

High-Perceptual Quality JPEG Decoding via Posterior Sampling

Sean Man Guy Ohayon Theo Adrai Michael Elad

Technion - Israel Institute of Technology, Haifa, Israel

{sean.man, ohayonguy, theoad}@campus.technion.ac.il, elad@cs.technion.ac.il

Abstract

JPEG is arguably the most popular image coding format, achieving high compression ratios via lossy quantization that may create visual artifacts degradation. Numerous attempts to remove these artifacts were conceived over the years, and common to most of these is the use of deterministic post-processing algorithms that optimize some distortion measure (e.g., PSNR, SSIM). In this paper we propose a different paradigm for JPEG artifact correction: Our method is stochastic, and the objective we target is high perceptual quality – striving to obtain sharp, detailed and visually pleasing reconstructed images, while being consistent with the compressed input. These goals are achieved by training a stochastic conditional generator (conditioned on the compressed input), accompanied by a theoretically well-founded loss term, resulting in a sampler from the posterior distribution. Our solution offers a diverse set of plausible and fast reconstructions for a given input with perfect consistency. We demonstrate our scheme’s unique properties and its superiority to a variety of alternative methods on the FFHQ and ImageNet datasets.

1. Introduction

JPEG (Joint Photographic Experts Group) [53] is one of the most popular *lossy* image compression techniques, extensively used in digital cameras, internet communications and more. JPEG reduces image file-size by discarding information that is supposed to be less valuable to a human observer. To achieve high compression ratios, JPEG often discards noticeable visual details, which may lead to strong artifacts in the decompressed image, such as blockiness.

Since its conception in the late 80’s, numerous post-processing algorithms were proposed for removing JPEG artifacts.¹ Such methods start with the given compressed image and somehow provide an estimation of the source image. A common estimation approach attempts to minimize the average discrepancy between the source image

¹While it is outside the scope of this paper to reference this vast literature, we do provide links to leading such techniques in [Section 2](#).

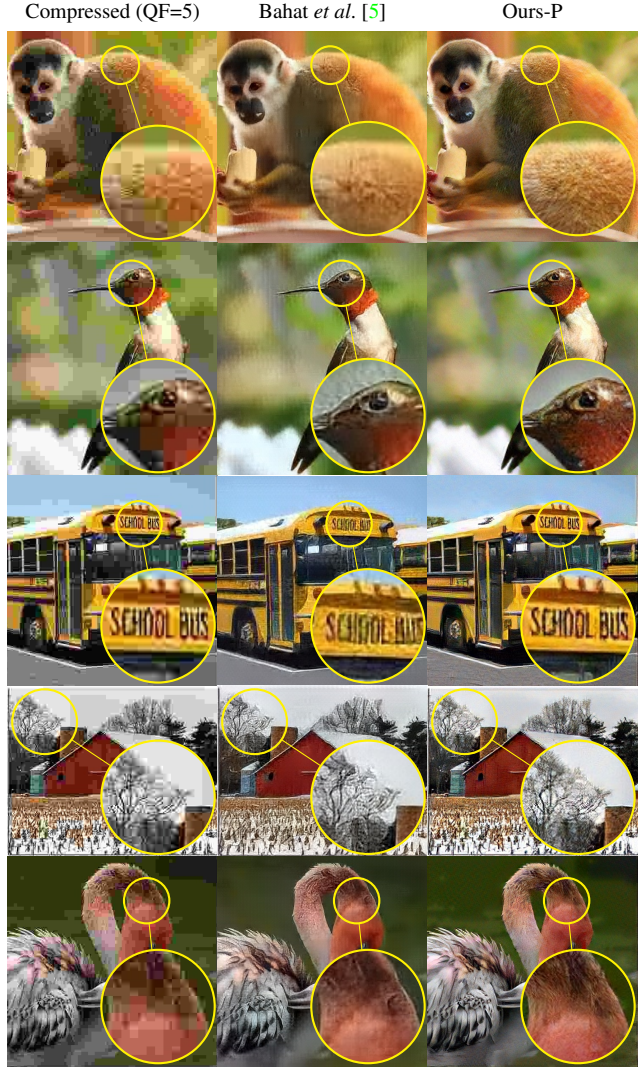


Figure 1. Reconstruction examples of highly compressed JPEG images using [5] and our proposed method. Our method produces stochastic outputs, all of which are perfectly consistent with the compressed image, like [5], but with far better perceptual quality.

and the recovered one, *e.g.*, minimizing the Mean-Squared-Error (MSE). While this strategy may improve the quality of the compressed input and remove the blockiness effect, it still evidently leads to visually unpleasing images that are often accompanied by blurriness. As shown in [6], this is a direct manifestation of the perception-distortion tradeoff, which is apparently noticeable and strict for natural images and common distortion measures – low MSE contradicts high perceptual quality.

Another key aspect in JPEG restoration processes is the consistency of the reconstructed image with the compressed input: The compressed version of the output image should be identical to the compressed input. This is a reasonable requirement as any inconsistent image could never be the true source. Interestingly, while older classical methods tend to be consistent (*e.g.* using iterative projection algorithms [31, 38, 39, 49]), this feature has been overlooked in recent learning-based recovery techniques (such as [8, 14, 15, 18]).

In this work we propose a new JPEG reconstruction algorithm that produces consistent outputs, while also attaining very high perceptual quality. The produced images by our algorithm are free of artifacts such as blockiness, blurriness, etc. However, due to the perception-distortion tradeoff [6], this comes at the unavoidable cost of reduced distortion performance. Example reconstructions from our method are presented in Figure 1. Our approach is based on a Generative Adversarial Network (GAN) [17], conditioned on the compressed image, and having a stochastic synthesis process. As such, the generator is able to produce a variety of plausible and consistent reconstructions for a given compressed input. This conditional GAN is trained with a loss comprised of (i) an adversarial term that promotes high perceptual quality, (ii) a penalty term adopted from [37] that promotes output variability (per input), and (iii) a consistency term that promotes a compliance with the compressed input. Our optimization task theoretically admits the posterior sampler as a unique solution, leading to perfect perceptual quality and perfect consistency with the compressed input.

The contributions of this paper are the following: (i) We propose a novel algorithm that produces consistent JPEG reconstructions with high perceptual quality, leading to state-of-the-art (SoTA) results on ImageNet; (ii) We identify an empirical perception-consistency tradeoff of stochastic estimators, extending the results presented in [36]. Through a consistency penalty term (instead of enforcing perfect consistency) we are able to traverse this empirical tradeoff; (iii) We analyze and compare our method with previous approaches that targetted the same problem, while not being aware of this empirical tradeoff.

2. Related work

A plethora of ideas on JPEG’s artifacts removal were published since the late 80’s, typically offering a post-processing mechanism within the decoding stage. These include techniques such as simple spatial filters [11, 28, 41], MSE optimization of codebooks [22, 57], MAP optimization using Gibbs priors [29, 39, 49], sparse representations [10, 25] and many other techniques (*e.g.*, see the review by Shen and Kuo [46]).

In recent years, deep-learning-based methods took the lead in post-processing JPEG compressed images. Dong *et al.* [14] were the first to suggest a CNN-based regression model, followed by work such as [8] that improved the architecture and the loss function. Guo *et al.* [18] combined pixel domain and DCT domain sub-networks to take advantage of JPEG’s mode of operation, followed by work such as [56, 58] that improved the dual-domain concept and [15, 50] that took a step further by operating solely in the DCT domain. All of these methods adopt a supervised approach, training a network to best fit a given compressed image to its desired ideal output.

As noted in Section 1, a common difficulty of these techniques (classical or deep-learning based) is the overly smoothed reconstructed images that lack fine texture and high frequency details, which occurs due to the inherent tradeoff between perceptual quality and an optimization of a distortion criterion [6]. Indeed, some methods [15, 16, 19] incorporated adversarial training [17] to produce sharp details while not being aware of this tradeoff, hence kept incorporating distortion measures as part of their loss, and thus compromising on the perceptual quality of their outcome as well.

Among the deep-learning-based works, only a few addressed the desire that the reconstructed images should be consistent with the compressed inputs. The work reported in [50] enforced this requirement as a constraint on the resulting DCT coefficients, while [19] encouraged the reconstructions to obey this requirement using a penalty as part of the optimization criterion.

We now turn to discuss two recent and relevant works that inspired this project [5, 37]. PSCGAN [37] is a recently proposed image denoiser that attempts to attain a stochastic estimator that samples from the posterior distribution. While [37] bares some similarities to our work, it differs in the task addressed (denoising vs. JPEG reconstruction), the conditional GAN loss used [17, 34], the fact that it disregards consistency, and the very different neural network architectures deployed. Nevertheless, [37] introduce a first moment loss that encourages acceptable distortion but bypasses the perception-distortion tradeoff, which we leverage in our JPEG recovery algorithm.

The work by Bahat *et al.* [5] deserves a special mention as it is the only prior work, to the best of our knowl-

edge, that addresses both the perception-distortion tradeoff and the consistency with the compressed input. [5] imposes consistency by predicting bounded residuals in the DCT domain and optimizing a GAN loss that bypasses a distortion criteria, hence avoiding the pitfall of the mentioned tradeoff. Nonetheless, our proposed approach achieves much better perceptual quality results compared to [5] while retaining perfect consistency. We give an explanation to this improvement in Section 4 based on the empirical perception-consistency tradeoff, and conduct extensive experiments to compare both methods in Section 5.

3. Method: fundamentals

We assume that a natural image X is a multivariate random variable with probability density function p_X . We denote by Y the JPEG-compressed-decompressed version of X , and from Y we provide an estimate of X , denoted by \hat{X} . We do so by attempting to sample from the posterior distribution $p_{X|Y}$. Such an estimator would be highly effective since (i) it's outputs are consistent with the measurements², i.e. $JPEG(\hat{X}) = Y$; and (ii) it attains perfect perceptual quality [6], i.e., $p_{\hat{X}} = p_X$.

3.1. JPEG

The JPEG compression algorithm [53] operates using a pre-defined quantization matrices $Q \in \mathbb{Z}^{8 \times 8}$ in the DCT domain on an image X in the YCbCr color-space, which, for simplicity, we assume is of size $8m \times 8n \times 3$. The algorithm starts by dividing the image's channels into 8×8 blocks $\{X_i\}^{m \times n \times 3}$. Each block X_i is transformed by the 2D-DCT transformation $X_i^D = \text{2D-DCT}(X_i)$ and divided element-wise by the corresponding quantization matrix Q to achieve $X_i^Q = X_i^D \oslash Q$. Finally, each entry of the block is rounded to achieve $X_i^R = \lfloor X_i^Q \rfloor$. All the blocks $\{X_i^R\}^{m \times n \times 3}$ are stored using a lossless entropy-coding algorithm alongside the matrix Q . Decompression is obtained by multiplying these rounded values by Q and applying an inverse 2D-DCT on the resulting blocks. We denote the above compression-decompression process by $JPEG_Q(\cdot)$.

Note that rounding is a lossy operation, and the matrices Q controls the amount of lost information. These matrices are a function of the quality-factor (QF) chosen by the user – an integer between 0 to 100, where 0 means maximum compression (in this paper we use the baseline matrices defined in [2]).

²In our notations throughout the paper, $Y = JPEG(X)$ implies a compression-decompression operation, and thus Y is an image of the same size as X .

3.2. Consistency

To measure the consistency of \hat{X} with Y , we define the consistency index by

$$c(p_{\hat{X}}, p_Y) = \mathbb{E}_Y \left[\mathbb{E}_{\hat{X}|Y} \left[\left\| Y - JPEG_Q(\hat{X}) \right\| | Y \right] \right], \quad (1)$$

The estimator \hat{X} is perfectly consistent with the compressed inputs Y if $Y = JPEG_Q(\hat{X})$, or equivalently, iff $c(p_{\hat{X}}, p_Y) = 0$. Interestingly, as we show next, prior JPEG restoration models that attempt to minimize the MSE loss implicitly striven (at least theoretically) to become perfectly consistent, since the MMSE estimator produces consistent restoration:

Theorem 1. *Let $\bar{X}(Y)$ be an MMSE estimator for JPEG artifact removal, i.e $\bar{X}(Y) = \mathbb{E}[X|Y]$. Then $\bar{X}(Y)$ is necessarily perfectly consistent with the compressed input Y .*

See proof in [Appendix B](#).

3.3. Perceptual quality

While there are several definitions for the notion of perceptual quality, we follow the notion developed in [6], which measures for an estimator \hat{X} the index

$$d(p_X, p_{\hat{X}}), \quad (2)$$

where $d(p, q)$ is some divergence between the two distributions, e.g., Kullback-Leibler's (KL), Jensen-Shanon's (JS) divergence, or Wasserstein distance. \hat{X} attains perfect perceptual quality when $d(p_X, p_{\hat{X}}) = 0$, i.e., $p_X = p_{\hat{X}}$.

3.4. Posterior sampling - our goal

The work in [36] ties the perceptual quality and the consistency of estimators via the following theorem:

Theorem 2. *Let $X \sim p_X$ be a random multivariate variable, and let $Y = D(X)$ be a deterministic degradation of X . If \hat{X} is an estimator such that $p_X = p_{\hat{X}}$ and $Y = D(\hat{X})$, then $p_{\hat{X}|Y} = p_{X|Y}$, i.e., \hat{X} is a necessarily posterior sampler.*

Proof. Found in [36]. □

In our case, $Y = D(X) = JPEG_Q(X)$. Recall that X cannot be uniquely recovered from Y since $JPEG_Q(\cdot)$ is a non-invertible degradation, and thus the support of $p_{X|Y}$ is not a singleton – there is a variety of possible sources that correspond to the same compressed image. Since we aim to sample from the posterior, our method must be stochastic, capable of providing many reconstructed samples given the same compressed image it. This is in opposition to most prior work, which adopts a deterministic recovery strategy.

Our solution is comprised of a conditional GAN [35] that takes Y as an input and produces high perceptual quality

outputs that are consistent with Y . Posing our task as finding a sampler from $p_{\hat{X}|Y}$, we form a loss function that has two ingredients. Just as with all GAN methods, we use an adversarial loss \mathcal{L}_{Adv} to minimize the divergence between the real and the generated data, matching p_X and $p_{\hat{X}}$ as best as possible. To promote consistency, we include another objective that penalizes any discrepancy between Y and $\text{JPEG}_Q(\hat{X})$, leading to the following final loss:

$$\min_{p_{\hat{X}|Y}} \mathcal{L}_{\text{Adv}}(p_X, p_{\hat{X}}) + \lambda \mathbb{E}_Y \mathbb{E}_{\hat{X}|Y} \left[\left\| Y - \text{JPEG}_Q(\hat{X}) \right\|_2^2 \right]. \quad (3)$$

According to [Theorem 2](#), [Eq. 3](#) admits a single optimal solution for any $\lambda > 0$: $p_{X|Y}$. This is true since any optimal solution would attain perfect perceptual quality and produce perfectly consistent restorations. Practically, however, we solve [Eq. 3](#) with parametric neural networks and with a data set of finite size, and thus our solution may not be the true posterior. Moreover, due to the nature of practical optimization, the choice of λ may also affect the obtained solution, as discussed in the following section.

4. Method: practice

4.1. Achieving the posterior

The work in [\[36\]](#) revealed that a posterior sampler is the only consistent restoration algorithm that attains perfect perceptual quality. This leads to a tradeoff between consistency and perceptual quality for deterministic estimators, as any such estimator cannot be a posterior sampler. This theoretical tradeoff **does not** affect our method, as our algorithm is stochastic. In practice, however, it is very likely that a suboptimal optimization procedure, a highly non-convex loss surface, a finite size data set, and a limited capacity architecture, would all make it extremely hard to attain the perfect solution – a sampler from the posterior. Thus, we expect to improve both the consistency index ([Eq. 1](#)) and the perceptual index ([Eq. 2](#)) up to a certain point, beyond which we shall observe an *empirical* tradeoff, where the improvement of one quality comes at the expense of the other. This empirical tradeoff is not revealed or discussed in [\[36\]](#). By changing λ in our optimization task ([Eq. 3](#)), such a tradeoff can be controlled, i.e., we can decide to attain higher perceptual quality or better consistency.

As such an empirical tradeoff has not been revealed in prior work on stochastic estimators, the balance between perceptual quality and consistency has been implicitly addressed. Bahat *et al.* [\[5\]](#) imposed a consistency requirement as a constraint using their generator architecture and not as a penalty. This can be interpreted as choosing $\lambda \rightarrow \infty$, requiring perfect consistency at the cost of lower perceptual quality. On the other hand, prior work that attempted to at-

tain the posterior without paying attention to consistency in any way, such as [\[37\]](#), controlled the tradeoff implicitly by the choice of architecture, loss, and optimizers. In fact, we argue that any attempt to sample from the posterior would most likely compromise on either consistency or perceptual quality, since attaining the true posterior is highly challenging.

4.2. Training method

We denote our estimator by $G_\theta(Z, Y)$, where G_θ is a neural network, Y is the input JPEG-compressed image, and Z is a random seed that enables a diverse set of outputs for any input image Y . Our training procedure consists of a weighted sum of several objectives. First, a non-saturating adversarial loss term [\[17\]](#), accompanied by an R_1 gradient penalty for the critic [\[34\]](#). We denote these GAN losses by $V(D_\omega, G_\theta)$ and $R_1(D_\omega)$, where D_ω is our critic. Second, we use a consistency penalty term $C(G_\theta)$, as in [Eq. 3](#):

$$C(G_\theta) = \mathbb{E}_{Y,Z} \left[\left\| Y - \text{JPEG}_Q(G_\theta(Z, Y)) \right\|_2^2 \right]. \quad (4)$$

Third, we incorporate a first-moment penalty term $FM(G_\theta)$, originally proposed in [\[37\]](#):

$$FM(G_\theta) = \mathbb{E}_{X,Y} \left[\left\| X - \mathbb{E}_Z[G_\theta(Z, Y)|Y] \right\|_2^2 \right], \quad (5)$$

which specifies that the average of many outputs $G_\theta(Z, Y)$ that refer to a fixed Y while varying Z should be close to the ideal image X . If indeed this multitude of outputs form a fair sampling from the posterior, this penalty leads exactly to the MMSE estimation $-\mathbb{E}_Z[G_\theta(Z, Y)|Y] = \mathbb{E}_X[X|Y]$. This term is a natural force that replaces the more intuitive supervised distortion penalty $\mathbb{E}_{X,Z} \left[\left\| X - G_\theta(Z, Y) \right\|_2^2 \right]$. As we have already mentioned, a distortion penalty typically hinders the perceptual quality, while the alternative in [Eq. 5](#) does not, further strengthening the overall optimization. More specifically, without using [Eq. 5](#) we observe mode-collapse during training – as we only have one X per given Y in our dataset, the generator is not incentivized to generate stochastic reconstructions, hence, it almost completely ignores Z (as explained in [\[37\]](#)).

On top of the above, in complex general content datasets such as ImageNet [\[13\]](#) we empirically find that further guidance is required in order to achieve satisfactory results. In these scenarios we include a VGG “perceptual” loss [\[24, 48\]](#), which promotes the generation of fine details in severely compressed images:

$$P(G_\theta) = \mathbb{E}_{X,Y,Z} \left[\left\| \text{VGG}_{5,4}(X) - \text{VGG}_{5,4}(G_\theta(Z, Y)) \right\|_2^2 \right]. \quad (6)$$

Where $\text{VGG}_{5,4}(\cdot)$ are the features of a trained VGG-19 network at the specified convolutional layer. Lastly, in order to

increase the variation in the estimator’s reconstructions, we introduce a new second-moment penalty:

$$SM(G_\theta) = \mathbb{E}_{X,Y} \left[\left\| (X - \bar{X}(Y))^2 - \text{Var}_Z[G_\theta(Z, Y)|Y] \right\| \right]. \quad (7)$$

This term specifies that the variance of the generated images for a given Y should be close to the sample variance using a single **ground-truth** sample and a pre-trained MMSE estimator $\bar{X}(Y)$. In [Appendix C](#) we give further rational behind this penalty.

In [Subsection 5.3](#) we present an ablation study to show the importance of the VGG loss and the second-moment penalty. Note that while we do not have a theoretical guaranty that the VGG loss does not introduce a perception-distortion tradeoff, we see in our experiments that both the perceptual quality and the consistency of the generated images improves.

All of the above forces result in the following unified minimax game:

$$\begin{aligned} \min_{\theta} \max_{\omega} & V(D_{\omega}, G_{\theta}) + \lambda_{R_1} R_1(D_{\omega}) + \lambda_C C(G_{\theta}) \\ & + \lambda_{FM} FM(G_{\theta}) + \lambda_P P(G_{\theta}) + \lambda_{SM} SM(G_{\theta}). \end{aligned} \quad (8)$$

We solve this task using a block-coordinate optimization, resulting in alternating optimization tasks:

$$\max_{\omega} V(D_{\omega}, G_{\theta}) + \lambda_{R_1} R_1(D_{\omega}), \quad (9)$$

$$\begin{aligned} \min_{\theta} & V(D_{\omega}, G_{\theta}) + \lambda_C C(G_{\theta}) + \lambda_{FM} FM(G_{\theta}) \\ & + \lambda_P P(G_{\theta}) + \lambda_{SM} SM(G_{\theta}). \end{aligned} \quad (10)$$

We should note that the consistency penalty term requires a differentiable implementation of JPEG, a concept introduced in prior work such as [\[32, 47, 51\]](#). We opt to approximate the backward pass of the rounding operation using $\nabla \lfloor X \rfloor = 1$.

4.3. Projection

As we enforce consistency through the training objective and not via the architecture, we can expect the results of our trained models to be inconsistent to some degree. Post training, though, we can still produce perfect consistency by projecting the DCT coefficients of any reconstructed image \hat{X} to the range permitted by the rounding operation that created Y . We denote the projected results by \tilde{X} . Per block, and in the DCT domain, the projection operation is defined as

$$\tilde{X}_i^Q = Y_i^Q + \max \left(\min \left(\hat{X}_i^Q - Y_i^Q, 0.5 \right), -0.5 \right). \quad (11)$$

Note that it is expected that our perceptual quality would degrade as a result of the projection operation, since it is not guaranteed to result in a natural image and also due to

empirical perception-consistency tradeoff for stochastic estimators. However, for a model that attains high perceptual quality and satisfying consistency, this degradation should be minor. We demonstrate the projection’s effect on our models in [Section 5](#).

5. Experiments

Methods: We compare our method (denoted **Ours**) with several alternative post-processing methods: (1) **QGAC** and **QGAC-GAN** [\[15\]](#), a SoTA regression method and a GAN method fine-tuned from it; (2) **SwinIR** [\[30\]](#), a SoTA regression methods trained separately for each QF; (3) **FBCNN** [\[23\]](#), a SoTA regression method; (4) **Bahat et al.** [\[5\]](#), as noted in [Section 2](#), the only other work that attempts to attain perfect consistency and perceptual quality; and (5) **Ours-MSE**, our very same architecture trained as a regression model using solely an MSE loss as a baseline. Following [Theorem 1](#), this model should produce consistent reconstructions. Moreover, we test the results obtained by our method after projection (as explained in [Subsection 4.3](#)), and denote these by **Ours-P**.

Unless mentioned otherwise, when possible we use checkpoints as published by the authors. In other cases that require training, we use commonly available hardware.³ Please refer to [Appendix A](#) for details on training and architectures.

Metrics: To measure consistency, we compute the RMSE between the compressed-decompressed versions of the original and the restored images using the same JPEG settings. Note that the value shown is per-pixel and in units of gray-levels.

To evaluate perceptual quality we compute the FID [\[21, 43\]](#) between the real uncompressed images and the restored ones. For stochastic methods (ours and Bahat *et al.*’s) we present the mean and standard deviation of 64 repeated FID evaluations, where in each we generate one restoration per compressed test image. This makes sure that our model consistently performs with high perceptual quality regardless of the seed. We also present a “ground truth” score, which is the FID between the training and the validation images (in [Subsection 5.1](#)) or between the validation and the test images (in [Subsection 5.2](#)).

5.1. Results: FFHQ

To showcase the empirical perception-consistency trade-off we use the FFHQ [\[26\]](#) thumbnails dataset, in which each image is of size 128×128 . We compress the images with QF=5 and use the same train-validation-test split as in [\[37\]](#). We compare our method to Bahat *et al.* (trained using their official implementation and following the training method described in their paper) and the base-

³The networks were trained on a single NVIDIA A6000 GPU.

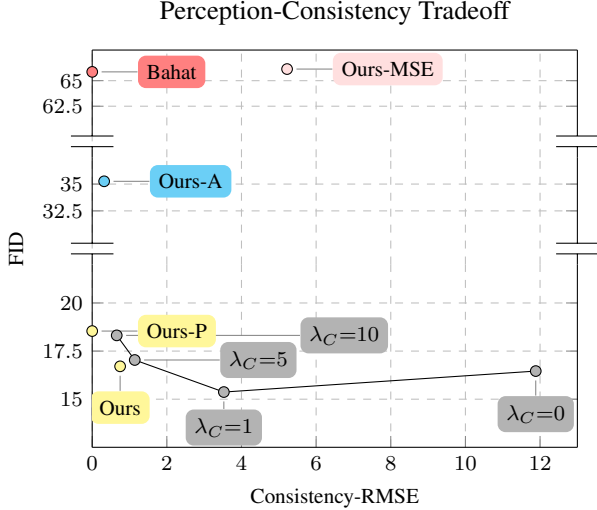


Figure 2. FID versus Consistency of different methods on FFHQ-128 with QF=5. By adjusting λ_C , our method exposes explicit control over the perception-consistency tradeoff of the trained estimator. Observe that when λ_C is increased, the consistency of the estimators improves up to a point where consistency comes at the expense of perceptual quality. By carefully tuning λ_C during training we are able to reach an estimator (Ours) with similar perceptual quality as the quality attained when using $\lambda_C=0$, but with a much improved consistency. The consistency of constrained methods (Ours-P, Bahat [5]) are practically zero – see Subsection 5.1 for more details.

line Ours-MSE method. To show the importance of the consistency regularization, we also compare the performance of our method trained by setting $\lambda_C = 0$ in the loss formulation (Eq. 10). As a proxy to an MMSE estimator, we also compare the performance of an average of 64 realizations per input from our model, denoted as **Ours-A**. The empirical perception-consistency tradeoff of the different methods on the FFHQ-128 dataset is visualized in Figure 2. Further visual results and quantitative FID, consistency and PSNR results are summarized in Appendix E (Table 4).

Consistency: As shown in Figure 2, Ours with $\lambda_C=0$ and Ours-MSE both produce unsatisfactory consistency levels, which suggests that to attain high consistency, some type of supervision is required. Just by activating our penalty term we are able to improve the consistency-RMSE by more than 11 gray-levels.

Note that both Ours-A and Ours-MSE are approximations of an MMSE estimator, and they differ only in their loss functions. While Ours-MSE attains slightly higher PSNR (see Appendix E), Ours-A attains significantly better consistency. Following Theorem 1 we know that an MMSE estimator should be perfectly consistent, so Ours-A is closer to a true MMSE estimator in that sense, even though it attains a slightly lower PSNR compared to the regression

model.

By projecting the restored images (Ours-P) we improve the consistency significantly, bringing our method to be on-par with Bahat’s. While the projection should have resulted in perfect-consistency, we get a slight deviation due to numerical approximations in the JPEG algorithm (in the color-space conversion). This is also apparent in the results of Bahat’s method that enforce the consistency as part of the architecture. In Appendix D we further investigate this phenomenon and show that it affects even the standard JPEG implementation *libjpeg* [1].

In Appendix E we present more visual and quantitative results regarding the consistency of the different methods.

Perception-consistency tradeoff: By controlling λ_C in Eq. 10 we can incentivize our generator to produce more consistent reconstructions with the compressed input, as evident in Figure 2. Starting with $\lambda_C=0$, we converge to an estimator with good perceptual quality but with lacking consistency. By choosing a small penalty, such as $\lambda_C=1$, we are able to converge to an estimator that achieves both better perceptual quality and better consistency. Yet, the consistency of the results are far from being satisfactory. Cranking the penalty coefficient up to $\lambda_C=5$ improves the consistency significantly, but the perceptual-quality deteriorate compared to $\lambda_C=0$. This trend continues as we further increase λ_C .

By carefully adjusting λ_C during training (please refer to Appendix A.2 for details) we are able to achieve a balanced result – the perceptual quality of $\lambda_C=0$ with the improved consistency of a large penalty term. This exploration demonstrates our main point in Subsection 4.1 – the chosen penalty term gives us explicit control over the empirical perception-consistency tradeoff, which lets us converge to a better local-minima, as opposed to the implicit or no control in previous methods. We can expect further improvements with more hyper-parameter tuning.

5.2. Results: general content

To showcase the performance of our method we train our model on 128×128 patches extracted from DIV2k [3] and Flickr2k [52] datasets at multiple QFs in the range [5, 50] and test all the mentioned methods on the ImageNet-ctest10k dataset, as proposed in [42]. QGAC [15] and FBCNN [23] are also trained on DIV2K and Flickr2K, SwinIR is also trained on BSDS500 [4] and WED [33], while Bahat *et al.* [5] is trained on ImageNet [13]. In Figure 3 we present the FID, consistency and PSNR of the different methods across a range of QFs from 5 to 50 on ImageNet-ctest10k, and in Figures 1 and 4 we present reconstruction examples from ImageNet-ctest10k of some of the different methods. Further visual results on ImageNet, LIVE1 [44, 45] and BSDS500 are presented in Appendix E.

Perceptual quality & consistency: Following the trends

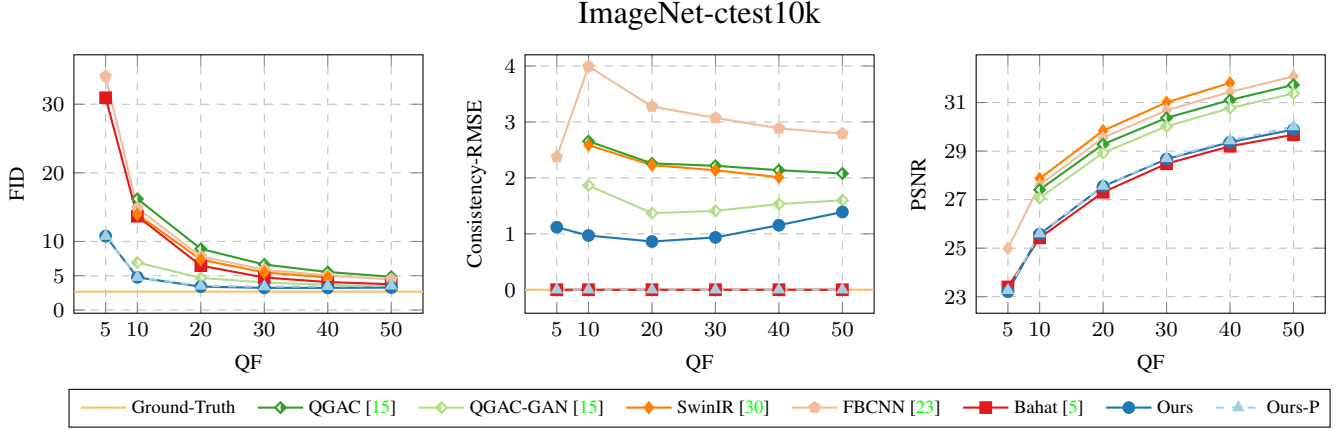


Figure 3. FID, Consistency and PSNR results of different methods on ImageNet-ctest10k across multiple QFs. By using explicit control over the perception-consistency tradeoff and by avoiding the perception-distortion tradeoff, our method provides SoTA FID results across a wide-range of QFs while being consistent with the compressed inputs. Each method is evaluated only on QFs it was trained on. The consistency of Ours-P and Bahat are practically zero – see [Subsection 5.1](#) for more details.

shown on FFHQ-128, Ours-P creates consistent reconstructions across a wide range of QFs similarly to Bahat’s, but with much improved perceptual quality.

When comparing to QGAC, SwinIR and FBCNN we see a predictable result – the regression models achieves the worst perceptual quality and QGAC-GAN model shows better results compared to Bahat’s model, but they all fail to generate consistent reconstructions. This is a direct result of the perception-distortion tradeoff [6] as also manifested in the PSNR results.

While QGAC, SwinIR and FBCNN are not trained to minimize an MSE loss, they achieve SoTA results in terms of PSNR on JPEG reconstruction, hence, they could be seen as a compelling MMSE estimator candidates. As such, they should have near-perfect consistency following [Theorem 1](#), but clearly this is not the case. Similarly to the experiments shown in [Subsection 5.1](#) on Ours-A and Ours-MSE, it seems that explicit supervision on the consistency of the reconstructed images, or enforcing it via the architecture, is necessary to achieve consistent results.

5.3. Ablation

As explained in [Subsection 4.2](#) our method is trained using several loss terms. For the FFHQ-128 data set we use a GAN loss, a consistency loss, and the first-moment penalty, while for ImageNet we also use a VGG “perceptual” loss and a second-moment penalty. In [Table 1](#) we present an ablation study of training the model with the different terms on the ImageNet data set, and here we provide a detailed explanation for each row in the table: **Baseline**: By using GAN and consistency losses alone ([Eq. 3](#)) we attain better perceptual quality (FID) than Bahat [5] but with lower variability (Per-Pixel STD); + **First-Moment**: As noted in [37], the

first-moment penalty alleviates the mode-collapse issue of conditional GANs, and indeed, it significantly improves the output variation of our model; + **VGG**: A perceptual loss further improves the perceptual quality at the cost of lower variability; + **Second-Moment**: The new penalty increases output variation without hindering perceptual quality, hence suggesting that the increased output variation is of meaningful details; + **Ours**: With increased second-moment coefficient we further improve the output variability; + **Ours-P**: By enforcing perfect consistency via projection we still attain much better perceptual quality and output variation compared to Bahat [5]; **Perceptual Baseline**: Without explicitly requiring consistency we achieve similar FID to our method but without consistency with the measurements.

6. Summary

In this work we approach the JPEG decompression task from an uncommon direction – generate visually pleasing and consistent reconstructions by leveraging recent advancements in image restoration, such as the perception-distortion and the perception-consistency theoretical trade-offs. Using these tools we surpass prior work and provide decompressed JPEG images of tunable consistency and high perceptual quality.

The proposed solution is based on a stochastic conditional GAN with carefully tailored loss function that promotes detailed and vivid results, consistency to the measurements, proximity to the training data without sacrificing quality, and a spread of the randomized results. Our future work will focus on better diversifying the obtained solutions and on a quest for a tractable computational method for evaluating the proximity of the obtained model to the ideal posterior sampler.

Table 1. Ablation study on the ImageNet-ctest10k dataset at QF=10. The VGG loss is crucial for good-perceptual quality without hurting consistency. The second-moment penalty increases per-pixel standard deviation (a value between 0 and 1) without hurting either perceptual quality nor consistency. The consistency of constrained methods, marked by *, are practically zero – see [Subsection 5.1](#) for more details.

Method	GAN	C	FM	P	SM	FID (\downarrow)	Consistency (\downarrow)	PSNR (\uparrow)	Per-Pixel STD (\uparrow)
Perceptual Baseline	✓		✓	✓		4.54	6.9724	26.0612	0.0111
Baseline	✓	✓				9.12	0.9650	25.5061	0.0021
+ First-Moment	✓	✓	✓			8.83	0.9694	25.5767	0.0064
+ VGG	✓	✓	✓	✓		4.71	0.9411	25.8906	0.0026
+ Second-Moment	✓	✓	✓	✓	✓ (10^1)	4.70	0.9452	25.8539	0.0063
Ours	✓	✓	✓	✓	✓ (10^3)	4.76	0.9696	25.5758	0.0194
Ours-P	✓	✓	✓	✓	✓ (10^3)	4.78	$\approx 0^*$	25.6008	0.0188
Bahat						13.71	$\approx 0^*$	25.4243	0.0040

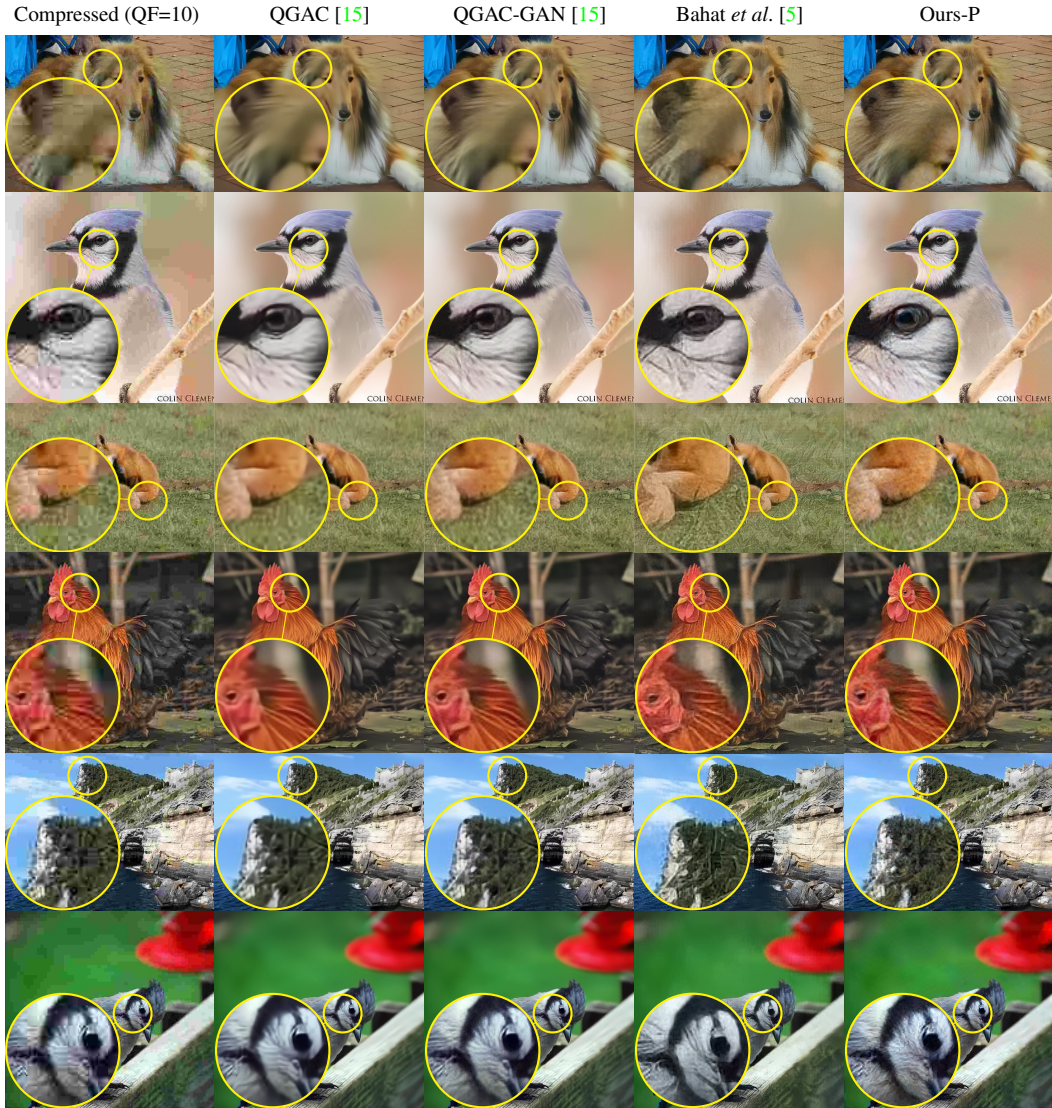


Figure 4. Reconstruction examples on ImageNet-ctest10k at QF=10; Bahat is the only other method that generates consistent reconstructions and aims for high perceptual quality; QGAC-GAN’s reconstructions are sharp but not consistent, hence could not have created the input image; Our method generates finer details while still being consistent, as seen in [Figure 3](#); Zoom-in on interesting regions are shown.

Acknowledgement This research was partially supported by the Israel Science Foundation (ISF) under Grant 335/18 and the Council For Higher Education - Planning & Budgeting Committee.

References

- [1] Independent JPEG Group. <http://www.ijg.org/>. 6, 13
- [2] Recommendation T.81: Information technology - Digital compression and coding of continuous-tone still images - Requirements and guidelines, Sept. 1992. 3
- [3] Eirikur Agustsson and Radu Timofte. NTIRE 2017 Challenge on Single Image Super-Resolution: Dataset and Study. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 126–135, 2017. 6
- [4] Pablo Arbeláez, Michael Maire, Charless Fowlkes, and Jitendra Malik. Contour Detection and Hierarchical Image Segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(5):898–916, May 2011. 6, 15
- [5] Yuval Bahat and Tomer Michaeli. What’s in the Image? Explorable Decoding of Compressed Images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2908–2917, 2021. 1, 2, 3, 4, 5, 6, 7, 8, 13, 16
- [6] Yochai Blau and Tomer Michaeli. The Perception-Distortion Tradeoff. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6228–6237, June 2018. 2, 3, 7
- [7] G. Bradski. The OpenCV library. *Dr. Dobb’s Journal of Software Tools*, 2000. 13
- [8] Lukas Cavigelli, Pascal Hager, and Luca Benini. CAS-CNN: A deep convolutional neural network for image compression artifact suppression. In *2017 International Joint Conference on Neural Networks (IJCNN)*, pages 752–759, May 2017. 2
- [9] Lucy Chai, Michaël Gharbi, Eli Shechtman, Phillip Isola, and Richard Zhang. Any-Resolution Training for High-Resolution Image Synthesis. In Shai Avidan, Gabriel Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner, editors, *Computer Vision – ECCV 2022*, volume 13676, pages 170–188. Springer Nature Switzerland, Cham, 2022. 12
- [10] Huibin Chang, Michael K. Ng, and Tiejong Zeng. Reducing Artifacts in JPEG Decompression Via a Learned Dictionary. *IEEE Transactions on Signal Processing*, 62(3):718–728, Feb. 2014. 2
- [11] T. Chen, H.R. Wu, and B. Qiu. Adaptive postfiltering of transform coefficients for the reduction of blocking artifacts. *IEEE Transactions on Circuits and Systems for Video Technology*, 11(5):594–602, May 2001. 2
- [12] Alex Clark. Pillow (PIL fork), 2015. 13
- [13] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, June 2009. 4, 6
- [14] Chao Dong, Yubin Deng, Chen Change Loy, and Xiaoou Tang. Compression Artifacts Reduction by a Deep Convolutional Network. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 576–584, Dec. 2015. 2
- [15] Max Ehrlich, Larry Davis, Ser-Nam Lim, and Abhinav Shrivastava. Quantization Guided JPEG Artifact Correction. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision – ECCV 2020*, Lecture Notes in Computer Science, pages 293–309, Cham, 2020. Springer International Publishing. 2, 5, 6, 7, 8, 14
- [16] Leonardo Galteri, Lorenzo Seidenari, Marco Bertini, and Alberto Del Bimbo. Deep Universal Generative Adversarial Compression Artifact Removal. *IEEE Transactions on Multimedia*, 21(8):2131–2145, Aug. 2019. 2
- [17] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative Adversarial Nets. In *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014. 2, 4
- [18] Jun Guo and Hongyang Chao. Building Dual-Domain Representations for Compression Artifacts Reduction. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Computer Vision – ECCV 2016*, Lecture Notes in Computer Science, pages 628–644, Cham, 2016. Springer International Publishing. 2
- [19] Jun Guo and Hongyang Chao. One-To-Many Network for Visually Pleasing Compression Artifacts Reduction. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4867–4876, July 2017. 2
- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. 12
- [21] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. 5
- [22] Sung-Wai Hong, Yuk-Hee Chan, and Wan-Chi Siu. A practical real-time post-processing technique for block effect elimination. In *Proceedings of 3rd IEEE International Conference on Image Processing*, volume 2, pages 21–24 vol.2, Sept. 1996. 2
- [23] Jiaxi Jiang, Kai Zhang, and Radu Timofte. Towards Flexible Blind JPEG Artifacts Removal. Sept. 2021. 5, 6, 7
- [24] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual Losses for Real-Time Style Transfer and Super-Resolution. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Computer Vision – ECCV 2016*, Lecture Notes in Computer Science, pages 694–711, Cham, 2016. Springer International Publishing. 4
- [25] Cheolkon Jung, Licheng Jiao, Hongtao Qi, and Tian Sun. Image deblocking via sparse representation. *Signal Processing: Image Communication*, 27(6):663–677, July 2012. 2
- [26] Tero Karras, Samuli Laine, and Timo Aila. A Style-Based Generator Architecture for Generative Adversarial Networks. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4396–4405, June 2019. 5

- [27] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 12
- [28] A. Kundu. Enhancement of JPEG coded images by adaptive spatial filtering. In *Proceedings., International Conference on Image Processing*, volume 1, pages 187–190 vol.1, Oct. 1995. 2
- [29] Jin Li and C.-C.J. Kuo. Coding artifact removal with multi-scale postprocessing. In *Proceedings of International Conference on Image Processing*, volume 1, pages 45–48 vol.1, Oct. 1997. 2
- [30] Jingyun Liang, Jiezhong Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. SwinIR: Image Restoration Using Swin Transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1833–1844, 2021. 5, 7
- [31] Jiebo Luo, Chang Wen Chen, K.J. Parker, and T.S. Huang. Artifact reduction in low bit rate DCT-based image compression. *IEEE Transactions on Image Processing*, 5(9):1363–1368, Sept. 1996. 2
- [32] Xiyang Luo, Hossein Talebi, Feng Yang, Michael Elad, and Peyman Milanfar. The Rate-Distortion-Accuracy Tradeoff: JPEG Case Study. *arXiv:2008.00605 [cs, eess]*, Aug. 2020. 5
- [33] Kede Ma, Zhengfang Duanmu, Qingbo Wu, Zhou Wang, Hongwei Yong, Hongliang Li, and Lei Zhang. Waterloo Exploration Database: New Challenges for Image Quality Assessment Models. *IEEE Transactions on Image Processing*, 26(2):1004–1016, Feb. 2017. 6
- [34] Lars Mescheder, Andreas Geiger, and Sebastian Nowozin. Which Training Methods for GANs do actually Converge? In *Proceedings of the 35th International Conference on Machine Learning*, pages 3481–3490. PMLR, July 2018. 2, 4
- [35] Mehdi Mirza and Simon Osindero. Conditional Generative Adversarial Nets. *arXiv:1411.1784 [cs, stat]*, Nov. 2014. 3
- [36] Guy Ohayon, Theo Adrai, Michael Elad, and Tomer Michaeli. Reasons for the Superiority of Stochastic Estimators over Deterministic Ones: Robustness, Consistency and Perceptual Quality. *arXiv:2211.08944 [cs, eess]*, Nov. 2022. 2, 3, 4
- [37] Guy Ohayon, Theo Adrai, Gregory Vaksman, Michael Elad, and Peyman Milanfar. High Perceptual Quality Image Denoising With a Posterior Sampling CGAN. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops*, pages 1805–1813, 2021. 2, 4, 5, 7
- [38] T.P. O’Rourke and R.L. Stevenson. Improved image de-compression for reduced transform coding artifacts. *IEEE Transactions on Circuits and Systems for Video Technology*, 5(6):490–499, Dec. 1995. 2
- [39] T. Ozcelik, J.C. Brailean, and A.K. Katsaggelos. Image and video compression algorithms based on recovery techniques using mean field annealing. *Proceedings of the IEEE*, 83(2):304–316, Feb. 1995. 2
- [40] Ethan Perez, Florian Strub, Harm de Vries, Vincent Dumoulin, and Aaron Courville. FiLM: Visual Reasoning with a General Conditioning Layer. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1), Apr. 2018. 12, 13
- [41] H. Reeve and Jae Lim. Reduction of blocking effect in image coding. In *ICASSP ’83. IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume 8, pages 1212–1215, Apr. 1983. 2
- [42] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-Image Diffusion Models. In *ACM SIGGRAPH 2022 Conference Proceedings*, SIGGRAPH ’22, pages 1–10, New York, NY, USA, July 2022. Association for Computing Machinery. 6
- [43] Maximilian Seitzer. Pytorch-fid: FID score for PyTorch, Aug. 2020. Version 0.2.1. 5
- [44] Hamid R Sheikh, Muhammad F Sabir, and Alan C Bovik. A statistical evaluation of recent full reference image quality assessment algorithms. *IEEE Transactions on image processing*, 15(11):3440–3451, 2006. 6, 15
- [45] HR Sheikh, Z Wang, L Cormack, and AC Bovik. LIVE image quality assessment database. <http://live.ece.utexas.edu/research/quality>. 6, 15
- [46] Mei-Yin Shen and C. C. Jay Kuo. Review of Postprocessing Techniques for Compression Artifact Removal. *Journal of Visual Communication and Image Representation*, 9(1):2–14, Mar. 1998. 2
- [47] Richard Shin and Dawn Song. JPEG-resistant Adversarial Images. In *NIPS 2017 Workshop on Machine Learning and Computer Security*, volume 1, page 8, 2017. 5
- [48] Karen Simonyan and Andrew Zisserman. Very Deep Convolutional Networks for Large-Scale Image Recognition, Apr. 2015. 4
- [49] R.L. Stevenson. Reduction of coding artifacts in transform image coding. In *1993 IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume 5, pages 401–404 vol.5, Apr. 1993. 2
- [50] Mengdi Sun, Xiaohai He, Shuhua Xiong, Chao Ren, and Xinglong Li. Reduction of JPEG compression artifacts based on DCT coefficients prediction. *Neurocomputing*, 384:335–345, Apr. 2020. 2
- [51] Hossein Talebi, Damien Kelly, Xiyang Luo, Ignacio Garcia Dorado, Feng Yang, Peyman Milanfar, and Michael Elad. Better Compression With Deep Pre-Editing. *IEEE Transactions on Image Processing*, 30:6673–6685, 2021. 5
- [52] Radu Timofte, Eirikur Agustsson, Luc Van Gool, Ming-Hsuan Yang, and Lei Zhang. NTIRE 2017 Challenge on Single Image Super-Resolution: Methods and Results. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 114–125, 2017. 6
- [53] G.K. Wallace. The JPEG still picture compression standard. *IEEE Transactions on Consumer Electronics*, 38(1):xviii–xxxiv, Feb. 1992. 1, 3
- [54] Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. Real-ESRGAN: Training Real-World Blind Super-Resolution With Pure Synthetic Data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1905–1914, 2021. 12, 13
- [55] Xintao Wang, Liangbin Xie, Ke Yu, Kelvin C.K. Chan, Chen Change Loy, and Chao Dong. BasicSR: Open source image and video restoration toolbox, 2022. 12

- [56] Zhangyang Wang, Ding Liu, Shiyu Chang, Qing Ling, Yingzhen Yang, and Thomas S. Huang. D3: Deep Dual-Domain Based Fast Restoration of JPEG-Compressed Images. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2764–2772, June 2016. [2](#)
- [57] S.-W. Wu and A. Gersho. Improved decoder for transform coding with application to the JPEG baseline system. *IEEE Transactions on Communications*, 40(2):251–254, Feb. 1992. [2](#)
- [58] Xiaoshuai Zhang, Wenhan Yang, Yueyu Hu, and Jiaying Liu. Dmccnn: Dual-Domain Multi-Scale Convolutional Neural Network for Compression Artifacts Removal. In *2018 25th IEEE International Conference on Image Processing (ICIP)*, pages 390–394, Oct. 2018. [2](#)

A. Implementation details

A.1. Architecture

Our GAN architecture is composed of an RRDB based generator and a ResNet critic. The generator largely follows the architecture proposed in Real-ESRGAN [54, 55], where we use a pixel-unshuffle block to rearrange a $m \times n \times 3$ input into a $m/2 \times n/2 \times 12$ input to reduce the computational complexity of the following 23 RRDB blocks, followed by an upscale of the result to restore the original input shape. One major difference from [54] is the injection of noise channels along the network – Each “Noise Injection” block concatenates a new channel filled with Gaussian noise, which is used by the network to generate stochastic details (this is Z , mentioned in the loss formulation). In cases supporting multiple QFs we implement a FiLM [40] block that modulates the results of an RRDB or upsampling block using a learned affine transformation, conditioned on the quantization table that is embedded in the JPEG file. We present the generator architecture in Figure 5. The critic is a plain ResNet34 [20].

A.2. Training

We train our models using the Adam [27] optimizer with batch-size of 32 to alternately update the generator and the critic networks, using Eq. 10 and Eq. 9. We use exponential moving average on the generator weights with decay factor of 0.999.

FFHQ-128: The learning rate starts at 1×10^{-4} and annealed to 1×10^{-7} after 400,000 steps using cosine annealing. We scale $V(D_\omega, G_\theta)$ in Eq. 10 by 1×10^{-3} and set $\lambda_{R_1} = 1$, $\lambda_{FM} = 1$ and λ_C according to the version reported in Figure 2. For the versions denoted as **Ours** and **Ours-P** we use cosine annealing of λ_C from 1×10^{-1} to 1×10^1 .

General-Content: The learning rate starts at 1×10^{-4} and annealed to 1×10^{-6} after 400,000 steps using cosine annealing. The weights are initialized from a model trained for 50,000 steps as a regression model using an MSE loss, as this was found to be important for stabilizing the training. We scale $V(D_\omega, G_\theta)$ in Eq. 10 by 1×10^{-3} and set $\lambda_{R_1} = 1$, $\lambda_{FM} = 1$, $\lambda_P = 1 \times 10^{-2}$, $\lambda_{SM} = 1 \times 10^3$ and use cosine annealing of λ_C from 1×10^0 to 1×10^3 .

To estimate the mean and variance of generated images for use in $FM(G_\theta)$ and $SM(G_\theta)$ we generate 16 different restorations (using different Z s) for each of the first 8 images in a batch. In order to save train time we preform this every 8 iterations as we found negligible performance differences compared to performing this each iteration.

As the reference MMSE estimator for $SM(G_\theta)$ we use a regression model with the same architecture, trained for 650,000 steps using a simple MSE loss.

As training data, we extract square patches with random

scale (between 128×128 and the image resolution) from the training set at random position and rescale them to 128×128 pixels. This is inspired by [9] and it exposes our GAN to more diverse set of patches.

B. MMSE estimator is consistent

Theorem 1. *Let $\bar{X}(Y)$ be an MMSE estimator for JPEG artifact removal, i.e. $\bar{X}(Y) = \mathbb{E}[X|Y]$. Then $\bar{X}(Y)$ is necessarily perfectly consistent with the compressed input Y .*

Proof. Denote by D the matrix that performs block-wise 2D-DCT and elementwise division by the matrix Q . Then $\bar{X}(Y)$ is consistent with Y iff $\|D\bar{X}(Y) - DY\|_\infty \leq \frac{1}{2}$. And indeed,

$$\begin{aligned} \|D\bar{X}(Y) - DY\|_\infty &= \|D\mathbb{E}[X|Y] - DY\|_\infty \\ &= \|\mathbb{E}[DX - DY|Y]\|_\infty \\ &\leq \mathbb{E}[\|DX - DY\|_\infty | Y] \\ &\leq \mathbb{E}\left[\frac{1}{2}|Y\right] = \frac{1}{2}, \end{aligned}$$

where we used the triangle inequality and the fact that at any point where $p(X|Y) > 0$, the maximal difference in the DCT domain, before rounding, is $\frac{1}{2}$. \square

C. Second-moment penalty

The **per-pixel** conditional variance $\sigma_{X|Y}^2$ of a random-variable X can be estimated from samples $\{x_i\}_{i=1}^n$ sampled from $p_{X|Y}$:

$$\sigma_{X|Y}^2 \approx \frac{1}{n} \sum_{i=1}^n (x_i - \mu_{X|Y})^2, \quad (12)$$

where $\mu_{X|Y}$ is the mean of X conditioned on Y . In practice we do not have access to this mean, hence we can approximate it either using an MMSE estimator $\bar{X}(Y)$ (that at optimality becomes the conditional mean) or using the sample mean:

$$\tilde{\mu}_{X|Y} = \frac{1}{n} \sum_{i=1}^n x_i. \quad (13)$$

In our case, we have two different variables for which we would like to compare their variances – the ground-truth images X and our reconstructed images \hat{X} , both conditioned on the compressed images Y .

Recall that for a fixed Y we have only a single ground-truth sample X and thus we cannot compute a useful sample conditional mean, hence we opt to use a pre-trained regression model as our MMSE estimator $\bar{X}(Y)$:

$$\tilde{\sigma}_{X|Y}^2 = (x - \bar{X}(Y))^2. \quad (14)$$

Intuitively, $\bar{X}(Y)$ averages all possible values of each pixel in the reconstructed image based on the probability of the

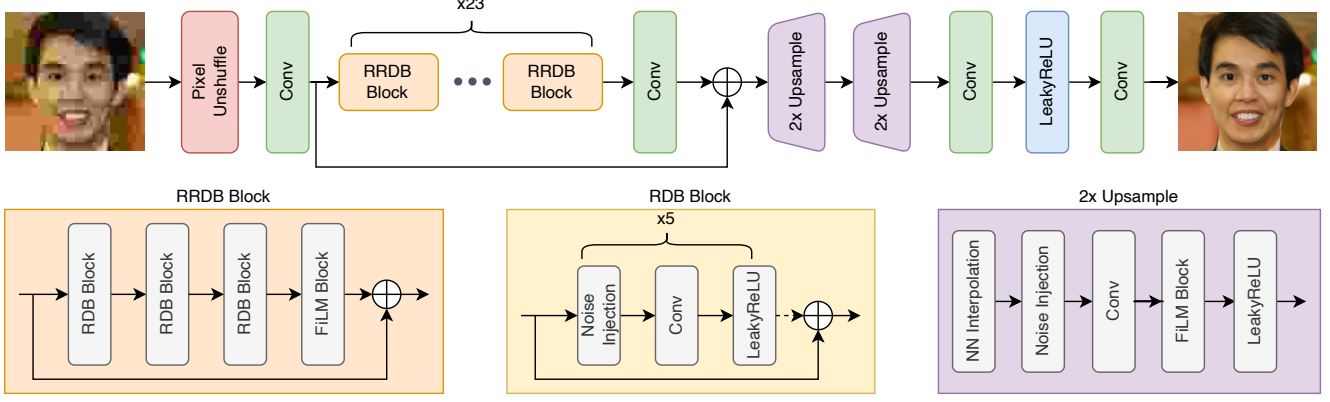


Figure 5. Our generator’s architecture is based on [54], where Residual-in-Residual Dense Blocks (RRDBs) are used to restore a degraded input. Due to the fact that we work on full resolution inputs, we trade spatial resolution with depth using a pixel-unshuffle block that rearrange the input from $m \times n \times 3$ tensor into a $m/2 \times n/2 \times 12$ tensor and in the end we upsample the resulted tensor. To achieve stochasticity, we inject a noise channel to the features at the beginning of every RDB block and in the upsampling blocks. To be agnostic to the QF of the input image, we condition the RRDB and upsampling blocks on the quantization table used to create the image using a FiLM [40] block.

value. Hence, pixels with similar values in $\bar{X}(Y)$ and in x are pixels with small ambiguity – all of the possible reconstructions agree on the pixel’s value, and thus we can expect small variance at such locations. On the other hand, pixels with large discrepancy between $\bar{X}(Y)$ and x are not necessarily an indication for large variance as the value in x might be rare and thus far from the mean. To better analyze the latter case we need further assumption on the conditional probability $X|Y$.

As $p_{\hat{X}|Y}$ changes during training, we cannot use a pre-trained regression model that was trained on $p_{X|Y}$, but we can generate as much samples as needed, hence we can compute a sample conditional mean that approximates well the true mean:

$$\tilde{\sigma}_{\hat{X}|Y}^2 = \frac{1}{n} \sum_{i=1}^n (\hat{x}_i - \tilde{\mu}_{\hat{X}|Y})^2, \quad (15)$$

To create the penalty mentioned in Subsection 4.2 we just compute the distance between those two, per-pixel, variance approximations over a batch of realizations of Y :

$$SM(\hat{X}) = \lambda_{SM} \mathbb{E}_{X, \hat{X}, Y} \left[\left\| \tilde{\sigma}_{\hat{X}|Y}^2 - \tilde{\sigma}_{X|Y}^2 \right\| \right]. \quad (16)$$

In Figure 6 we present a visual illustration of the variance estimator $\tilde{\sigma}_{\hat{X}|Y}^2$ using QGAC as our MMSE estimator alongside per-pixel sample variance of Ours’ and Bahat *et al.*’s methods.

D. JPEG numerical errors

As mentioned in Subsection 5.1 and Table 4, numerical approximations in the JPEG algorithm result in near,

but not perfect, consistent results. To showcase that this phenomena is not unique to the differentiable JPEG implementations used by us and [5] we test *libjpeg-9d* [1], a standard JPEG implementation used in packages such as OpenCV [7] and PIL [12]. We compress the FFHQ test set with QF=100 which corresponds to no quantization, at block size of 1 which means the DCT values equal the color values, and without chroma sub-sampling, meaning we should expect a perfect reconstruction of the input images theoretically. Table 2 present the RMSE between the ground-truth images and the compressed images, using different configurations of the encoder-decoder. As it can be clearly seen, as long as we do not skip the YCbCr color-space conversion, true lossless compression is not achieved due to numerical approximations. Table 3 presents the actual Consistency RMSE results of the projected methods. We can see that while all the methods are not perfectly consistent, they are extremely quite to zero and perform better than *libjpeg-9d* when operating at the YCbCr color space.

To reproduce the *libjpeg-9d* images, use the following snippet:

```
1 # 2D-DCT=float, Color-Space=YCbCr, Block-Size=1
2 cjpeg -block 1 -quality 100 -sample 1x1,1x1,1x1 -
   dct float <input_image> | djpeg -dct float >
   <output_image>
3
4 # 2D-DCT=float, Color-Space=RGB, Block-Size=1
5 cjpeg -block 1 -quality 100 -sample 1x1,1x1,1x1 -
   dct float -rgb <input_image> | djpeg -dct
   float > <output_image>
```

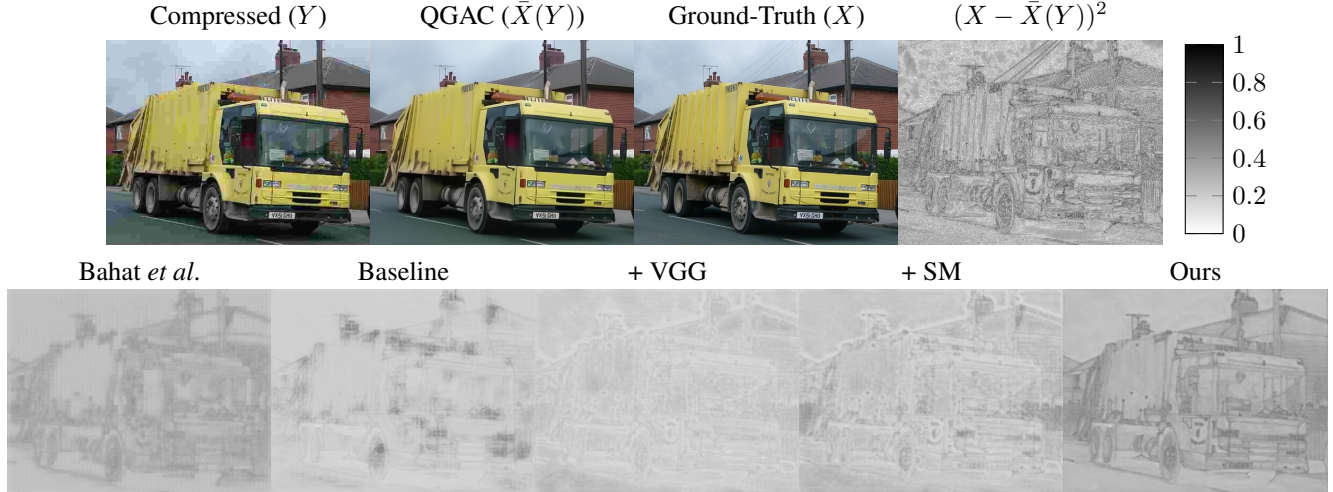



Figure 6. We approximate the conditional variance of clean images on a compressed input (Y) using a single ground-truth image (X) and an MMSE estimator ($\bar{X}(Y)$). Here we use QGAC [15] as our MMSE estimator on a JPEG compressed image with QF=10. On the bottom row we present per-pixel variance maps calculated on 64 samples of different methods (\hat{X}). For visualization purposes we use the 8th root of the variance and add a color bar (white and black correspond to low and high variance, respectively). The variance estimation $(X - \bar{X}(Y))^2$ indicates that we should expect more variance in regions with sharp transitions and this is indeed what we see in the variance maps of the different methods. Our method produce more variance compared to Bahat *et al.*'s method in regions with sharp transition and less variance in smooth regions. This is also supported quantitatively in Table 1 where we see larger average per-pixel standard deviation while achieving better FID.

Table 2. RMSE between clean FFHQ test images and compressed-decompressed FFHQ test images using the *libjpeg-9d* library. We use no chroma sub-sampling and quantization table of ones, hence, the process should be invertible. In practice, we see deviations due to numerical approximations.

2D-DCT	Color-Space	Block-Size	RMSE
float	YCbCr	1	0.5465
float	RGB	1	0

Table 3. Consistency results of different constrained methods on the FFHQ test dataset at QF=5. The inconsistencies stem from numerical approximations.

Method	Consistency
Bahat	0.0867
$\lambda_C=0$ -P	0.9466
Ours-P	0.1664
Ground Truth (theoretical)	0

E. More results

E.1. FFHQ

In Table 4 we present quantitative results of the different methods on the FFHQ test set at QF=5.

In Figure 7 we visually demonstrate the perceptual quality and the stochastic variation of our method by presenting several realizations of a given input and comparing them

to the other methods. As expected, the regression model generates overly-smoothed results and is unable to recover fine details such as hairs and wrinkles. Bahat’s method successfully recovers some fine details but suffers from severe color and grid-like artifacts. We find that the training of this method is highly unstable, and we hypothesize that this is partly due to the overly constrained optimization with perfect consistency. The results denoted Ours are highly visually appealing – fine details such as hair and wrinkles are generated in a reasonable manner.

In Figure 8 we present a couple of compressed images from the test set of FFHQ and the corresponding recompressed restoration from our method with and without consistency regularization and with projection. This showcases the effectiveness of our consistency regularization in improving the consistency of the reconstructed images without deteriorating their perceptual quality.

In Figure 9 we present the stochastic nature of the different methods (except the deterministic regression models). We expect to see different plausible details generated for the same compressed input image Y given different noise injection Z . Indeed, we see that our methods generate slight variations in the expression, in the beard and hair structure, in the background details and in the skin colors. Those variations are also indicated by the per-pixel standard deviation map we present for each method, where darker values represent more varying pixels in the restored images. While Bahat’s method also produces stochastic results, it can be

clearly seen that most of the variation comes from color artifacts.

In Figure 10 we present more results of the different methods on the test set of FFHQ.

Table 4. Quantitative results of different methods on the FFHQ test set at QF=5. The consistency of constrained methods, marked by *, are practically zero – please refer to Subsection 5.1 and Appendix D for more details.

Method	FID (\downarrow)	Consistency (\downarrow)	PSNR (\uparrow)
Regression	66.14 ± 0.00	5.2217	25.6579
Ours-A	35.26 ± 0.00	0.3203	25.4529
Bahat	65.86 ± 0.28	$\approx 0^*$	22.6807
$\lambda_C=0$	16.46 ± 0.16	11.8848	23.3457
$\lambda_C=0$ -P	20.60 ± 0.16	$\approx 0^*$	23.4896
Ours	16.70 ± 0.14	0.7481	23.7595
Ours-P	18.54 ± 0.14	$\approx 0^*$	23.7767
Ground Truth	10.28 ± 0.00	0	∞

E.2. ImageNet

In Table 5 we present quantitative results of the different methods on ImageNet-cstest10k at QF=10.

In Figure 11 and Figure 12 we present more results of the different methods on ImageNet-cstest10k. Note that QGAC and QGAC-GAN are not trained on QFs lower than 10, hence we do not show their results on QF=5 for fair comparison. The visual results further corroborate the quantitative results shown in Figure 3 – Our method provides the best perceptual results, creating more fine details and less artifacts and projecting our results does not deteriorate their perceptual quality. Note that while QGAC-GAN provide better visual results compared to Bahat’s, they are not consistent with the compressed inputs. This means that those are not valid reconstructions in the sense that they could not have created the compressed images.

E.3. LIVE1 & BSDS500

In Figure 13 and Figure 14 we present visual results of different methods on LIVE1 [44, 45] and BSDS500 [4] datasets.

Such small datasets (29 and 500 images, respectively) cannot be used for reliable deep-features-based, ensemble perceptual quality assessments (FID, KID, IS, etc.). From the official FID implementation⁴ “IMPORTANT: The number of samples [...] should be greater than the dimension of the coding layer, here 2048 [...]”. Hence, we do not include quantitative results for this datasets.

Table 5. Quantitative results of different methods on ImageNet-cstest10k at QF=10. The consistency of constrained methods, marked by *, are practically zero – please refer to subsection 5.1 and Appendix D for more details.

Method	FID (\downarrow)	Consistency (\downarrow)	PSNR (\uparrow)
SwinIR	13.93 ± 0.00	2.5858	27.8662
FBCNN	14.85 ± 0.00	3.9929	27.6000
QGAC	16.20 ± 0.00	2.6551	27.4091
QGAC-GAN	6.93 ± 0.00	1.8640	27.0681
Bahat	13.71 ± 0.03	0.3598*	25.4243
Ours	4.76 ± 0.01	0.9696	25.5758
Ours-P	4.78 ± 0.01	0.6411*	25.6008
Ground Truth	2.67 ± 0.00	0	∞

⁴<https://github.com/bioinf-jku/TTUR>



Figure 7. Zoom-in on the results of different recovery methods. Left column: The decompression of a single image from the FFHQ data set compressed using QF=5. Notice the smoothed result of Ours-MSE and the artifacts in Bahat’s solution, while our method produces sharp and realistic results. Right column: Four realizations from our method that further show the stochastic nature of the results. Note the different hair patterns and ear shapes.

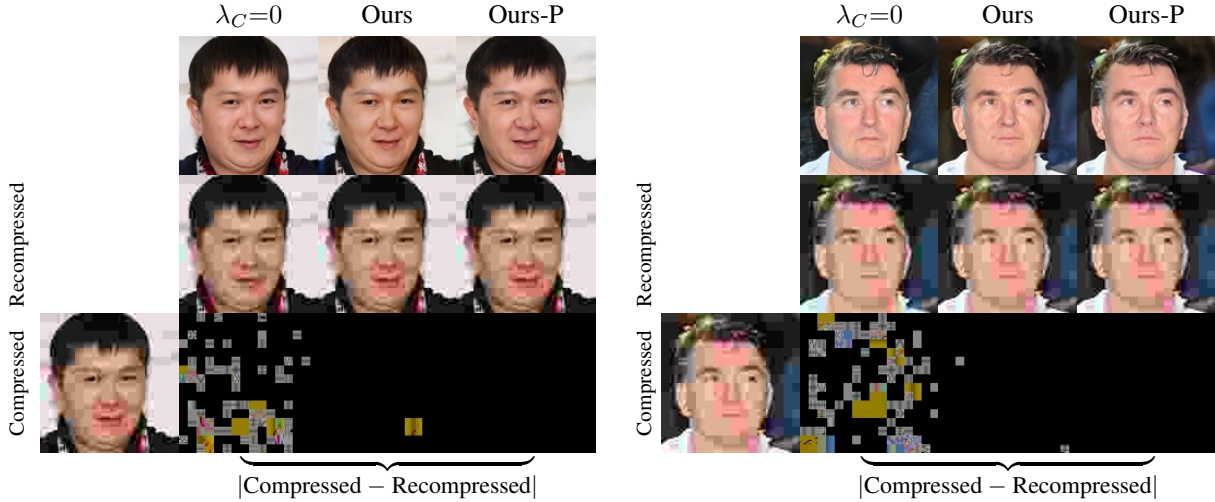


Figure 8. The difference between a compressed input and the recompressed outputs of our method with and without explicit consistency penalty and with projection. Values has been rescaled by taking the 4th root for visualization purposes. By adjusting λ_C we are able to produce reconstruction with near-perfect consistency (Ours). This allows us to project the results to achieve perfect consistency with minimal impact on the perceptual quality (Ours-P). Quantitative results can be seen in Table 4.

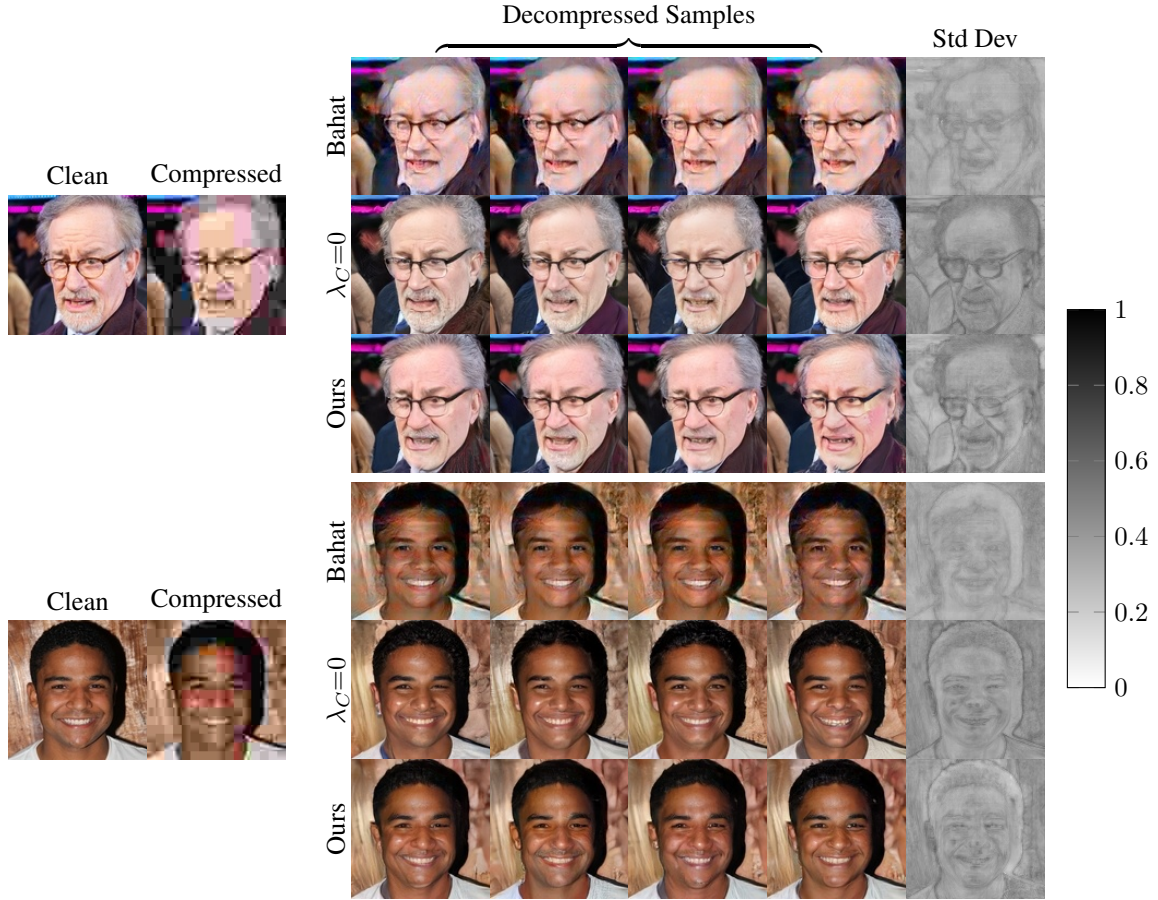


Figure 9. Stochastic variation of decompressed images using Bahat’s and our methods. To the left we present two clean images and their corresponding compressed JPEG version using QF=5. To the right we present 4 realizations using each method, along with per-pixel standard deviation map calculated on 32 samples. For visualization purposes we use the 4th root of the standard deviation and add a color bar (white and black correspond to low and high standard deviations, respectively). All decompressed images were obtained using the default noise injection scheme ($z \sim \mathcal{U}(-1, 1)$ for Bahat’s method and $z \sim \mathcal{N}(0, I)$ for the others).



Figure 10. Decompression results using different methods on FFHQ images compressed using JPEG with QF=5. For stochastic methods (all except for Regression), the default noise injection scheme ($z \sim \mathcal{U}(-1, 1)$ for Bahat’s method and $z \sim \mathcal{N}(0, I)$ for the others) was used during inference.

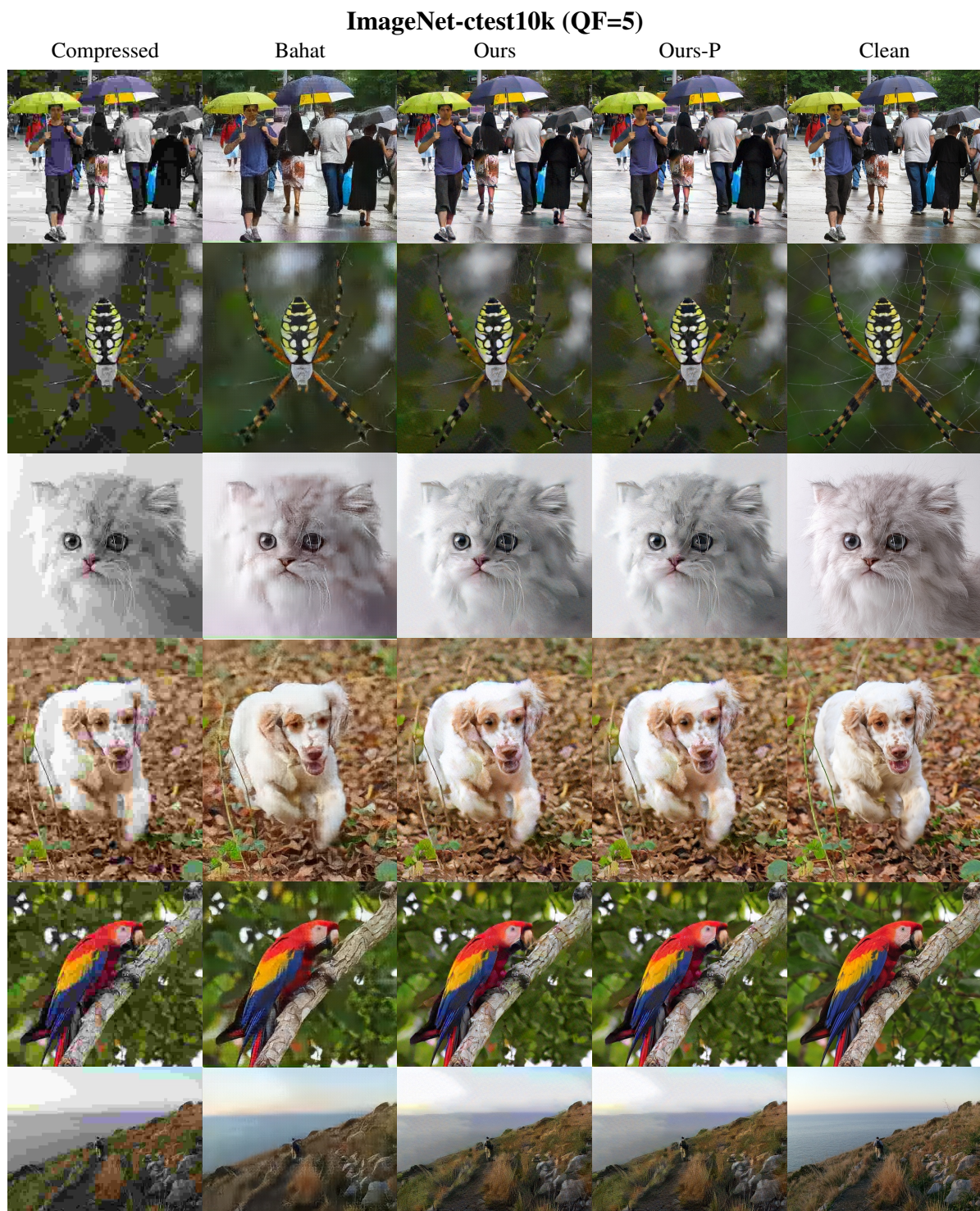


Figure 11. Decompression results using our method and Bahat *et al.*'s on ImageNet JPEG compressed images with QF=5.

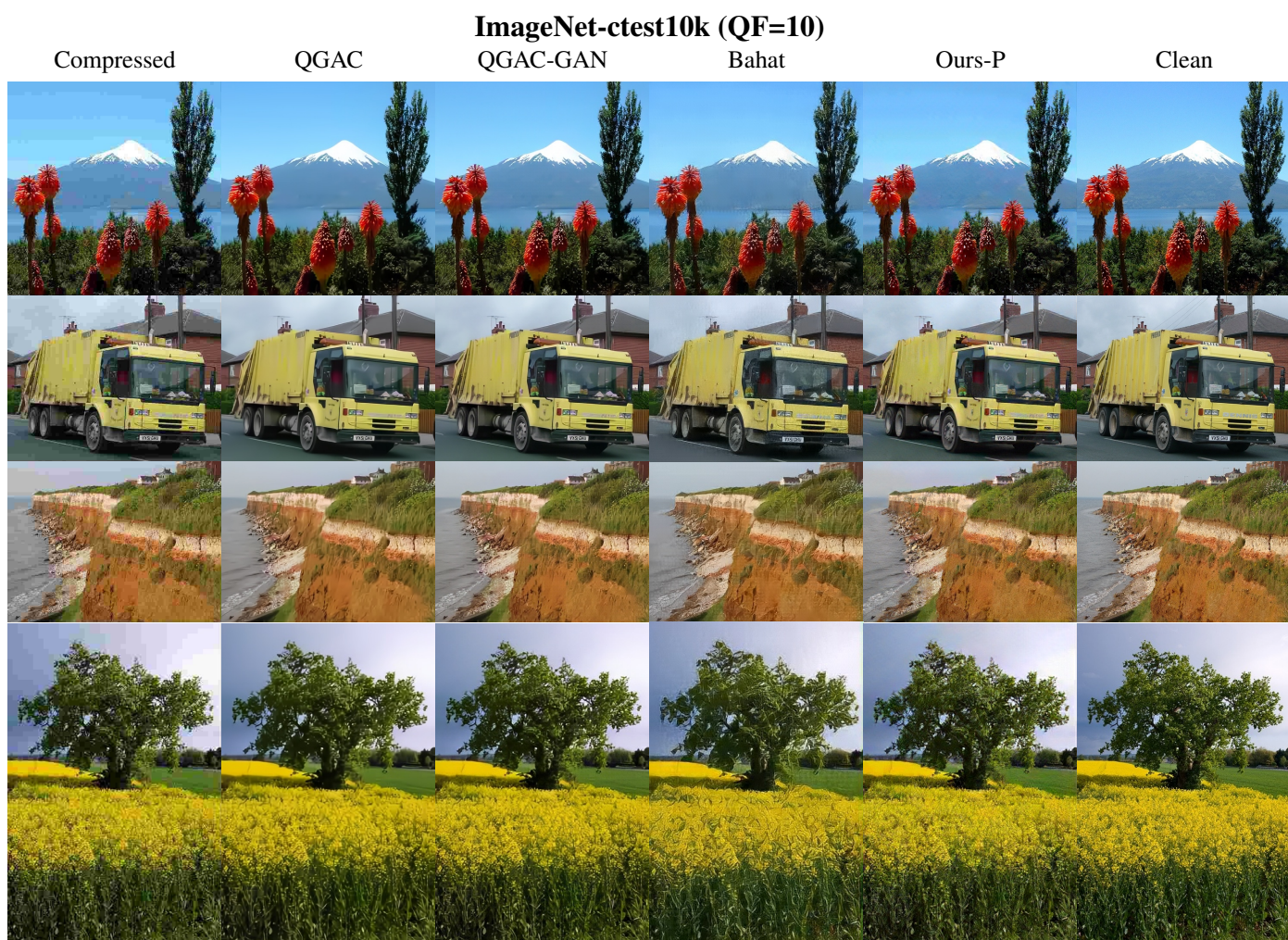


Figure 12. Decompression results using different methods on ImageNet JPEG compressed images with QF=10.



Figure 13. Decompression results using different methods on LIVE1 JPEG compressed images with QF=10.

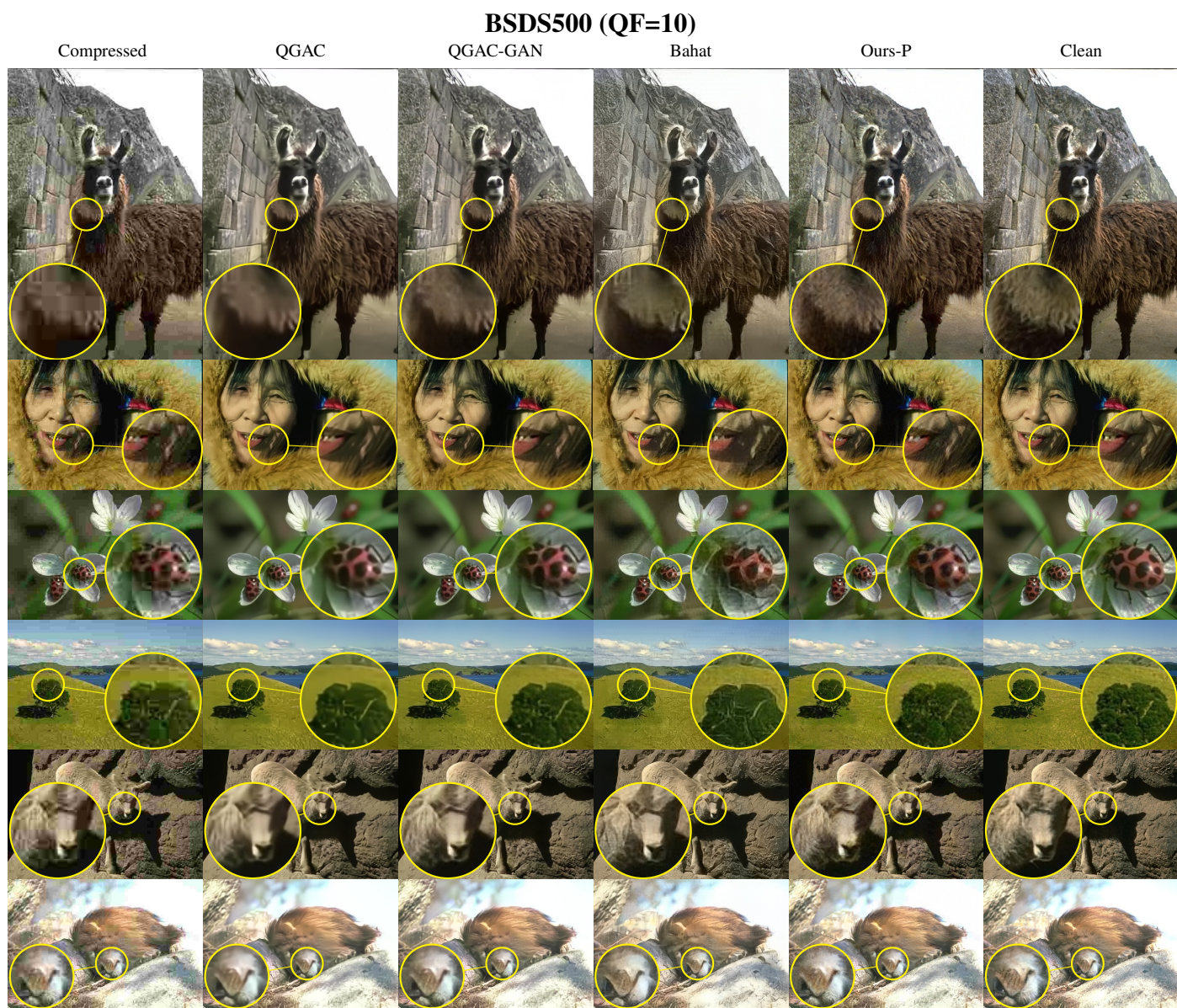


Figure 14. Decompression results using different methods on BSDS500 JPEG compressed images with QF=10.