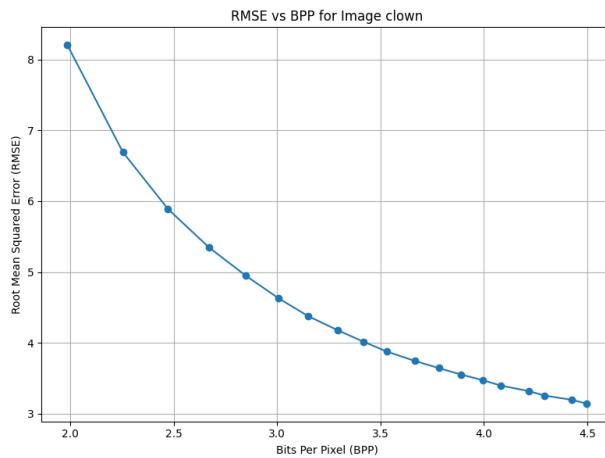


RMSE vs. BPP for varying quality factors (Q).

Clown



Original Image

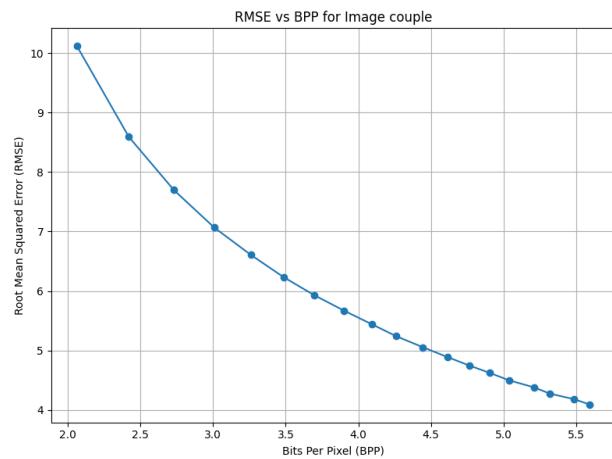


RMSE vs. BPP for varying quality factors (Q).

Couple



Original Image

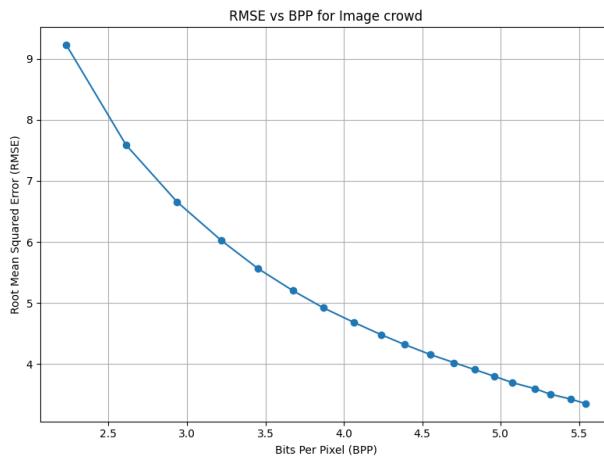


RMSE vs. BPP for varying quality factors (Q).

Crowd



Original Image

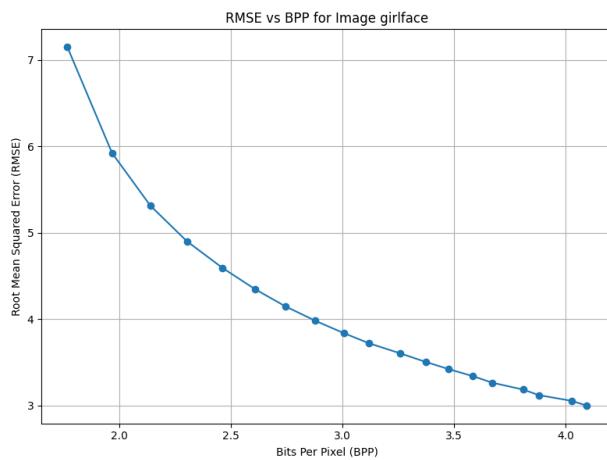


RMSE vs. BPP for varying quality factors (Q).

Girlface



Original Image

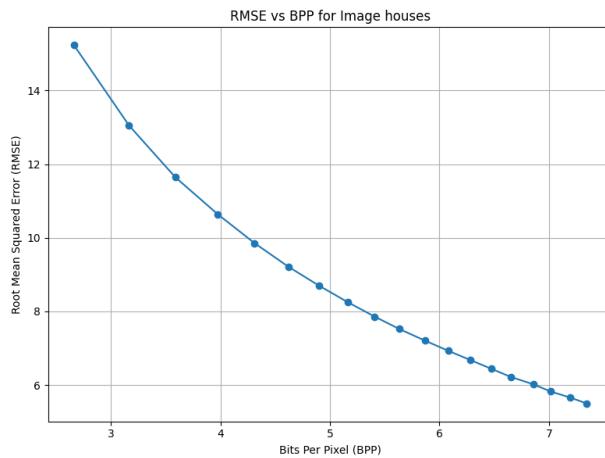


RMSE vs. BPP for varying quality factors (Q).

Houses



Original Image

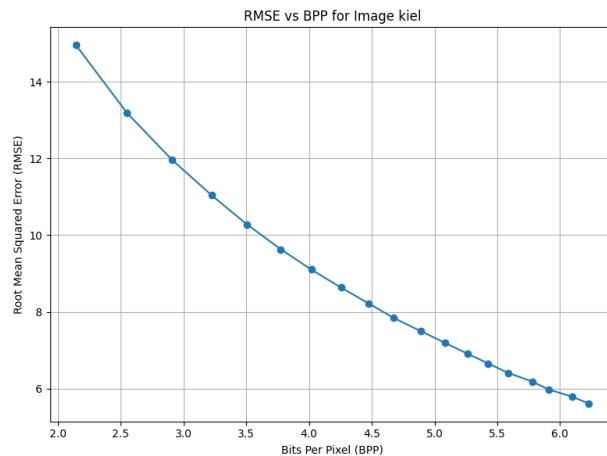


RMSE vs. BPP for varying quality factors (Q).

Kiel



Original Image

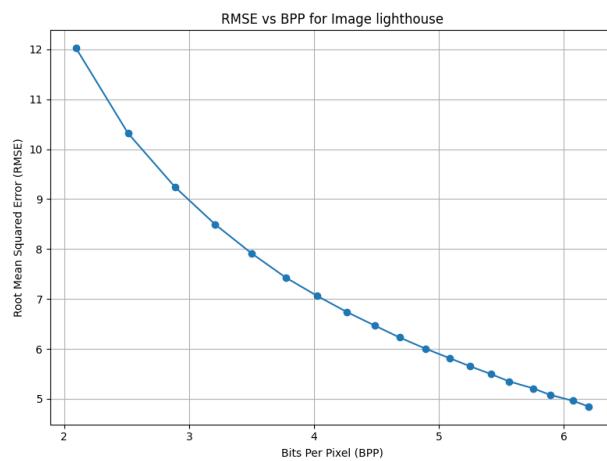


RMSE vs. BPP for varying quality factors (Q).

Lighthouse

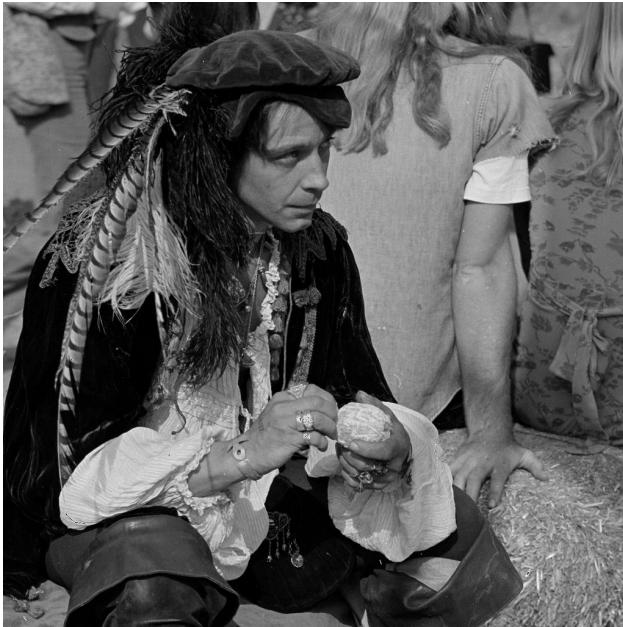


Original Image

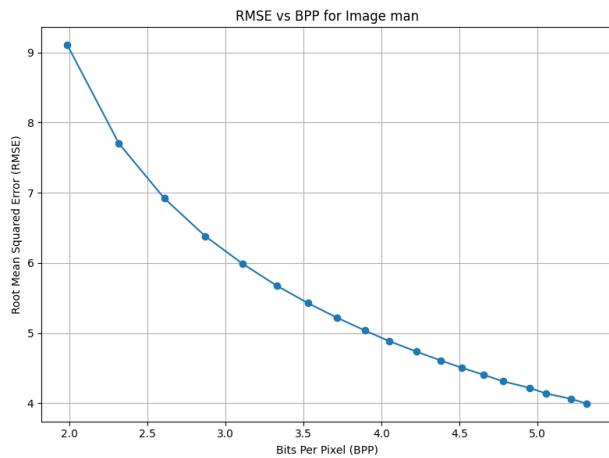


RMSE vs. BPP for varying quality factors (Q).

Man



Original Image

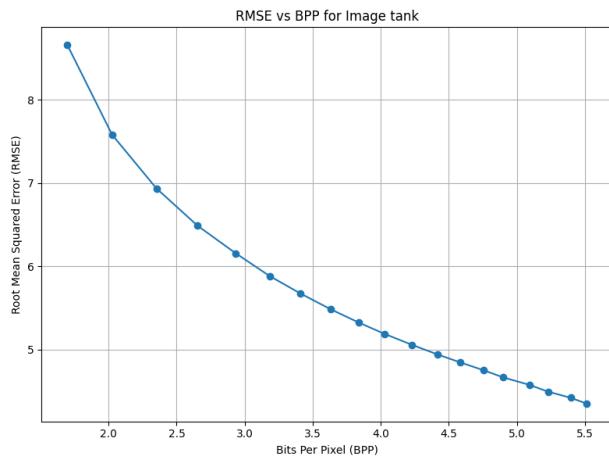


RMSE vs. BPP for varying quality factors (Q).

Tank



Original Image

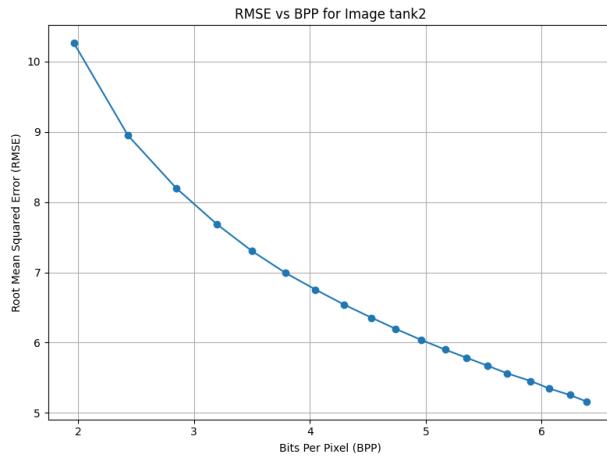


RMSE vs. BPP for varying quality factors (Q).

Tank2

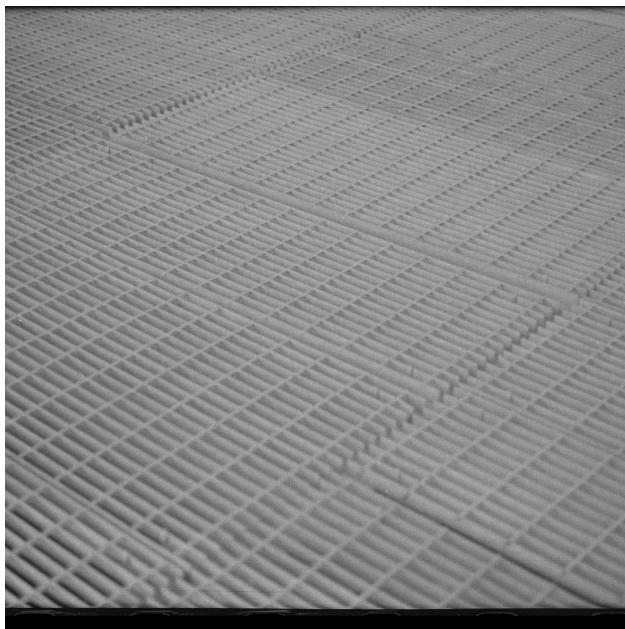


Original Image

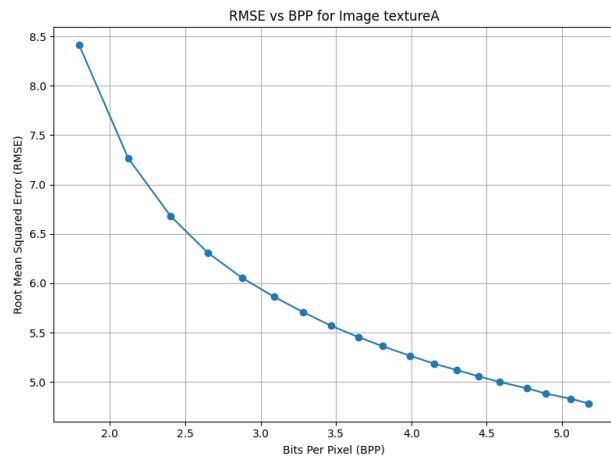


RMSE vs. BPP for varying quality factors (Q).

TextureA



Original Image

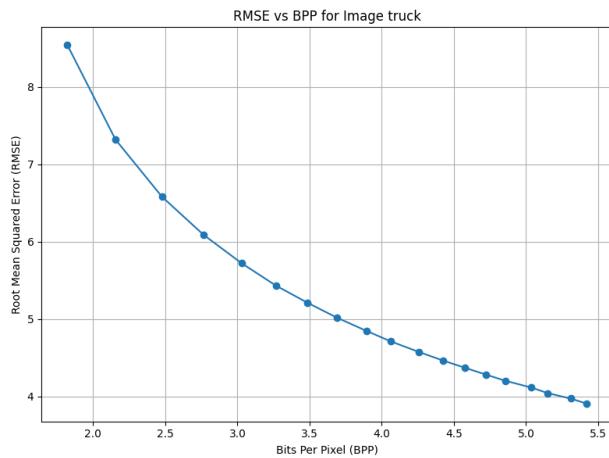


RMSE vs. BPP for varying quality factors (Q).

Truck



Original Image

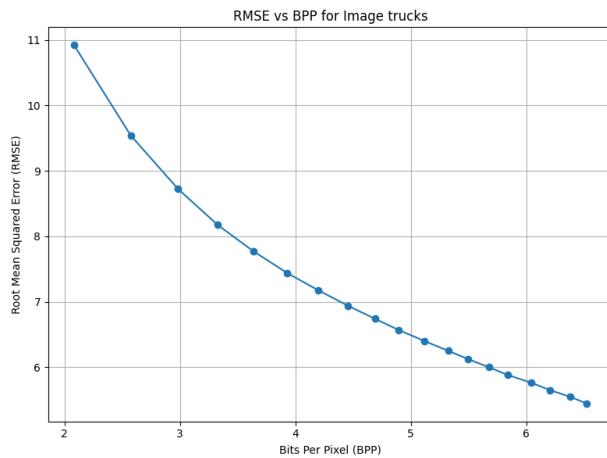


RMSE vs. BPP for varying quality factors (Q).

Trucks



Original Image

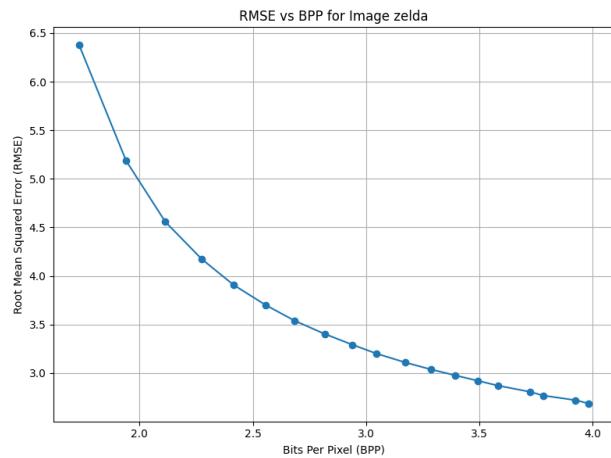


RMSE vs. BPP for varying quality factors (Q).

Zelda

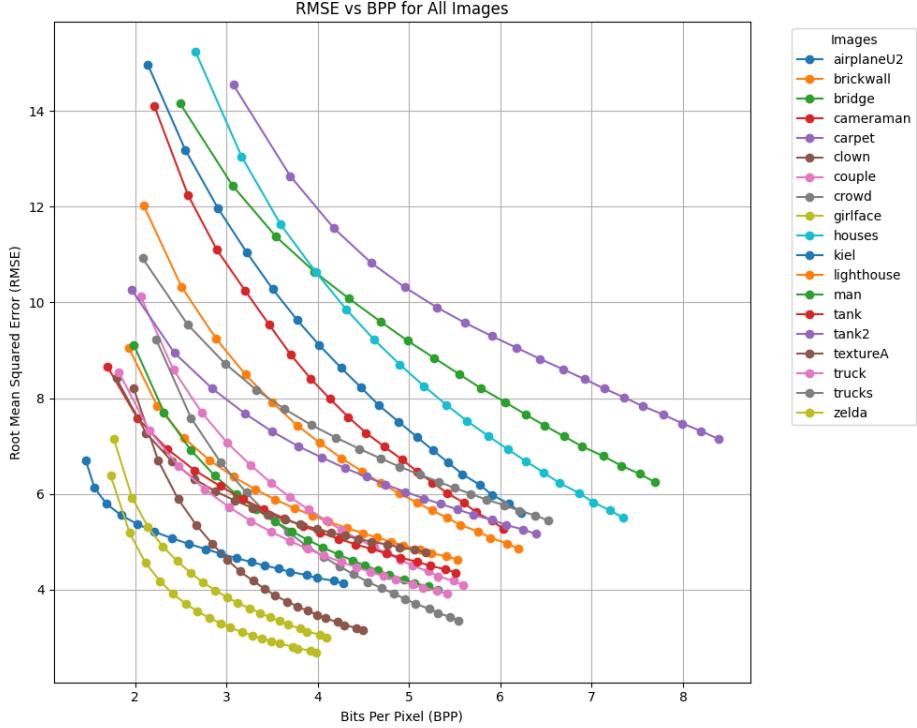


Original Image



RMSE vs. BPP for varying quality factors (Q).

All Images



RMSE vs. BPP for varying quality factors (Q).

The graph illustrates the relationship between **Root Mean Squared Error (RMSE)** and **Bits Per Pixel (BPP)** for various images. A general trend observed across all images is that the RMSE decreases as BPP increases. This aligns with expectations, as higher BPP allows for a higher fidelity representation, reducing the reconstruction error.

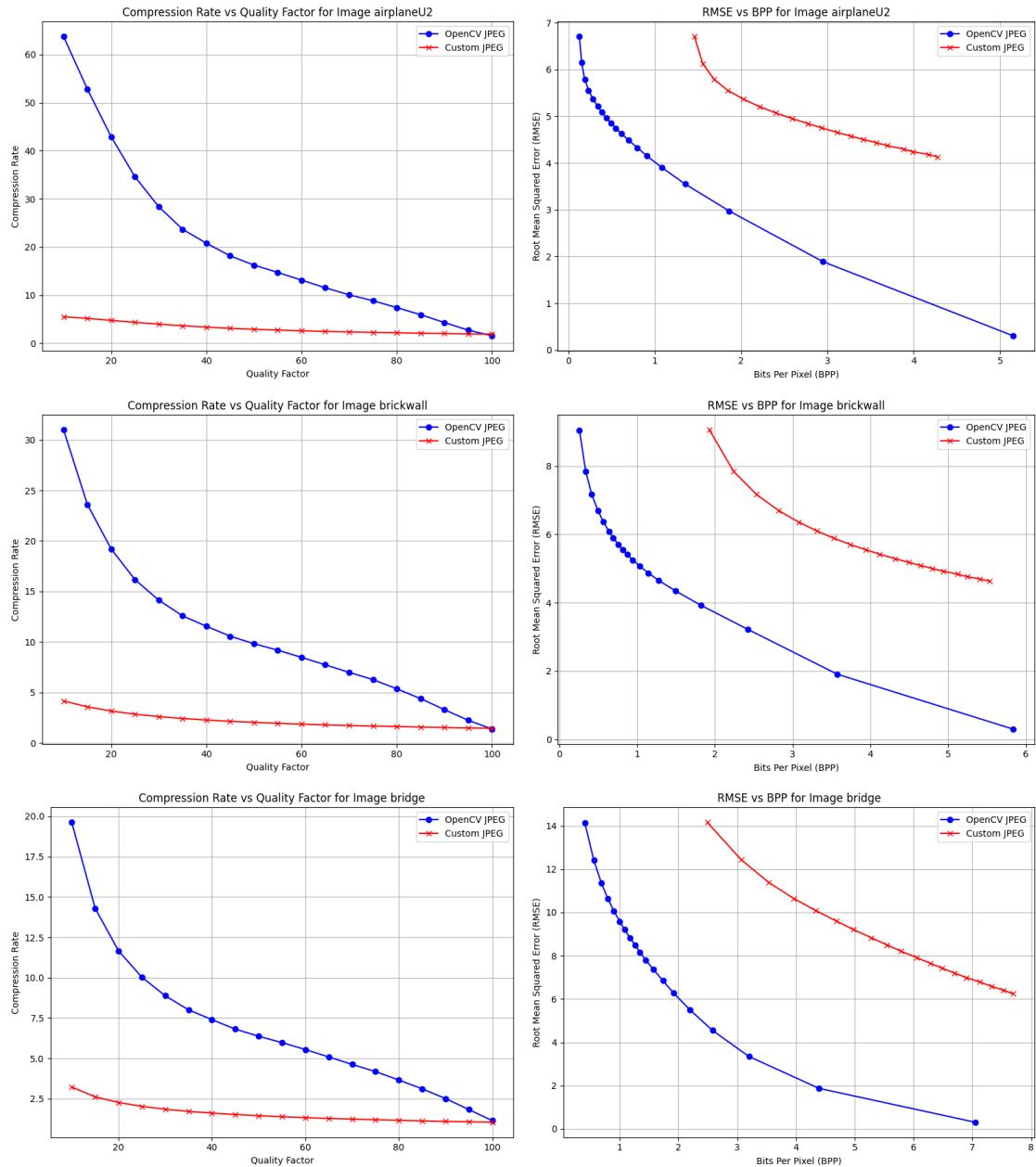
At low BPP values (e.g., $BPP < 3$), the RMSE decreases sharply for most images, indicating that even a modest increase in BPP leads to significant quality improvements. For example, images like *airplaneU2* and *zelda* show a steep decline in RMSE within this range, reflecting their efficient compression and reconstruction characteristics. In contrast, images with more texture or detail, such as *brickwall* and *textureA*, exhibit higher RMSE values at similar BPP, suggesting that they are more challenging to compress effectively.

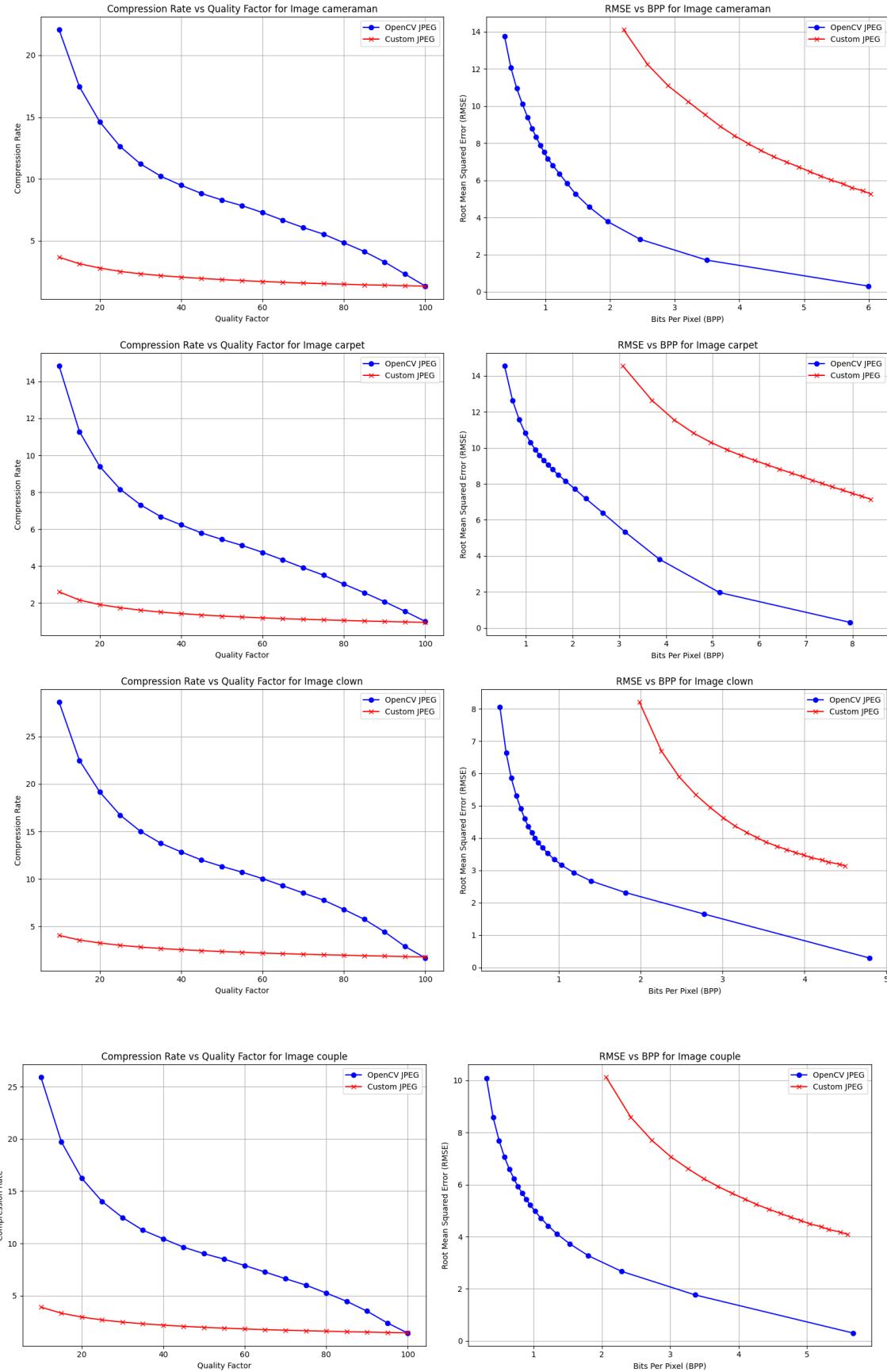
Beyond a certain threshold (e.g., $BPP > 6$), the curves exhibit asymptotic behavior, with RMSE improvements becoming negligible. This phenomenon demonstrates diminishing returns in quality with increased BPP, as the representation approaches near-lossless levels. For example, images like *brickwall* and *crowd* display relatively flat RMSE curves beyond $BPP = 6$, indicating that further increases in BPP would result in minimal perceptual improvement.

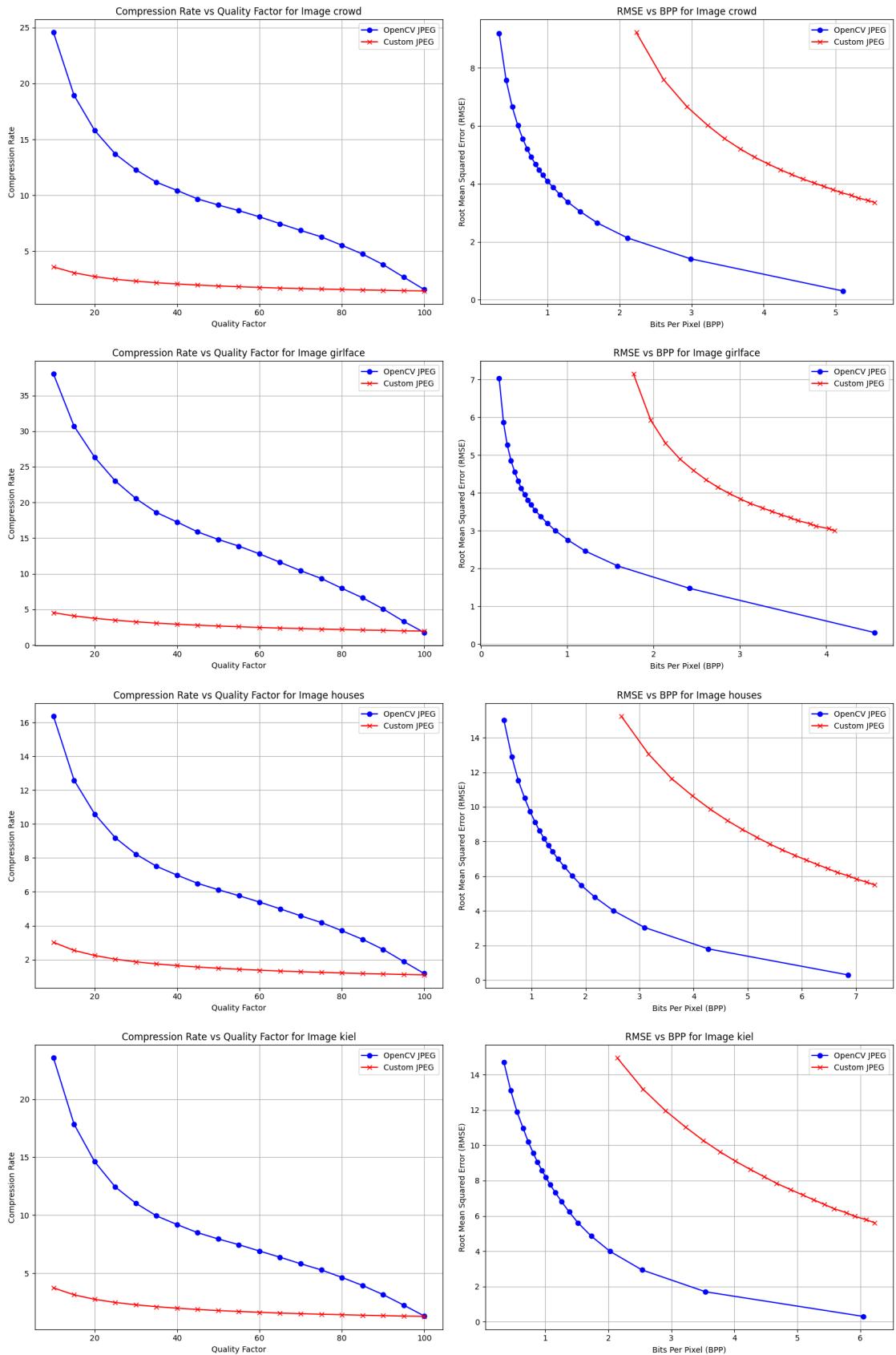
The variability between the curves highlights the impact of image characteristics on compression performance. Simpler images, such as *airplaneU2* and *zelda*, achieve lower RMSE values at the same BPP compared to more complex images like *brickwall* and *textureA*. Additionally, some images exhibit overlapping curves in the low BPP range (e.g., *clown*, *truck*, and *couple*), suggesting that certain features are handled similarly by the compression algorithm under stringent bit-budget constraints.

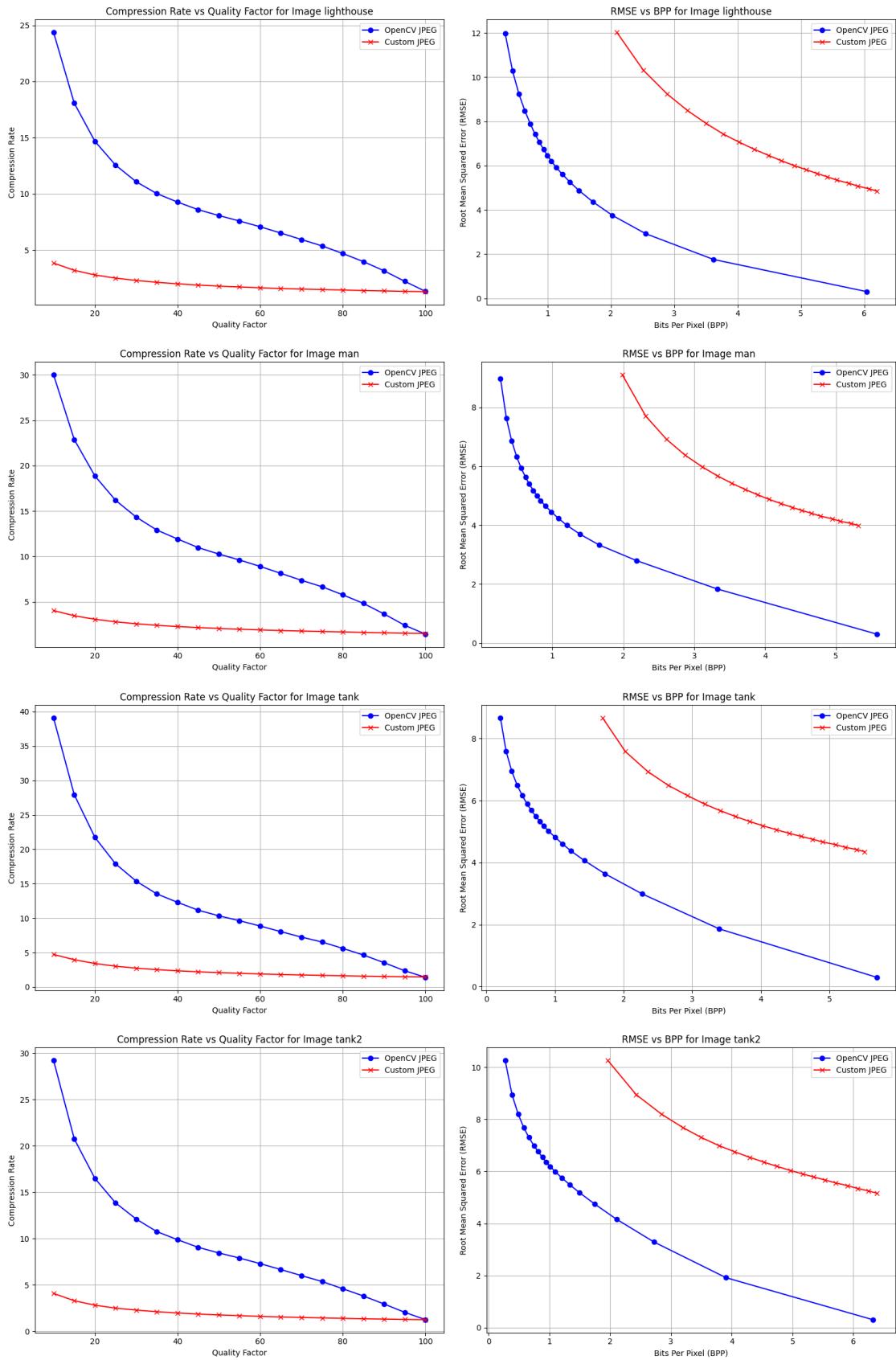
In summary, the analysis highlights the tradeoff between compression efficiency and reconstruction quality. For simpler images, lower BPP values suffice to maintain high quality, whereas textured images require higher BPP to achieve comparable RMSE. The diminishing returns beyond certain BPP values emphasize the importance of selecting an optimal bit allocation for practical applications.

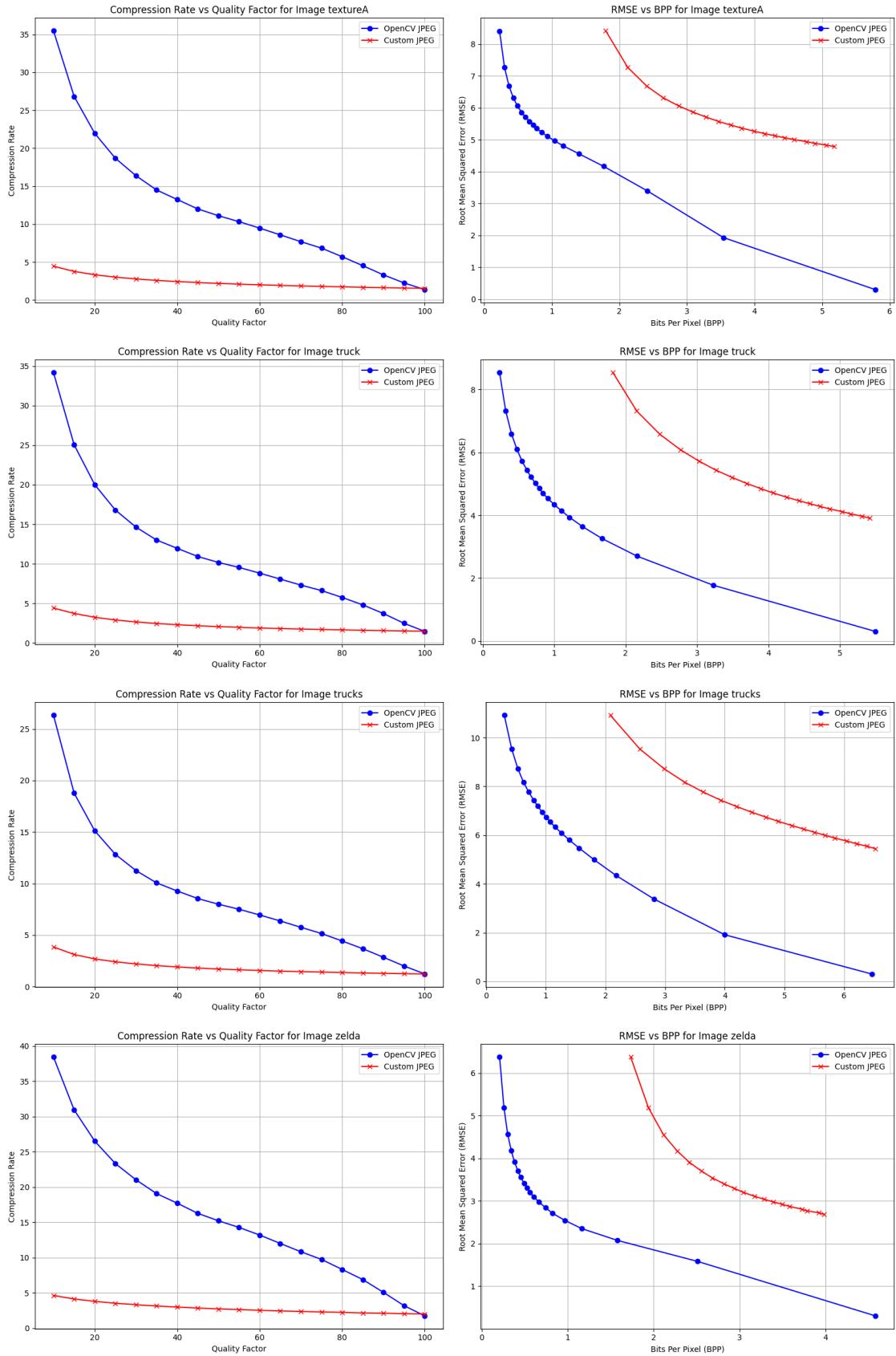
4 Comparison with CV2 Implementation











5 Experiments with Various Quantization Matrices

We run our custom JPEG algorithm using 3 different Quantization matrices and plot the "Compression Rate vs Quality Factor" and "RMSE vs BPP" curves for the same. We also display the image at Q=50 and try and observe the variations. The matrices are as follows:

$$\text{JPEG STANDARD} = \begin{bmatrix} 16 & 11 & 10 & 16 & 24 & 40 & 51 & 61 \\ 12 & 12 & 14 & 19 & 26 & 58 & 60 & 55 \\ 14 & 13 & 16 & 24 & 40 & 57 & 69 & 56 \\ 14 & 17 & 22 & 29 & 51 & 87 & 80 & 62 \\ 18 & 22 & 37 & 56 & 68 & 109 & 103 & 77 \\ 24 & 35 & 55 & 64 & 81 & 104 & 113 & 92 \\ 49 & 64 & 78 & 87 & 103 & 121 & 120 & 101 \\ 72 & 92 & 95 & 98 & 112 & 100 & 103 & 99 \end{bmatrix}$$

$$\text{HIGHLY DETAILED} = \begin{bmatrix} 4 & 3 & 3 & 4 & 6 & 10 & 13 & 15 \\ 3 & 3 & 4 & 5 & 7 & 16 & 16 & 14 \\ 4 & 4 & 5 & 7 & 10 & 15 & 18 & 15 \\ 4 & 5 & 6 & 8 & 14 & 24 & 23 & 19 \\ 6 & 7 & 11 & 17 & 20 & 32 & 30 & 23 \\ 10 & 16 & 25 & 29 & 37 & 48 & 53 & 41 \\ 13 & 16 & 20 & 23 & 28 & 56 & 55 & 46 \\ 15 & 14 & 15 & 19 & 23 & 41 & 43 & 41 \end{bmatrix}$$

$$\text{LOW DETAIL} = \begin{bmatrix} 50 & 40 & 35 & 50 & 60 & 80 & 90 & 100 \\ 40 & 40 & 50 & 60 & 70 & 100 & 110 & 90 \\ 45 & 50 & 60 & 75 & 85 & 100 & 120 & 100 \\ 50 & 60 & 70 & 90 & 110 & 150 & 140 & 120 \\ 60 & 70 & 90 & 110 & 130 & 180 & 170 & 140 \\ 80 & 100 & 130 & 140 & 160 & 200 & 210 & 160 \\ 90 & 110 & 130 & 150 & 170 & 220 & 240 & 190 \\ 100 & 120 & 140 & 160 & 180 & 220 & 240 & 210 \end{bmatrix}$$

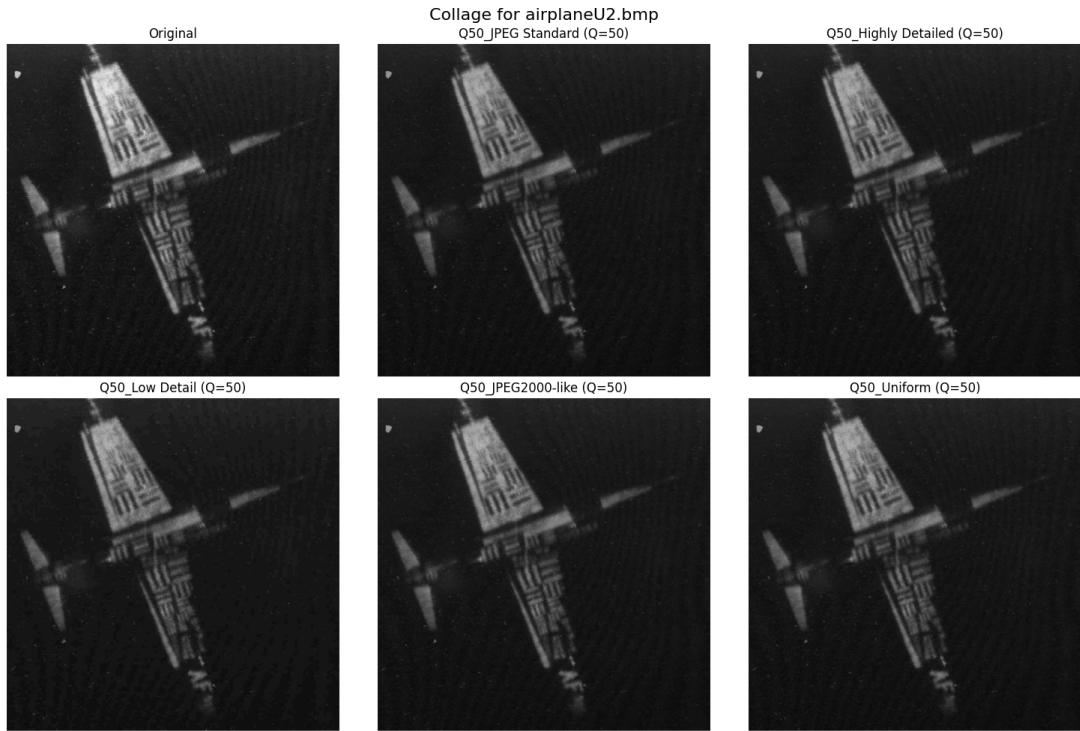
$$\text{JPEG2000 LIKE} = \begin{bmatrix} 9 & 5 & 4 & 6 & 8 & 12 & 16 & 18 \\ 5 & 5 & 6 & 8 & 10 & 18 & 19 & 16 \\ 6 & 6 & 8 & 10 & 14 & 21 & 24 & 20 \\ 6 & 8 & 10 & 13 & 17 & 27 & 26 & 21 \\ 8 & 10 & 14 & 17 & 21 & 31 & 29 & 23 \\ 12 & 18 & 21 & 24 & 30 & 39 & 42 & 32 \\ 16 & 19 & 24 & 26 & 33 & 43 & 44 & 36 \\ 18 & 16 & 20 & 21 & 26 & 36 & 39 & 35 \end{bmatrix}$$

$$\text{UNIFORM} = \begin{bmatrix} 20 & 20 & 20 & 20 & 20 & 20 & 20 & 20 \\ 20 & 20 & 20 & 20 & 20 & 20 & 20 & 20 \\ 20 & 20 & 20 & 20 & 20 & 20 & 20 & 20 \\ 20 & 20 & 20 & 20 & 20 & 20 & 20 & 20 \\ 20 & 20 & 20 & 20 & 20 & 20 & 20 & 20 \\ 20 & 20 & 20 & 20 & 20 & 20 & 20 & 20 \\ 20 & 20 & 20 & 20 & 20 & 20 & 20 & 20 \\ 20 & 20 & 20 & 20 & 20 & 20 & 20 & 20 \end{bmatrix}$$

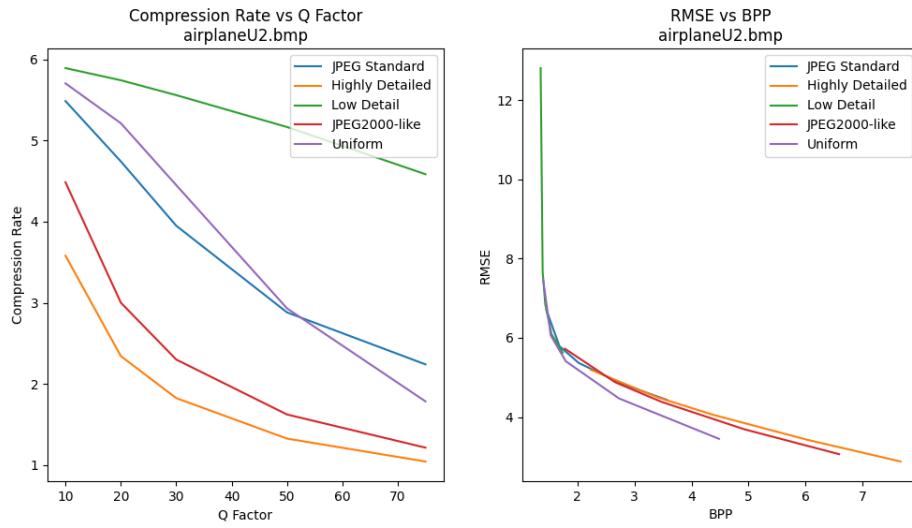
We present the results of image compression experiments for a subset of our images. Each page contains:

1. A collage of six images: the original image and Q=50 compressed images for various quantization matrices.
2. Two graphs comparing compression metrics mentioned above

AirplaneU2

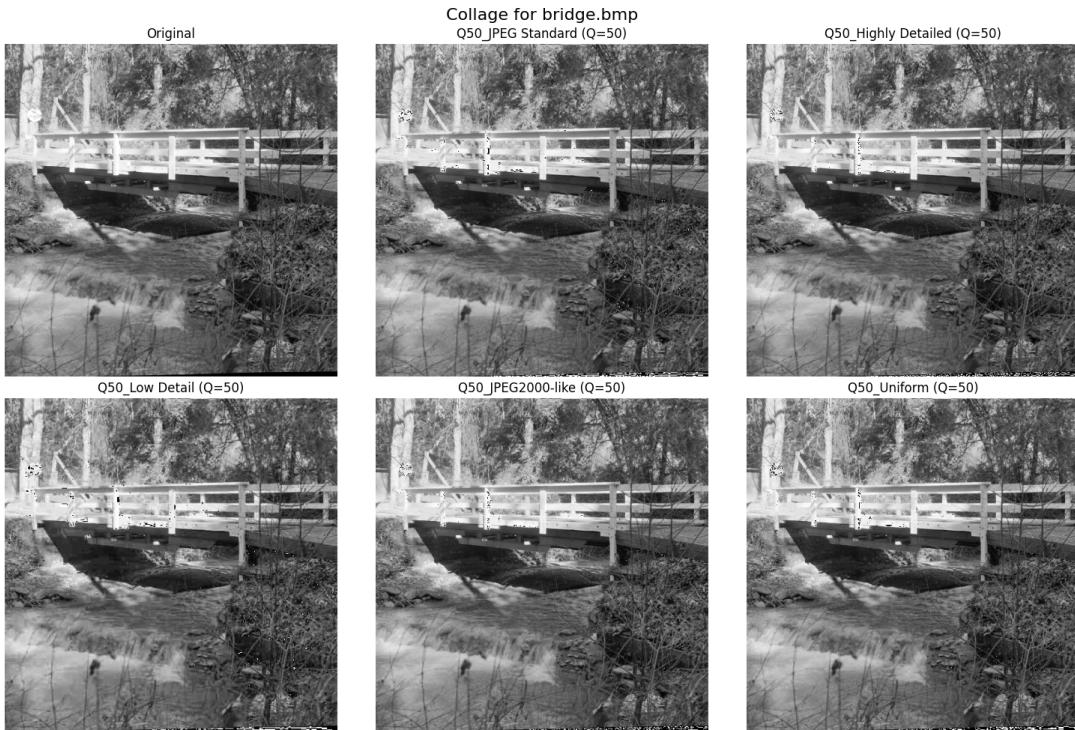


Original image and Q=50 compressed images.

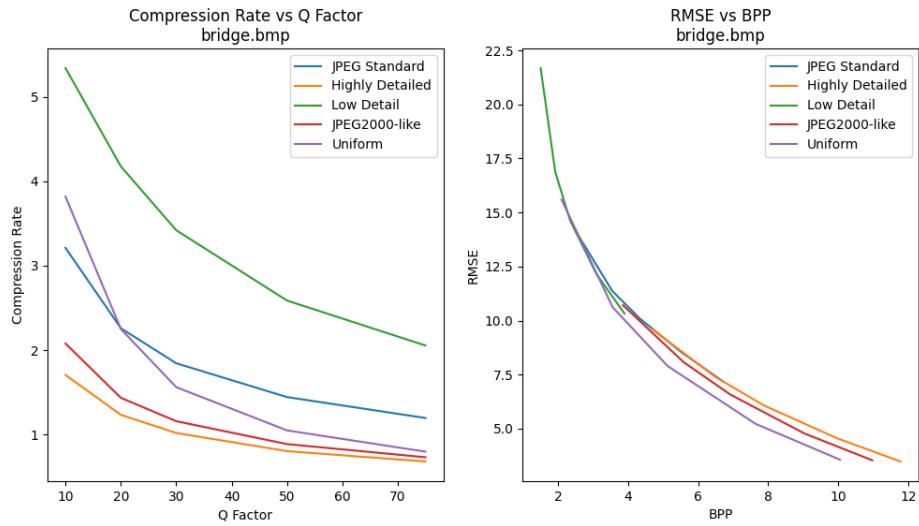


Graphs for metrics comparison.

Bridge

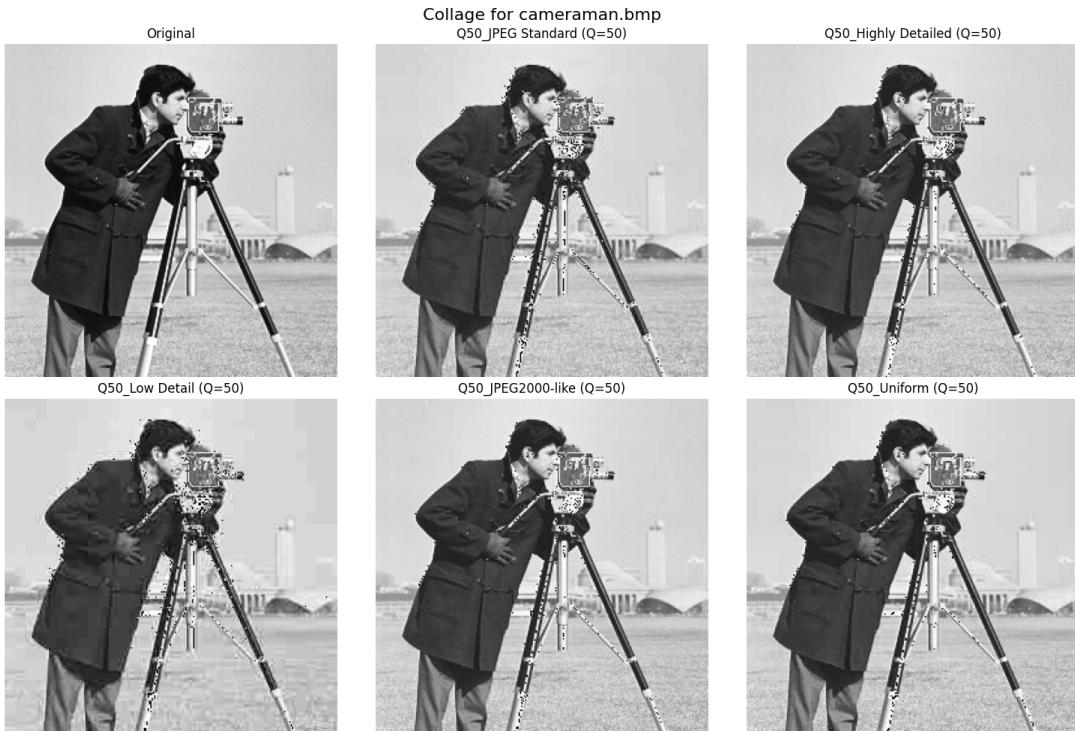


Original image and Q=50 compressed images.

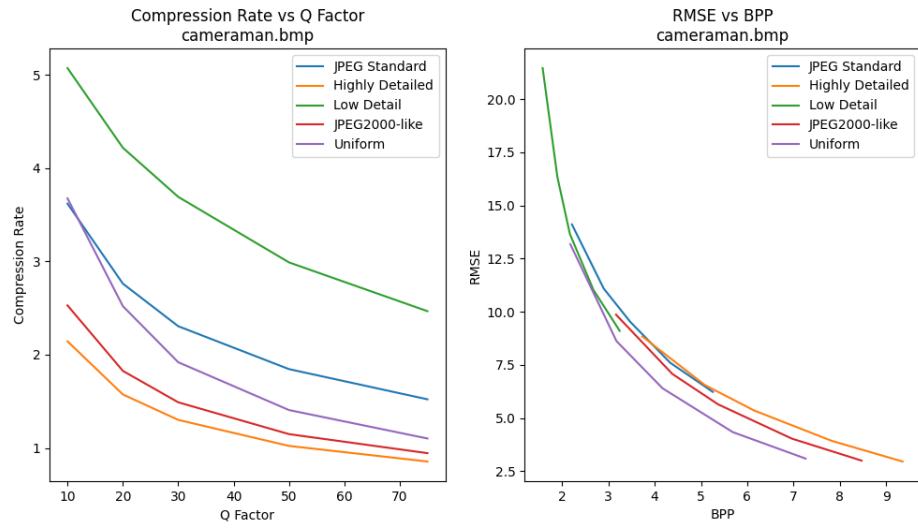


Graphs for metrics comparison.

Cameraman

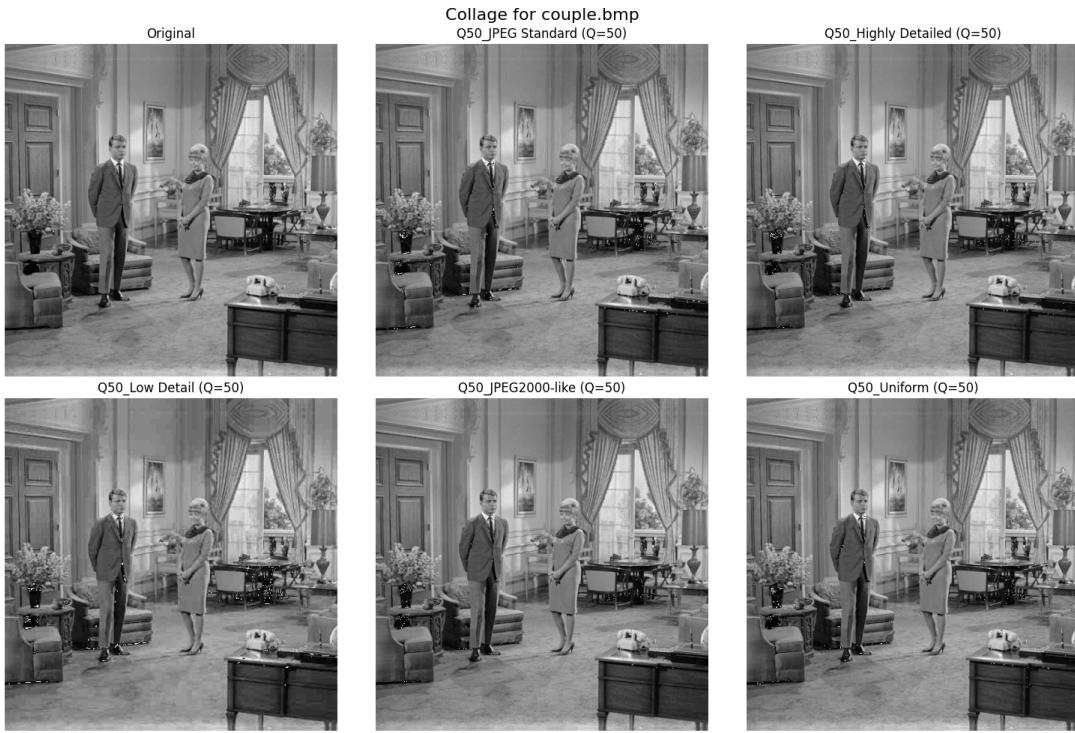


Original image and Q=50 compressed images.

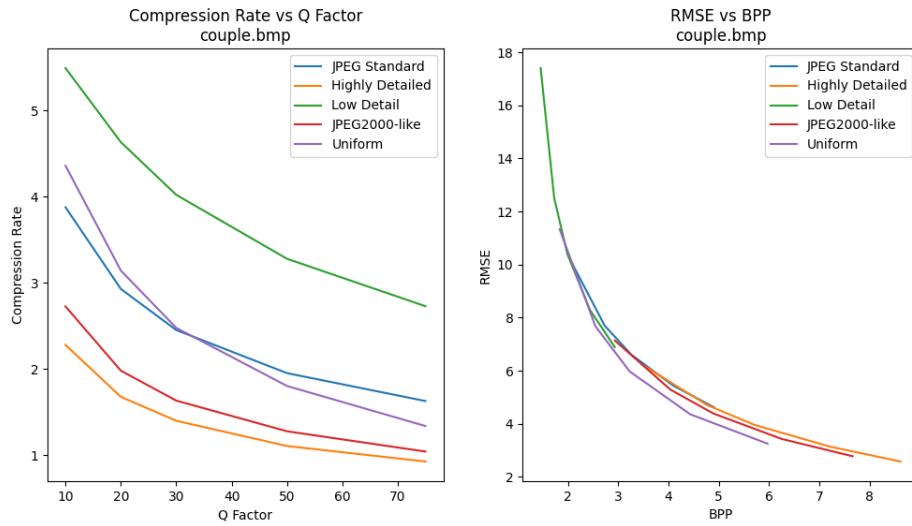


Graphs for metrics comparison.

Couple



Original image and Q=50 compressed images.

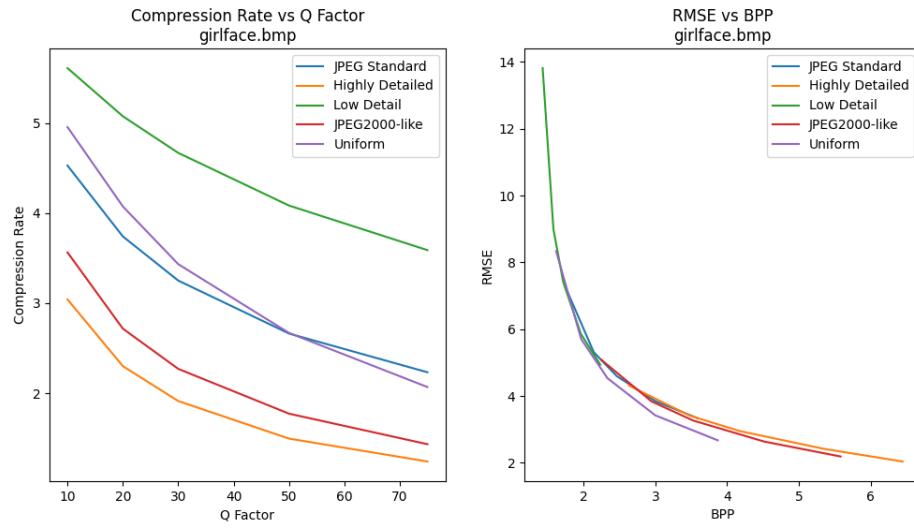


Graphs for metrics comparison.

Girlface

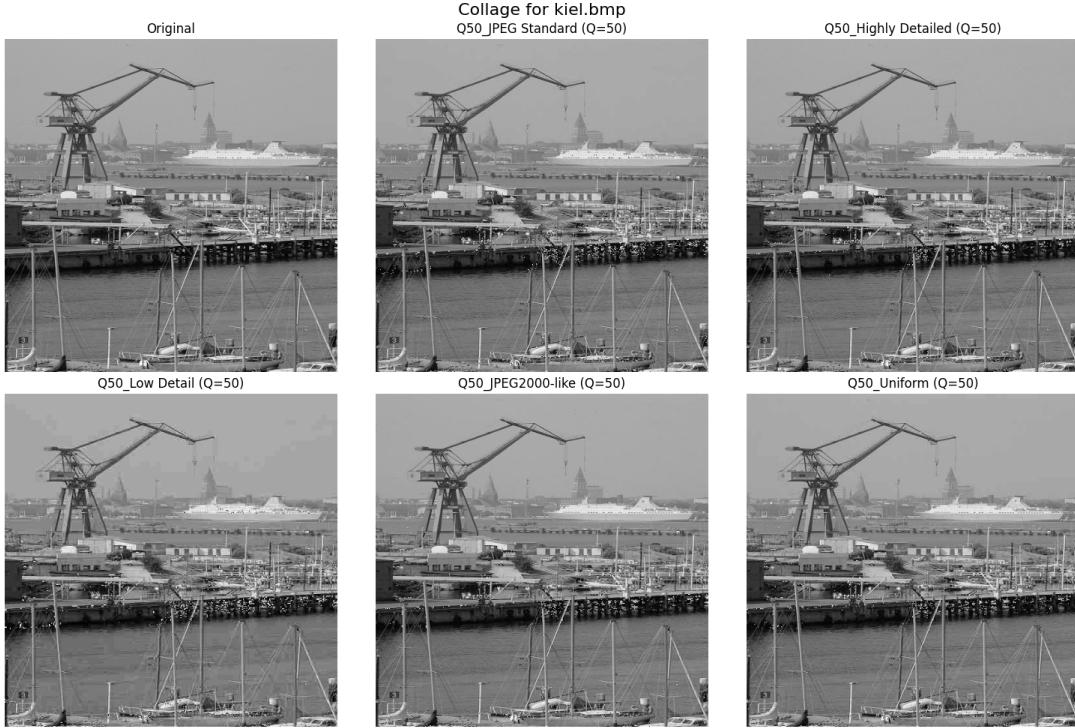


Original image and Q=50 compressed images.

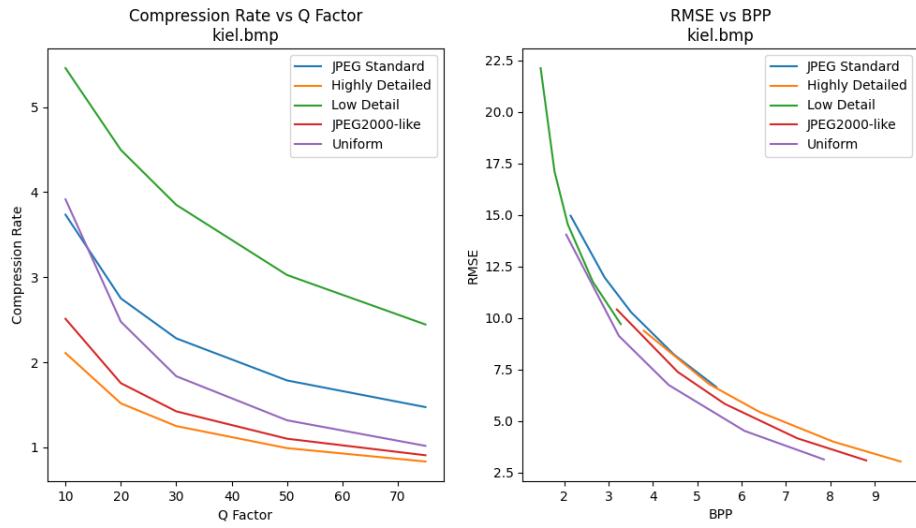


Graphs for metrics comparison.

Kiel

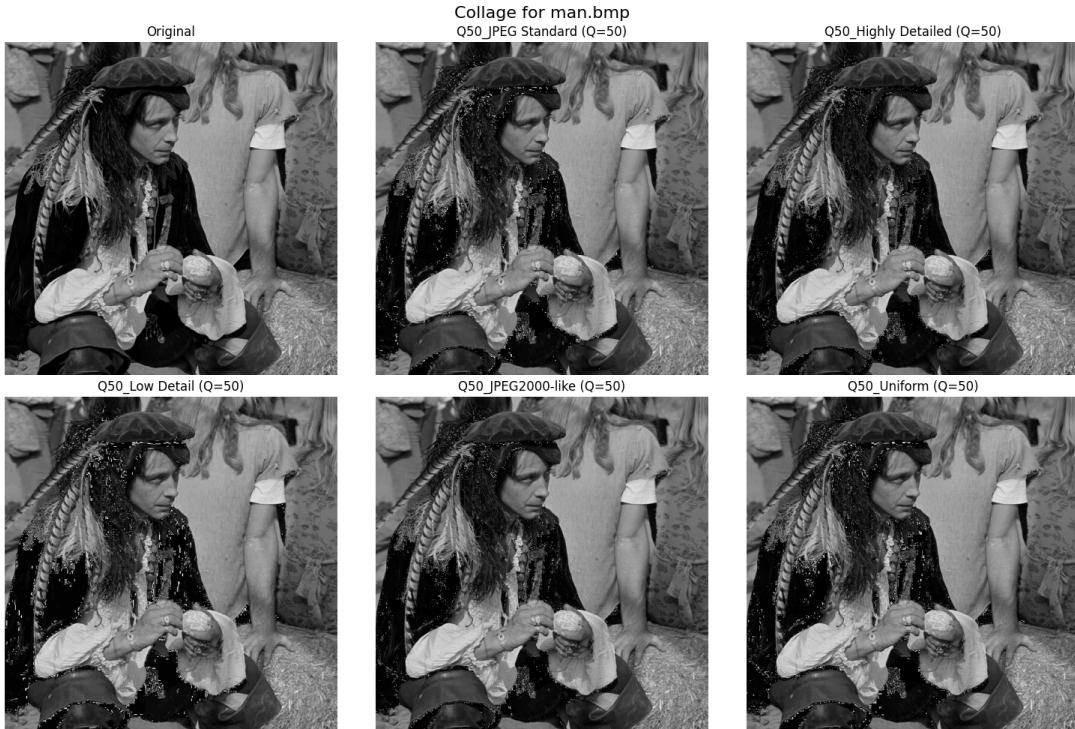


Original image and Q=50 compressed images.

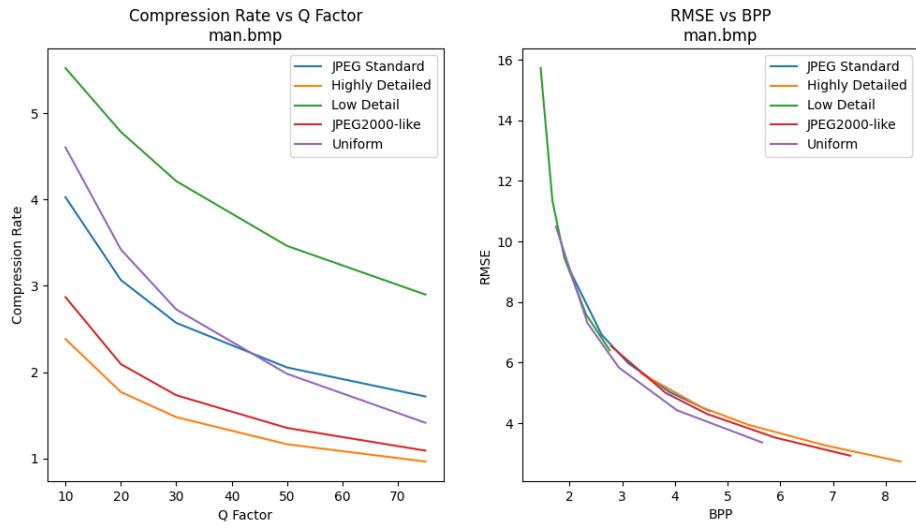


Graphs for metrics comparison.

Man

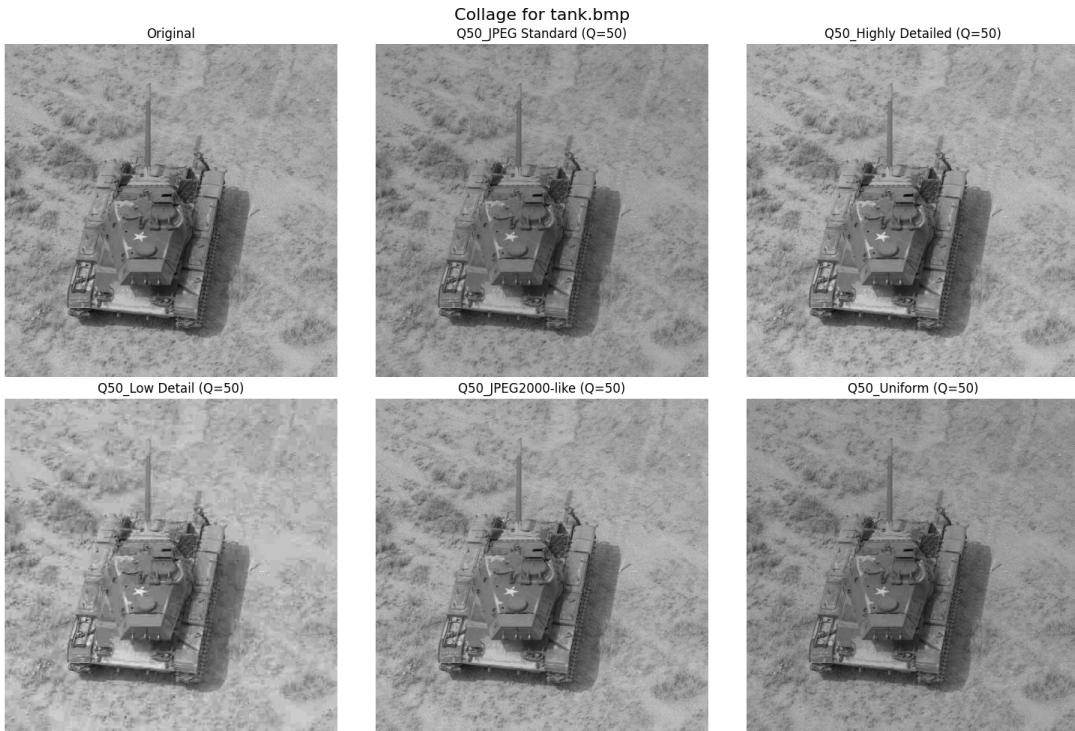


Original image and Q=50 compressed images.

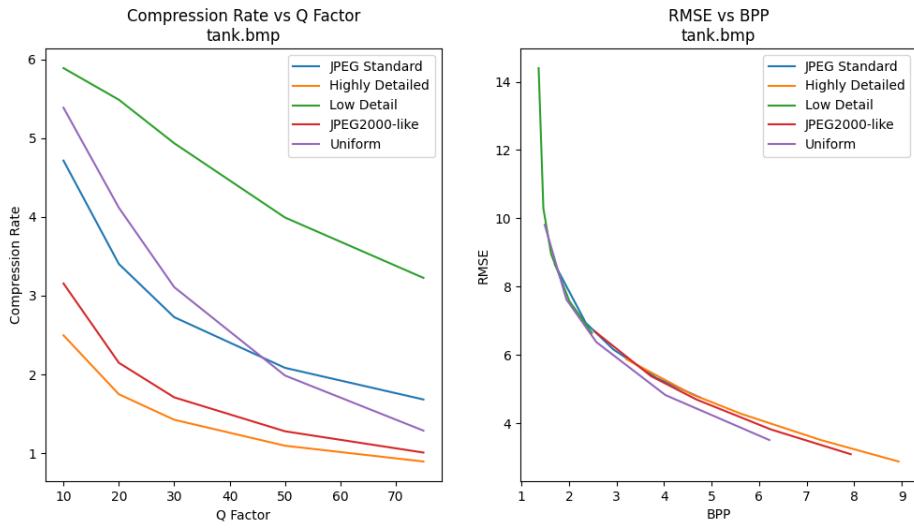


Graphs for metrics comparison.

Tank



Original image and Q=50 compressed images.

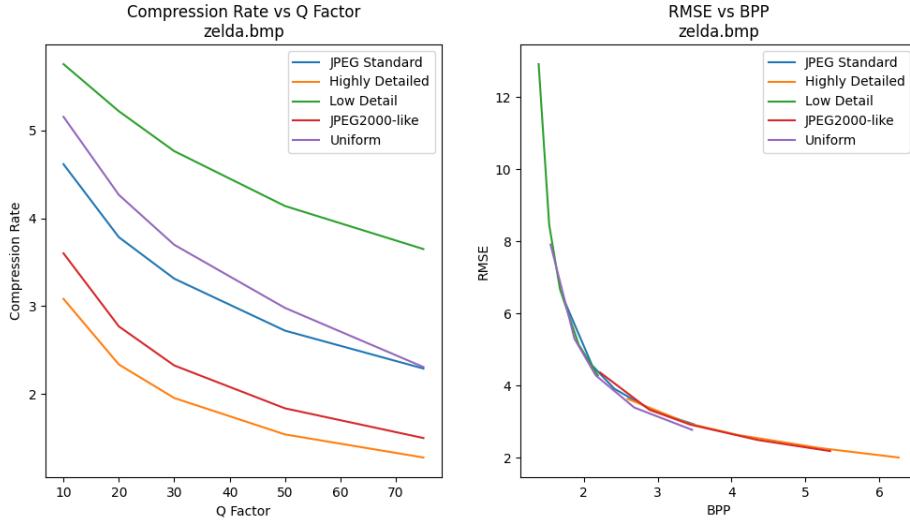


Graphs for metrics comparison.

Zelda



Original image and Q=50 compressed images.



Graphs for metrics comparison.

Analysis

The plots illustrate how the choice of quantization matrix impacts compression performance and quality, highlighting the differences across matrices in achieving a balance between compression rate and reconstruction accuracy.

Compression Rate vs. Q Factor

The left plot shows the compression rate as a function of the Q factor for various quantization matrices. Across all matrices, the compression rate decreases with increasing Q factor, but the extent of this decrease varies:

- The *JPEG Standard* quantization matrix performs the best, achieving a balance between high compression rates and adaptability across Q factors. This reflects its general-purpose design and its ability to allocate bits efficiently across image regions.
- The *Low Detail* matrix exhibits higher compression rates at lower Q factors, as it aggressively removes less critical information. However, its compression advantage diminishes at higher Q factors, where the matrices converge.
- The *Highly Detailed* matrix achieves the lowest compression rates across all Q factors, as it prioritizes retaining intricate details over aggressive compression.
- The *Uniform* and *JPEG2000-like* matrices display intermediate behavior, maintaining compression rates between the extremes. They perform closer to the *JPEG Standard* matrix but exhibit slightly less efficiency.

At low Q factors (10–30), there is greater divergence in compression rates, reflecting the matrices' distinct quantization strategies. At higher Q factors (above 50), the compression rates converge, as the emphasis shifts to preserving image details uniformly across all strategies.

RMSE vs. BPP

The right plot shows how RMSE varies with BPP. As BPP increases, the RMSE decreases for all matrices, but the *JPEG Standard* matrix consistently outperforms the others:

- The *JPEG Standard* matrix achieves the lowest RMSE at nearly all BPP values, demonstrating its superior balance between quality and compression efficiency. This reflects its optimized bit allocation to minimize reconstruction error across image regions.
- The *Highly Detailed* matrix also maintains a low RMSE, especially at low BPP, as it allocates bits to preserve fine details. However, its performance slightly lags behind the *JPEG Standard*.
- The *Low Detail* matrix starts with a significantly higher RMSE at low BPP values due to its aggressive compression. Its RMSE reduces more sharply as BPP increases, but it fails to catch up with the *JPEG Standard* or *Highly Detailed* matrices at higher BPP values.
- The *Uniform* and *JPEG2000-like* matrices perform moderately well, lying between the *Low Detail* and *Highly Detailed* strategies, but do not match the *JPEG Standard*'s quality.

At higher BPP values (above 4), differences in RMSE across matrices diminish, as sufficient bits are allocated to reconstruct most of the image, irrespective of the matrix used.

Impact of Quantization Matrices

The variation across quantization matrices underscores their distinct design goals:

- The *JPEG Standard* matrix optimizes for a general-purpose tradeoff, achieving both high compression rates and low RMSE across the Q factor and BPP ranges.
- The *Low Detail* matrix prioritizes compression efficiency by discarding high-frequency information but sacrifices quality at low BPP values.
- The *Highly Detailed* matrix minimizes RMSE, particularly at low BPP, but at the cost of reduced compression rates.
- Intermediate matrices like *Uniform* and *JPEG2000-like* balance these objectives but do not surpass the *JPEG Standard* in either metric.

6 PCA Implementation and Comparison with JPEG

6.1 Details of Implementation

The PCA class implements Principal Component Analysis (PCA) for compressing and decompressing image data. The class performs the following key operations:

1. `fit` method: Trains the PCA model on input data by computing the mean, covariance matrix, and eigenvectors.
2. `compress` method: Compresses the data by projecting it onto the top principal components.
3. `decompress` method: Reconstructs the original data by applying the inverse of the transformation.

Initialization (`__init__`)

The initialization method accepts the number of components (k) to retain in the compression. It sets this value in the class and prepares the PCA model.

```
1 def __init__(self, n_components):
2     self.n_components = n_components
```

Explanation: The constructor initializes the PCA model with a user-defined number of components. This value will be used during the dimensionality reduction (compression) process.

Fit Method (fit)

The `fit` method computes the principal components of the input data, which involves centering the data, calculating the covariance matrix, and extracting eigenvectors and eigenvalues.

```
1 def fit(self, data):
2     # Center the data by subtracting the mean
3     num_images, num_pixels = data.shape
4     self.mean = np.mean(data, axis=0)
5     centered_data = data - self.mean
6
7     # Calculate the covariance matrix
8     cov_matrix = np.cov(centered_data, rowvar=False)
9
10    # Compute the eigenvalues and eigenvectors
11    eigenvalues, eigenvectors = np.linalg.eigh(cov_matrix)
12
13    # Sort eigenvalues and eigenvectors in descending order
14    sorted_indices = np.argsort(eigenvalues)[::-1]
15    self.eigenvalues = eigenvalues[sorted_indices]
16    self.eigenvectors = eigenvectors[:, sorted_indices]
17
18    # Select the top n_components eigenvectors
19    self.components = self.eigenvectors[:, :self.n_components]
```

Explanation: - The data is centered by subtracting the mean of each feature (pixel). - The covariance matrix of the centered data is calculated to capture the variance and correlation between features. - Eigenvalues and eigenvectors are computed using `np.linalg.eigh`, which is specifically suited for symmetric matrices (like covariance matrices). - The eigenvectors are sorted in descending order based on their corresponding eigenvalues, which represent the variance explained by each principal component. - The top $n_{\text{components}}$ eigenvectors are selected as the principal components.

Mathematically, the process involves:

$$\text{Centering: } \mathbf{X}_{\text{centered}} = \mathbf{X} - \boldsymbol{\mu}$$

where $\boldsymbol{\mu}$ is the mean of the data matrix \mathbf{X} .

The covariance matrix Σ is calculated as:

$$\Sigma = \frac{1}{n} \mathbf{X}_{\text{centered}}^T \mathbf{X}_{\text{centered}}$$

The eigenvalue decomposition of the covariance matrix yields the eigenvectors \mathbf{v} and eigenvalues λ :

$$\Sigma \mathbf{v}_i = \lambda_i \mathbf{v}_i$$

where λ_i is the eigenvalue corresponding to the eigenvector \mathbf{v}_i .

Compress Method (compress)

The `compress` method reduces the dimensionality of the data by projecting it onto the principal components. It also saves the compressed data, components, and mean to a binary file for future use.

```
1 def compress(self, data, output_file="compression_output.npy"):
2     # Center the data (subtract the mean)
3     centered_data = data - self.mean
4
5     # Compress the data by projecting it onto the components
6     compressed_data = np.dot(centered_data, self.components)
7
8     compressed_data = compressed_data.astype("int16")
9     self.components = self.components.astype("float16")
10    self.mean = self.mean.astype("uint8")
11
12    # Store the data into a binary file using np.save
```

```

13     with open(output_file, 'wb') as f:
14         np.save(f, compressed_data) # Save compressed data
15         np.save(f, self.components) # Save components
16         np.save(f, self.mean) # Save mean

```

Explanation: - The data is centered by subtracting the mean, similar to the `fit` method. - The data is projected onto the principal components by performing a dot product with the eigenvectors:

$$\mathbf{X}_{\text{compressed}} = \mathbf{X}_{\text{centered}} \cdot \mathbf{V}$$

where \mathbf{V} is the matrix of the top $n_{\text{components}}$ eigenvectors. - The compressed data is then cast into a lower precision format to save memory space, specifically converting to `int16`, `float16`, and `uint8` for the compressed data, components, and mean, respectively. - The `np.save` function is used to save the compressed data, components, and mean to a binary file.

Decompress Method (`decompress`)

The `decompress` method reconstructs the original data from the compressed data by applying the inverse transformation (dot product with the transpose of the components and adding the mean).

```

1 def decompress(self, file_name="compression_output.npy"):
2     # Open the file in read mode and load the data
3     with open(file_name, 'rb') as f:
4         compressed_data = np.load(f) # Load compressed data
5         components = np.load(f) # Load components
6         mean = np.load(f) # Load mean
7
8     # Reconstruct the data
9     decompressed_data = np.dot(compressed_data, components.T) + mean
10    return decompressed_data

```

Explanation: - The compressed data, components, and mean are loaded from the binary file. - To reconstruct the data, the compressed data is multiplied by the transpose of the components matrix:

$$\mathbf{X}_{\text{reconstructed}} = \mathbf{X}_{\text{compressed}} \cdot \mathbf{V}^T + \mu$$

where \mathbf{V}^T is the transpose of the matrix of principal components, and μ is the mean.

6.2 Comparison for Original Dataset

The graphs compare the performance of PCA and a custom implementation of the JPEG compression algorithm in terms of RMSE versus BPP. The x-axis represents the BPP, which indicates the amount of information retained after compression. The y-axis denotes the RMSE, a measure of reconstruction error. A lower RMSE signifies higher image fidelity after compression. The two curves, plotted for PCA and the custom JPEG algorithm, reveal distinct characteristics of the two approaches to compression.

Performance of PCA

The blue curve represents the PCA-based compression. Initially, at very low BPP values, PCA shows very high RMSE, indicating poor reconstruction quality. However, as BPP increases, the RMSE drops significantly. PCA demonstrates a smooth curve, reflecting its ability to distribute the retained information across principal components effectively. The PCA curve converges to a relatively low RMSE at higher BPP, but it requires significantly higher BPP compared to JPEG to achieve comparable reconstruction quality.

Performance of Custom JPEG

The red curve represents the custom JPEG implementation. JPEG compression achieves a consistently lower RMSE than PCA at equivalent BPP values. This result highlights the efficiency of JPEG in leveraging the Discrete Cosine Transform (DCT) and quantization tables to optimize compression while retaining perceptual quality. The Custom JPEG curve flattens earlier than PCA, indicating that it achieves near-optimal quality at lower BPP values.

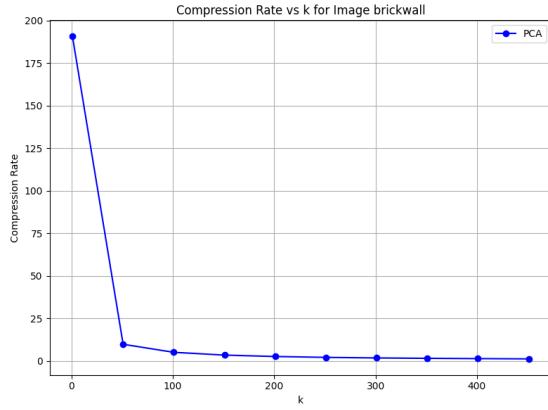
Comparison Between PCA and JPEG

For a given BPP, the custom JPEG consistently outperforms PCA in terms of RMSE. This is especially evident at lower BPP values, where PCA struggles with higher reconstruction errors. The superior performance of JPEG can be attributed to its design, which incorporates human visual system characteristics, prioritizing perceptually important information.

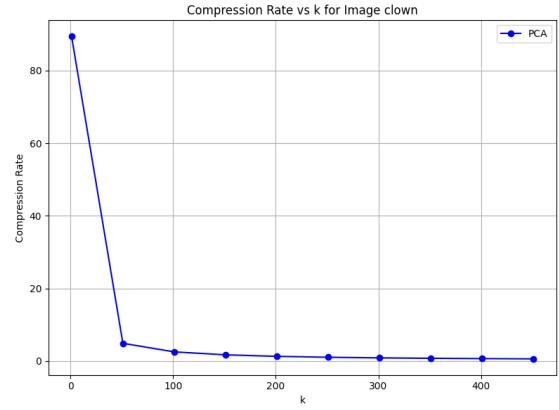
Trade-off Analysis

The trade-off between compression and reconstruction quality is more favorable for JPEG at lower BPP. PCA might be preferable in scenarios where high dimensionality reduction is required, but it is less effective for perceptual image quality preservation compared to JPEG.

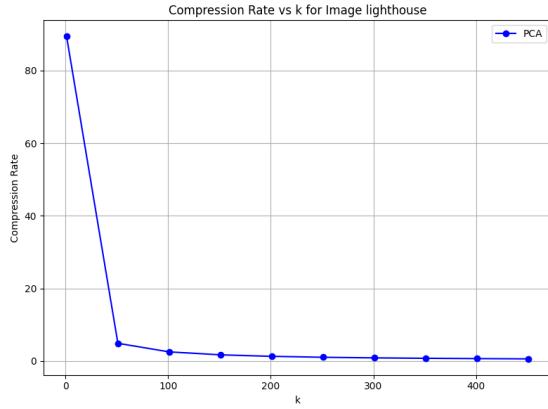
Compression Rate vs. k for a few images



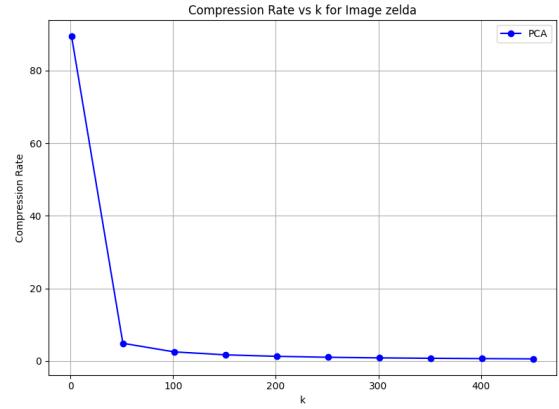
(a) Image 1: Compression Rate vs. k



(b) Image 2: Compression Rate vs. k

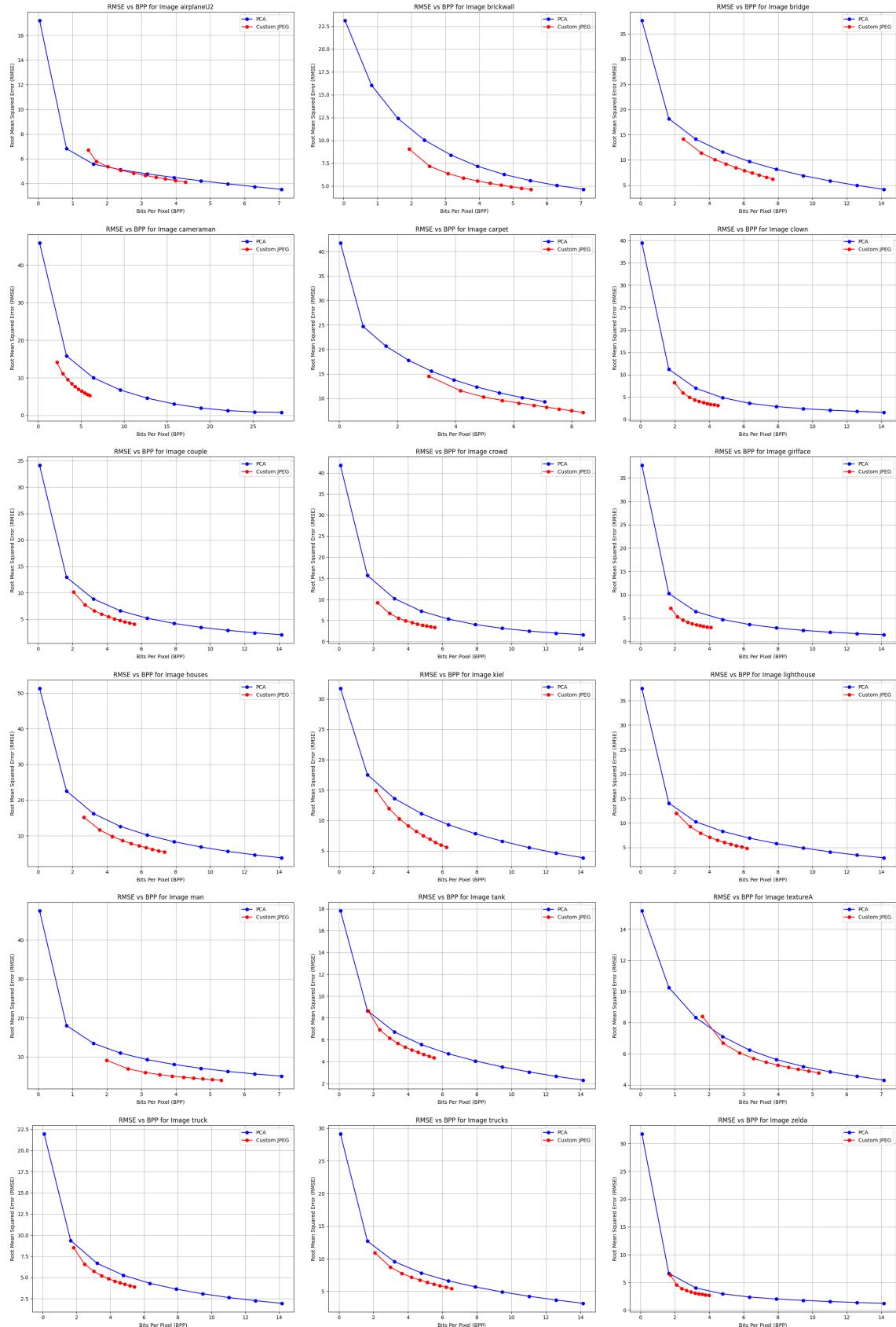


(c) Image 3: Compression Rate vs. k



(d) Image 4: Compression Rate vs. k

Figure 20: Compression Rate vs. k for a few images. Each graph shows how the compression rate evolves as k increases. The compression rate decreases sharply at first and stabilizes as k increases, indicating diminishing returns in compression efficiency with higher values of k .



7 Sparse Orthonormal Transforms and Comparison with Custom JPEG

7.1 Details of Implementation

The `JPEG_SOT` class is essentially the same as the `JPEG` class as defined before, with just a few changes, which will be discussed. The algorithms were trained on 8x8 blocks of two images, astronaut and gravel from the skimage data module.

Algorithm 1

To implement Algorithm 1, which is essentially convergence of a particular transformation matrix given blocks, we had to find the optimal coefficients for the blocks given a transformation matrix, and the optimal transformation matrix given these coefficients, and continue this cycle until the matrix converges. This is done as follows.

```
1 def get_opt_coeff_given_transform(G_i, blocks, lam):
2     """
3         G_i is NxN, blocks should be CxN (flattened from sqrtN x sqrtN)
4         return c_ij (should be CxN)
5     """
6     coeffs = blocks @ G_i
7     coeffs[np.abs(coeffs) < np.sqrt(lam)] = 0
8
9     return coeffs
10
11 def get_opt_transform_given_coeff(c_i, x):
12     """
13         x should be a numpy array CxN, where C is total no of blocks, c_i should be a numpy
14             array CxN, corresponding to x
15     """
16     Y = c_i.T @ x
17     U, _, V = np.linalg.svd(Y, full_matrices=False)
18
19     G_i = V @ U.T
20     return G_i
21
22 def algorithm_1(transform, lam, blocks):
23     """
24         transform NxN, coeffs CxN, blocks CxN
25     """
26     prev_transform = np.zeros(transform.shape)
27     coeffs = None
28     while(np.mean((transform - prev_transform)**2) > 0.05):
29         prev_transform = transform
30         coeffs = get_opt_coeff_given_transform(G_i=transform, blocks=blocks, lam=lam)
31         transform = get_opt_transform_given_coeff(c_i=coeffs, x=blocks)
32
33     return transform
```

The formulae required and the proof of optimality for the equations above have been covered in the paper.

Algorithm 2

Algorithm 2 uses Annealing of the lambda parameter, and Algorithm 1 as a sub-process, to perform a clustering sort of process for each block, and transformation matrix. Essentially, we find eight optimal transformations using algorithm 1, perform re-labelling of blocks according to the lowest cost, and then lower the value of lambda and repeat the process, until the labellings converge completely. The code has been added below.

```
1 def algorithm_2(transforms, blocks, labels, lam_max=200, lam_min=20, lam_diff=20):
2     lam = lam_max
3     old_labels = np.zeros_like(labels)
4     cnt = 0
```

```

5     while(not np.array_equal(old_labels, labels)):
6         cnt += 1
7         if cnt > 100:
8             break
9         old_labels = labels
10        while lam > lam_min:
11            for i in range(8):
12                reqd_blocks = blocks[np.where(labels == i)[0]]
13                transforms[i] = algorithm_1(transform=transforms[i], lam=lam, blocks=
14                    reqd_blocks)
15                for i in range(blocks.shape[0]):
16                    costs = np.zeros(8)
17                    for k in range(8):
18                        coeff = get_opt_coeff_given_transform(transforms[k], blocks[i], lam)
19                        costs[k] = np.sum((blocks[i]-coeff@transforms[k])**2)+lam*(np.
20                            count_nonzero(coeff))
21                    new_label = np.argmin(costs)
22                    labels[i] = new_label
23
24        lam -= lam_diff
25        lam = lam_max
26        print(cnt)
27        print(np.sum(old_labels != labels))
28    return transforms

```

We found that using lower values of lambda during training and higher values of lambda during compression is beneficial. Here, the loop which calls Algorithm 1 finds the optimal transformations, and the loop after that performs cost minimisation and re-labelling (essentially performing clustering, with transformation matrices as the centers).

Changes to JPEG Class

Now that we have our transformation matrices after training, only trivial changes need to be made to our original JPEG class. We need to find and store the optimal transformation matrix corresponding to each block, and perform the transformation and inverse transformation using that matrix. This is done by changing the `_dct2` and `_idct2` functions from the original class. The changes have been shown below. Our `JPEG_SOT` class uses a default lambda of 10000.

```

1 def t(self, block):
2     flat_block = block.flatten()
3     flat_block = np.reshape(flat_block, (1,64))
4     costs = np.zeros(8)
5     for i in range(8):
6         coeff = get_opt_coeff_given_transform(self.transforms[i], flat_block, self.lam)
7         costs[i] = np.sum((flat_block-coeff@self.transforms[i])**2)+self.lam*(np.
8             count_nonzero(coeff))
9     best_transform = self.transforms[np.argmin(costs)]
10    return np.reshape(get_opt_coeff_given_transform(best_transform, flat_block, self.lam),
11        (8,8)), best_transform
12
13 def it(self, block, transform):
14     return np.reshape((transform @ block.flatten()), (8,8))

```

Essentially, we are creating 64x64 transformation matrices, and hence we use the blocks as 64x1 matrices. We find the transformation matrix with the lowest cost, and return the best transformation matrix for that block, along with the optimal coefficients for it. During the inverse step, the optimal matrix simply needs to be multiplied with the coefficients to get back the image. The storage of matrices was done by adding the following lines. Note that the lines below are simply changes to the original class and not complete functions.

```

1 def __init__(self, Q = 50, lam=10000):
2     self.transforms = transforms
3     self.block_to_transform = {}
4     self.lam = lam
5
6 def compress(self, image, filename):

```

```

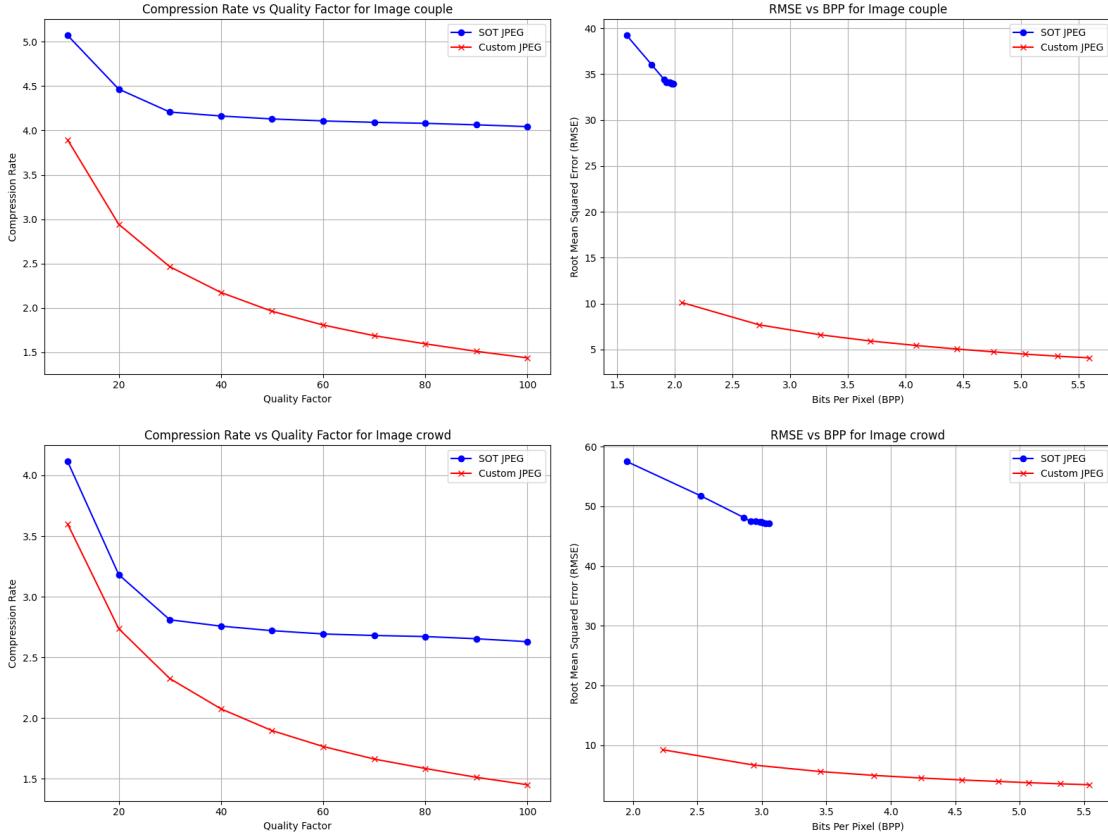
7     for i, chunk in enumerate(chunks):
8         val, transform = self.t(chunk)
9         quantized_matrix = self._quantize(val)
10        flattened_matrix = self._zigzag(quantized_matrix)
11        flattened_matrices.append(flattened_matrix)
12        self.block_to_transform[i] = transform
13
14
15    def decompress(self, filename):
16        for i in range(len(flattened_matrices)):
17
18            flattened_matrices[i] = self._reverse_zigzag(flattened_matrices[i], 8, 8)
19
20            # now de-quantize the matrix and rescale
21            flattened_matrices[i] = (flattened_matrices[i] * self.scaled_quantization_matrix)
22            flattened_matrices[i] = self.it(flattened_matrices[i], self.block_to_transform[i]) +
23                128

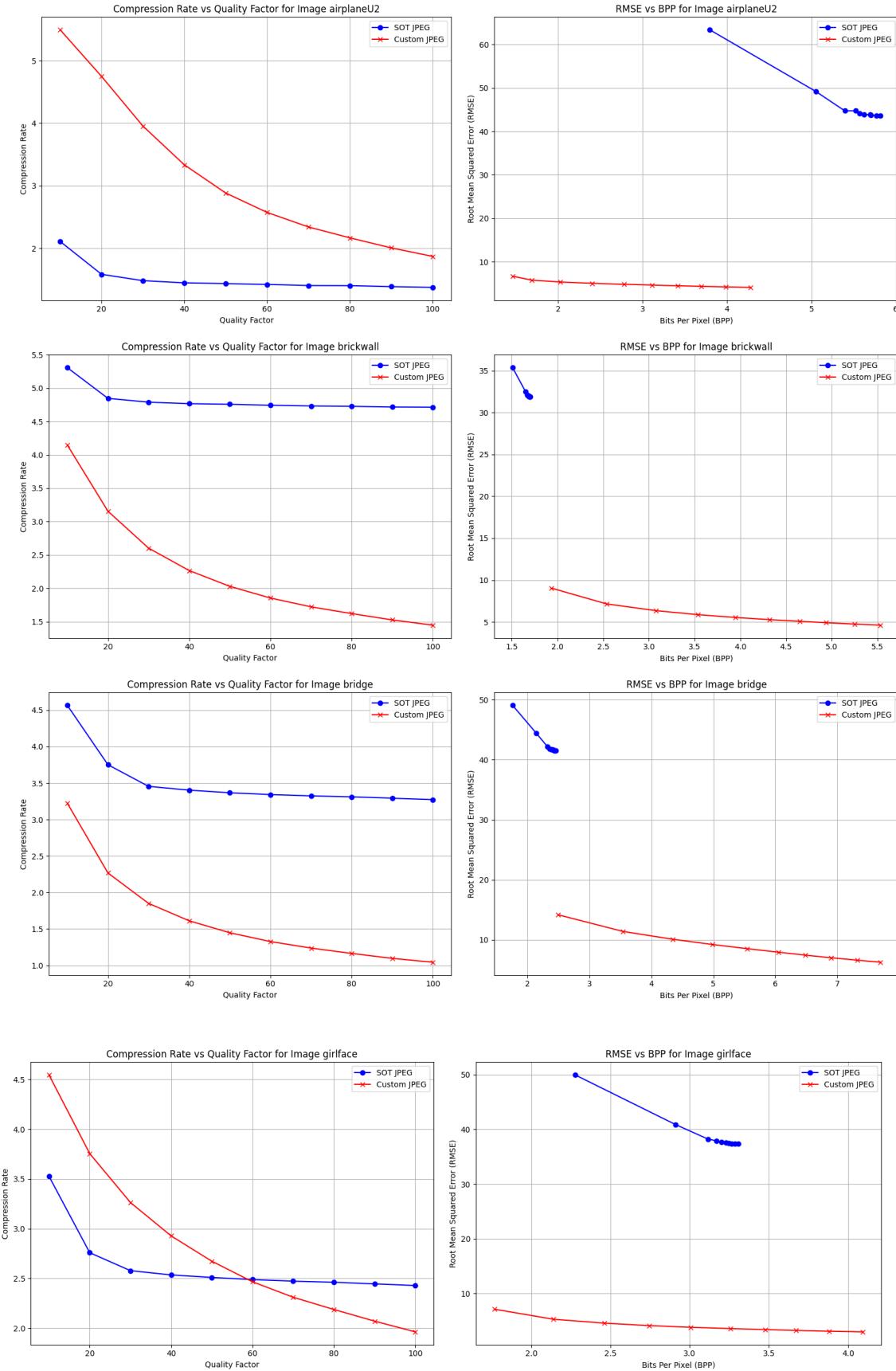
```

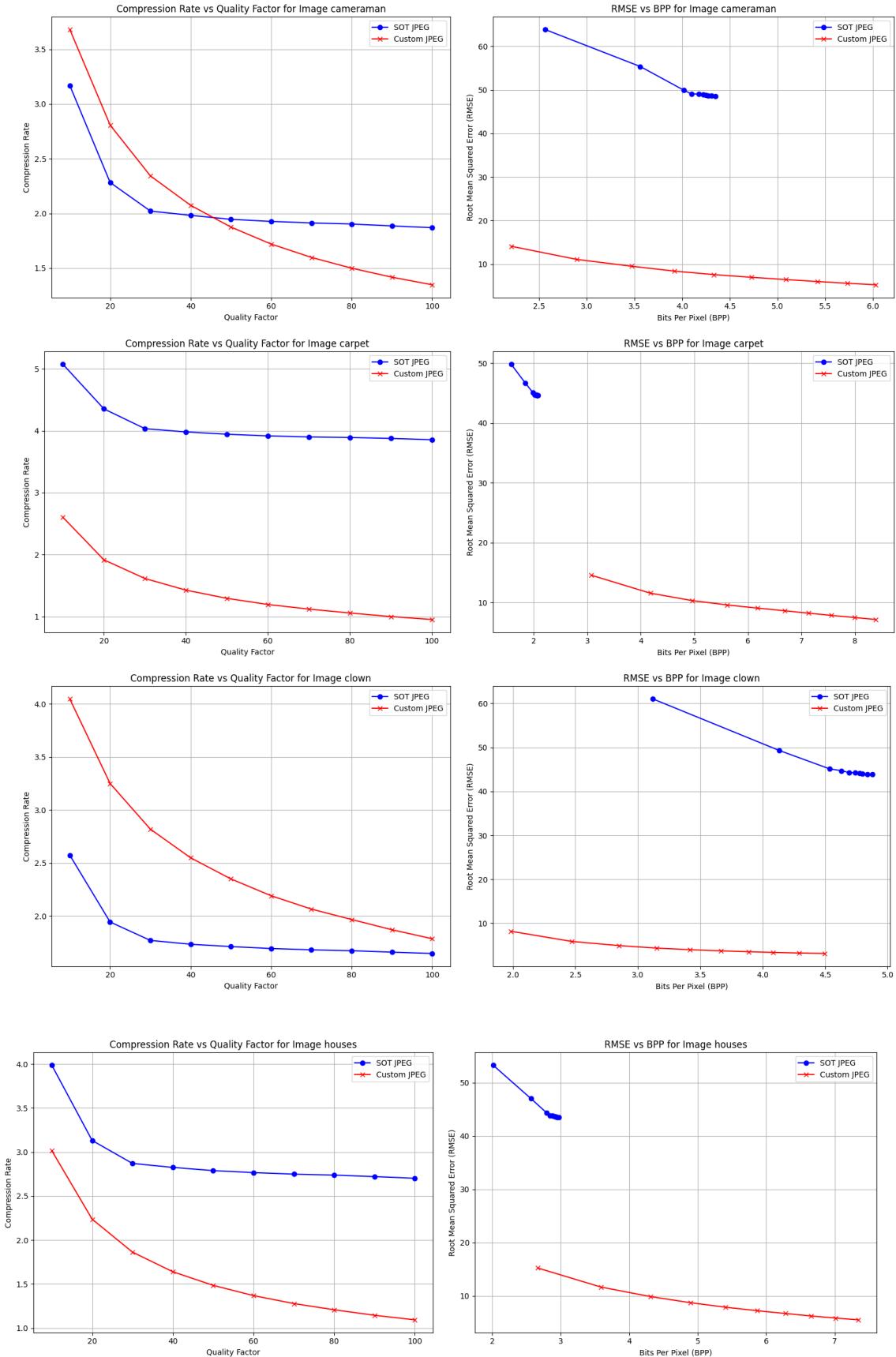
7.2 Comparison with Custom JPEG

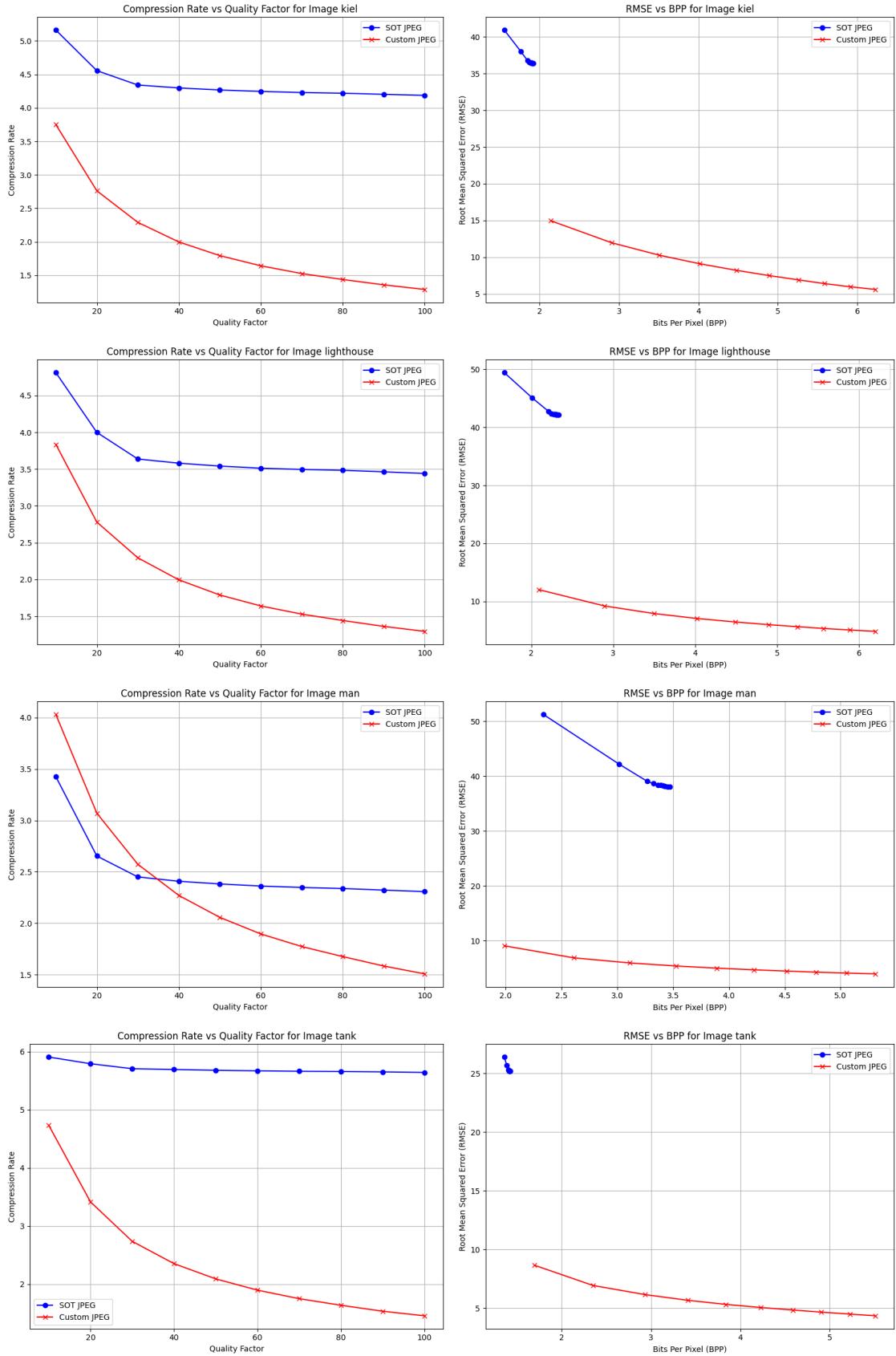
Since we use SOT on top of our custom JPEG class, we would want to compare the performance of trained transformations vs direct DCT applications. This comparison would give a direct answer as to whether there are improvements.

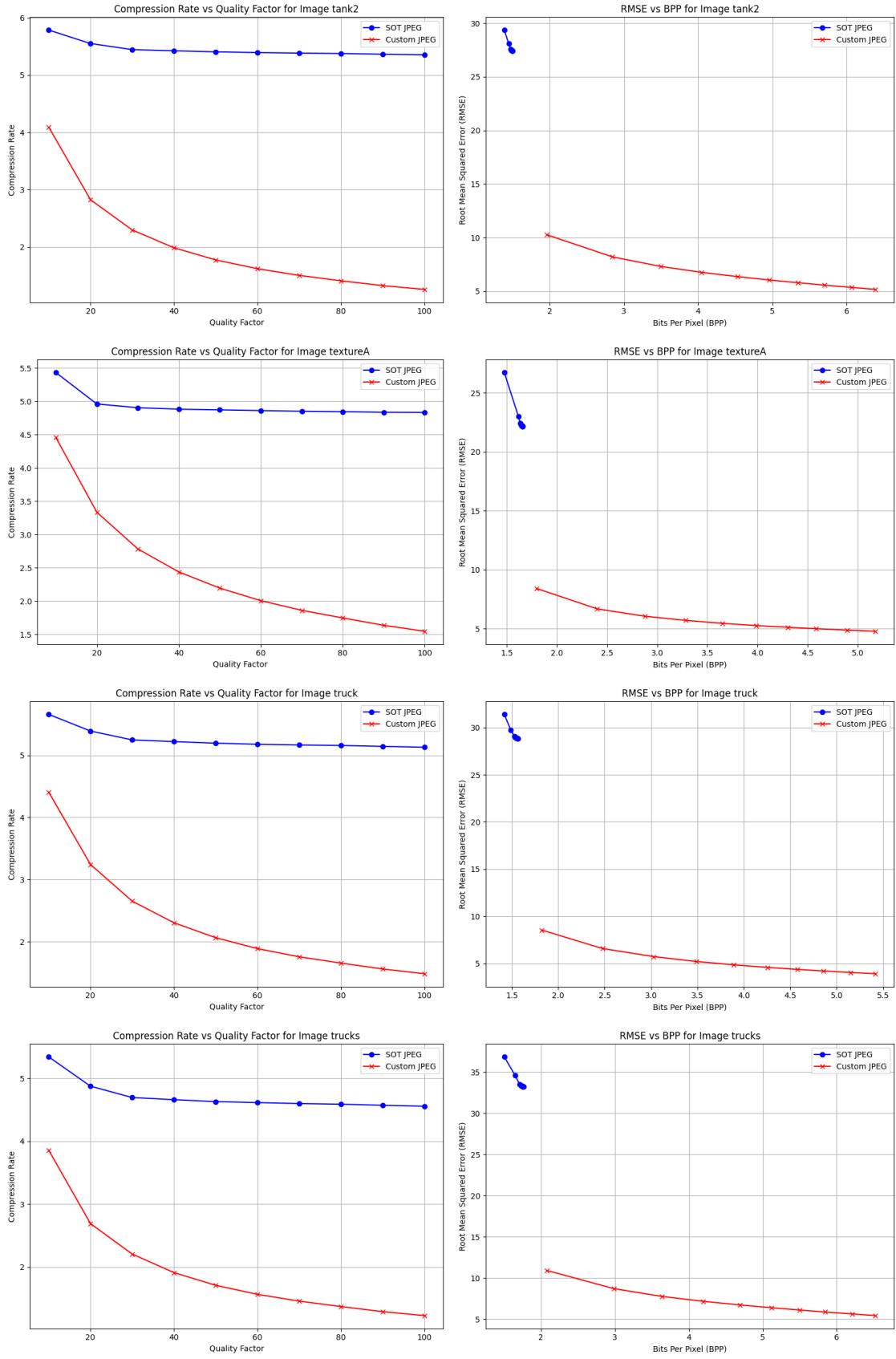
As we can see from the graphs below, on most occasions, SOT outperforms DCT in terms of compression rate, which means that it would give noticeable improvements over standard JPEG, especially with further hyperparameter tuning and more training time. However, the BPP for SOT does not vary much at all, and the RMSE is quite high as well. The results run on the 19 images are shown below.

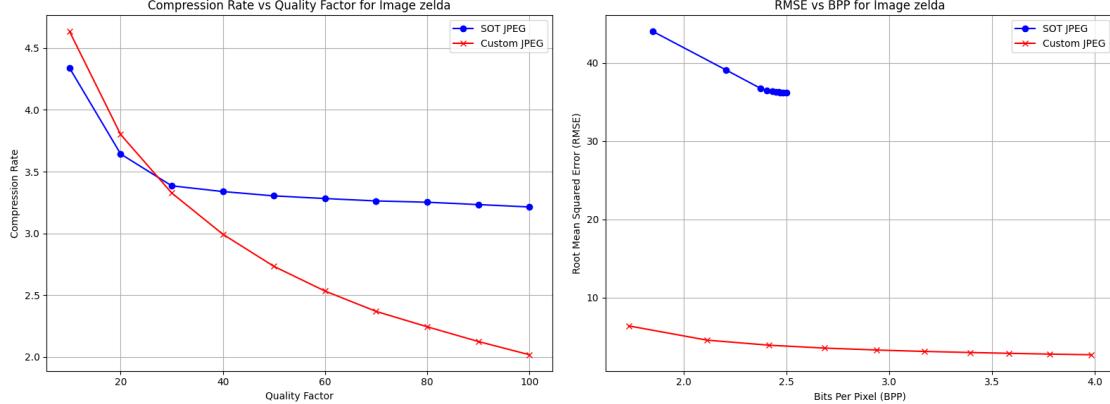












8 Conclusions and Wrapping Up

Our compression scheme doesn't compare too well to the standard cv2 implementation. This could be for a number of reasons. Most notably, the fact that we have a custom compression and encoding scheme which is most likely filled with inefficiencies. We have also left out the AC/DC coefficients in our implementation which just adds complexity to only marginally improve the compression rate.

On the bright side, we have RMSE values somewhat comparable to the standard cv2 implementation for varying Q. Our implementation is also robust across various different kinds of images, achieving reasonable compression even on "challenging" images where even cv2 fails to get good compression.

Our JPEG implementation also outperforms a naive implementation of PCA Compression in most natural images, which is a good sign.

In addition, we saw positive results with the SOT implementation, which could be added to a well performing JPEG implementation and hence, improve it.

9 Contributions

Everyone contributed equally towards all parts of the project, however the main implementation work was split equally, which has been listed below.

- Brian Mackwan
 1. PCA Implementation
 2. Extensive Analysis and Experiments in all implementations
- Ekansh Ravi Shankar
 1. Sparse Orthogonal Transformation Implementation
 2. Extensive Analysis and Experiments in all implementations
- Ved Danait
 1. Basic JPEG Implementation
 2. Extensive Analysis and Experiments in all implementations