# The Perception-Distortion Tradeoff

Yochai Blau and Tomer Michaeli

**Abstract**—Image restoration algorithms are typically evaluated by some distortion measure (e.g. PSNR, SSIM, IFC, VIF) or by human opinion scores that quantify perceived perceptual quality. In this paper, we prove mathematically that distortion and perceptual quality are at odds with each other. Specifically, we study the optimal probability for correctly discriminating the outputs of an image restoration algorithm from real images. We show that as the mean distortion decreases, this probability must increase (indicating worse perceptual quality). As opposed to the common belief, this result holds true for any distortion measure, and is not only a problem of the PSNR or SSIM criteria. We also show that generative-adversarial-nets (GANs) provide a principled way to approach the perception-distortion bound. This constitutes theoretical support to their observed success in low-level vision tasks. Based on our analysis, we propose a new methodology for evaluating image restoration methods, and use it to perform an extensive comparison between recent super-resolution algorithms.

✦

## 1 INTRODUCTION

THE last decades have seen continuous progress in image restoration algorithms (e.g. for denoising, deblurring, super-resolution) both in visual quality and in distortion measures like peak signal-to-noise ratio (PSNR) and structural similarity index (SSIM) [2]. However, in recent years, it seems that the improvement in reconstruction accuracy is not always accompanied by an improvement in visual quality. In fact, and perhaps counter-intuitively, algorithms that are superior in terms of perceptual quality, are often inferior in terms of e.g. PSNR and SSIM [3], [4], [5], [6], [7], [8], [9]. This phenomenon is commonly interpreted as a shortcoming of the existing distortion measures [10], which fuels a constant search for alternative "more perceptual" criteria.

In this paper, we offer a complementary explanation for the apparent tradeoff between perceptual quality and distortion measures. Specifically, we prove that there exists a region in the perception-distortion plane, which cannot be attained regardless of the algorithmic scheme (see Fig. 1). Furthermore, the boundary of this region is monotone. Therefore, in its proximity, it is only possible to improve *either perceptual quality or distortion*, one at the expense of the other. The perception-distortion tradeoff exists for *all distortion measures*, and is not only a problem of the mean-square error (MSE) or SSIM criteria.

Let us clarify the difference between distortion and perceptual quality. The goal in image restoration is to estimate an image $x$ from its degraded version $y$ (e.g. noisy, blurry, etc.). Distortion refers to the dissimilarity between the reconstructed image $\hat{x}$ and the original image $x$. Perceptual
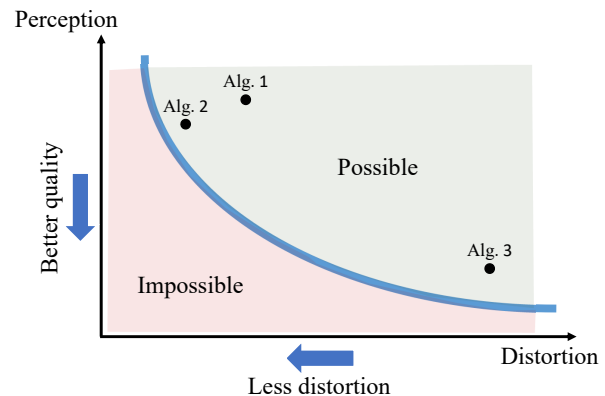


Fig. 1. **The perception-distortion tradeoff.** Image restoration algorithms can be characterized by their average distortion and by the perceptual quality of the images they produce. We show that there exists a region in the perception-distortion plane which cannot be attained, regardless of the algorithmic scheme. When in proximity of this unattainable region, an algorithm can be potentially improved only in terms of its distortion *or* in terms of its perceptual quality, one at the expense of the other.

quality, on the other hand, refers only to the visual quality of $\hat{x}$, regardless of its similarity to $x$. Namely, it is the extent to which $\hat{x}$ looks like a valid natural image. An increasingly popular way of measuring perceptual quality is by using real-vs.-fake user studies, which examine the ability of human observers to tell whether $\hat{x}$ is real or the output of an algorithm [5], [11], [12], [13], [14], [15], [16], [17] (similarly to the idea underlying generative adversarial nets [18]). Therefore, perceptual quality can be defined as the best possible probability of success in such discrimination experiments, which as we show, is proportional to the distance between the distribution of reconstructed images and that of natural images.

Based on these definitions of perception and distortion, we follow the logic of rate-distortion theory [19]. That is, we seek to characterize the behavior of the best attainable perceptual quality (minimal deviation from natural image

---

*Y. Blau and T. Michaeli are with the Technion – Israel Institute of Technology, Haifa, Israel. E-mail: {yochai@campus, tomer.m@ee}.technion.ac.il*

statistics) as a function of the maximal allowable average distortion, for any estimator. This perception-distortion function (wide curve in Fig. 1) separates between the attainable and unattainable regions in the perception-distortion plane and thus describes the fundamental tradeoff between perception and distortion. Our analysis shows that algorithms cannot be simultaneously very accurate *and* produce images that fool observers to believe they are real, no matter what measure is used to quantify accuracy. This tradeoff implies that optimizing distortion measures can be not only ineffective, but also potentially damaging in terms of visual quality. This has been empirically observed e.g. in [3], [4], [5], [6], [7], but was never established theoretically.

From the standpoint of algorithm design, we show that generative adversarial nets (GANs) provide a principled way to approach the perception-distortion bound. This gives theoretical support to the growing empirical evidence of the advantages of GANs in image restoration [3], [6], [7], [11], [20], [21], [22].

The perception-distortion tradeoff has major implications on low-level vision. In certain applications, reconstruction accuracy is of key importance (e.g. medical imaging). In others, perceptual quality may be preferred. The impossibility of simultaneously achieving both goals calls for a new way for evaluating algorithms: By placing them on the perception-distortion plane. We use this new methodology to conduct an extensive comparison between recent super-resolution (SR) methods, revealing which SR methods lie closest to the perception-distortion bound.

## 2 DISTORTION AND PERCEPTUAL QUALITY

Distortion and perceptual quality have been studied in many different contexts, and are sometimes referred to by different names. Let us briefly put past works in our context.

### 2.1 Distortion (full-reference) measures

Given a distorted image $\hat{x}$ and a ground-truth reference image $x$, full-reference distortion measures quantify the quality of $\hat{x}$ by its discrepancy to $x$. These measures are often called full reference image quality criteria because of the reasoning that if $\hat{x}$ is similar to $x$ and $x$ is of high quality, then $\hat{x}$ is also of high quality. However, as we show in this paper, this logic is not always correct. We thus prefer to call these measures distortion or dissimilarity criteria.

The most common distortion measure is the MSE, which is quite poorly correlated with semantic similarity between images [10]. Many alternative, more perceptual, distortion measures have been proposed over the years, including SSIM [2], MS-SSIM [23], IFC [24], VIF [25], VSNR [26] and FSIM [27]. Recently, measures based on the $\ell_2$-distance between deep feature maps of a neural-net have been shown to capture more semantic similarities [3], [4], [28].

### 2.2 Perceptual quality

The perceptual quality of an image $\hat{x}$ is the degree to which it looks like a natural image, and has nothing to do with its similarity to any reference image. In many image processing domains, perceptual quality has been associated with deviations from natural image statistics.

### Human opinion based quality assessment

Perceptual quality is commonly evaluated empirically by the mean opinion score of human subjects [29], [30]. Recently, it has become increasingly popular to perform such studies through real vs. fake questionnaires [5], [11], [12], [13], [14], [15], [16], [17]. These test the ability of a human observer to distinguish whether an image is real or the output of some algorithm. The probability of success $p_{\text{success}}$ of the optimal decision rule in this hypothesis testing task is known to be (see Appendix A in the Supplementary Material)

$$p_{\text{success}} = \tfrac{1}{2}d_{\text{TV}}(p_X, p_{\hat{X}}) + \tfrac{1}{2}, \tag{1}$$

where $d_{\text{TV}}(p_X, p_{\hat{X}})$ is the total-variation (TV) distance between the distribution $p_{\hat{X}}$ of images produced by the algorithm in question, and the distribution $p_X$ of natural images [31]. Note that $p_{\text{success}}$ decreases as the deviation between $p_{\hat{X}}$ and $p_X$ decreases, becoming $\tfrac{1}{2}$ (no better than a coin toss) when $p_{\hat{X}} = p_X$.

### No-reference quality measures

Perceptual quality can also be measured by an algorithm. In particular, no-reference measures quantify the perceptual quality $Q(\hat{x})$ of an image $\hat{x}$ *without* depending on a reference image. These measures are commonly based on estimating deviations from natural image statistics. The works [32], [33], [34] proposed perceptual quality indices based on the KL divergence between the distribution of the wavelet coefficients of $\hat{x}$ and that of natural scenes. This idea was further extended by the popular methods DIIVINE [29], BRISQUE [30], BLIINDS-II [35] and NIQE [36], which quantify perceptual quality by various measures of deviation from natural image statistics in the spatial, wavelet and DCT domains.

### GAN-based image restoration

Most recently, GAN-based methods have demonstrated unprecedented perceptual quality in super-resolution [3], [6], [9], [37], inpainting [7], [20], [38], compression [21], [39], [40], deblurring [41] and image-to-image translation [11], [22], [42]. This was accomplished by utilizing an adversarial loss, which minimizes some distance $d(p_X, p_{\hat{X}_{\text{GAN}}})$ between the distribution $p_{\hat{X}_{\text{GAN}}}$ of images produced by the generator and the distribution $p_X$ of images in the training dataset. A large variety of GAN schemes have been proposed, which minimize different distances between distributions. These include the Jensen-Shannon divergence [18], the Wasserstein distance [43], and any $f$-divergence [44].

### Single image quality vs. image ensemble quality

A common measure of quality is the log-likelihood $Q_{\text{LL}}(\hat{x}) = \log(p_X(\hat{x}))$. However, this notion of quality evaluates each image individually, and therefore has shortcomings. As an example, consider a reconstruction algorithm that disregards the input image, and always outputs the same "good-looking" natural image that has a high likelihood. While this algorithm would rate very well by the average log-likelihood measure $\mathbb{E}[Q_{\text{LL}}(\hat{X})]$, it is obviously quite useless. An observer examining an *ensemble* of outputs, would easily notice this flawed behavior. Therefore, in this paper we are more interested in "ensemble quality". The

Original    Degraded    Reconstructed

$X$     $Y$     $\hat{X}$

$p_X$     $p_{\hat{X}}$

$p_{Y|X}$     $p_{\hat{X}|Y}$

Distortion:
$$\mathbb{E}[\Delta(X, \hat{X})]$$

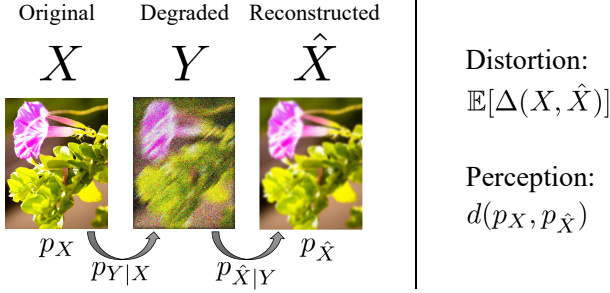Perception:
$$d(p_X, p_{\hat{X}})$$

Fig. 2. **Problem setting.** Given an original image $x \sim p_X$, a degraded image $y$ is observed according to some conditional distribution $p_{Y|X}$. Given the degraded image $y$, an estimate $\hat{x}$ is constructed according to some conditional distribution $p_{\hat{X}|Y}$. Distortion is quantified by the mean of some distortion measure between $\hat{X}$ and $X$. The perceptual quality index corresponds to the deviation between $p_{\hat{X}}$ and $p_X$.

relation between the average log-likelihood and "ensemble quality" can be understood by noting that

$$\mathbb{E}_{\hat{X} \sim p_{\hat{X}}}[Q_{\text{LL}}(\hat{X})] = -d_{\text{KL}}(p_{\hat{X}}, p_X) - H(p_{\hat{X}}). \quad (2)$$

Here $d_{\text{KL}}$ is the Kullback-Leibler divergence and $H$ denotes entropy. The second term in this decomposition discourages diversity. Choosing to drop it, results in the distributional-divergence based quality measures described above.

## 3 PROBLEM FORMULATION

In statistical terms, a natural image $x$ can be thought of as a realization from the distribution of natural images $p_X$. In image restoration, we observe a degraded version $y$ relating to $x$ via some conditional distribution $p_{Y|X}$. In this paper we focus on non-invertible settings[1], where $x$ cannot be estimated from $y$ with zero error. This is typically the case in denoising, deblurring, inpaitning, super-resolution, etc. Given $y$, an image restoration algorithm produces an estimate $\hat{x}$ according to some distribution $p_{\hat{X}|Y}$. Note that this description is quite general in that it does not restrict the estimator $\hat{x}$ to be a deterministic function of $y$. This problem setting is illustrated in Fig. 2.

Given a full-reference dissimilarity criterion $\Delta(x, \hat{x})$, the average distortion of an estimator $\hat{X}$ is given by

$$\mathbb{E}[\Delta(X, \hat{X})], \quad (3)$$

where the expectation is over the joint distribution $p_{X,\hat{X}}$. This definition aligns with the common practice of evaluating average performance over a database of degraded natural images. We assume that the dissimilarity criterion is such that $\Delta(x, \hat{x}) \geq 0$ with equality when $\hat{x} = x$. Note that some distortion measures, e.g. SSIM, are actually *similarity* measures (higher is better), yet can always be inverted (and shifted) to become dissimilarity measures.

As discussed in Sec. 2.2, the perceptual quality of an estimator $\hat{X}$ (as quantified e.g. by real vs. fake human opinion studies) is directly related to the distance between the distribution of its reconstructed images, $p_{\hat{X}}$, and the distribution of natural images, $p_X$. We thus define the

---

1. By invertible we mean that the support of $p_{X|Y}(\cdot|y)$ is a singleton for almost all $y$'s (see Appendix C for a formal definition).

perceptual quality index (lower is better) of an estimator $\hat{X}$ as

$$d(p_X, p_{\hat{X}}), \quad (4)$$

where $d(p, q)$ is some divergence between distributions that satisfies $d(p, q) \geq 0$ with equality if $p = q$, e.g. the KL divergence, TV distance, Wasserstein distance, etc. It should be pointed out that the divergence function $d(\cdot, \cdot)$ which best relates to human perception is a subject of ongoing research. Yet, our results below hold for (nearly) any divergence.

Notice that the best possible perceptual quality is obtained when the outputs of the algorithm follow the distribution of natural images (i.e. $p_{\hat{X}} = p_X$). In this situation, by looking at the reconstructed images, it is impossible to tell that they were generated by an algorithm. However, not every estimator with this property is necessarily accurate. Indeed, we could achieve perfect perceptual quality by randomly drawing natural images that have nothing to do with the original ground-truth images. In this case the distortion would be quite large.

Our goal is to characterize the tradeoff between (3) and (4). But let us first see why minimizing the average distortion (3), does not necessarily lead to a low perceptual quality index (4). We start by illustrating this with the square-error distortion $\Delta(x, \hat{x}) = \|x - \hat{x}\|^2$ and the $0 - 1$ distortion $\Delta(x, \hat{x}) = 1 - \delta_{x,\hat{x}}$ (where $\delta$ is Kronecker's delta). We specifically illustrate that those measures are generally not distribution preserving in the following sense.

**Definition 1.** *We say that a distortion measure $\Delta(\cdot, \cdot)$ is distribution preserving at $p_{X,Y}$ if the estimator $\hat{X}$ minimizing the mean distortion (3) satisfies $p_{\hat{X}} = p_X$.*

More details about those examples are provided in Appendix B (Supplementary Material). We then proceed to discuss this phenomenon for arbitrary distortions in Sec. 3.3.

### 3.1 The square-error distortion

The minimum mean square-error (MMSE) estimator is given by the posterior-mean $\hat{x}(y) = \mathbb{E}[X|Y = y]$. Consider the case $Y = X + N$, where $X$ is a discrete random variable with probability mass function

$$p_X(x) = \begin{cases} p_1 & x = \pm 1, \\ p_0 & x = 0, \end{cases} \quad (5)$$

and $N \sim \mathcal{N}(0, 1)$ is independent of $X$ (see Fig. 3). In this setting, the MMSE estimate is given by

$$\hat{x}_{\text{MMSE}}(y) = \sum_{n \in \{-1,0,1\}} n\, p(X = n|y), \quad (6)$$

where

$$p(X = n|y) = \frac{p_n \exp\{-\frac{1}{2}(y - n)^2\}}{\sum_{m \in \{-1,0,1\}} p_m \exp\{-\frac{1}{2}(y - m)^2\}}. \quad (7)$$

Notice that $\hat{x}_{\text{MMSE}}$ can take any value in the range $(-1, 1)$, whereas $x$ can only take the discrete values $\{-1, 0, 1\}$. Thus, clearly, $p_{\hat{X}_{\text{MMSE}}}$ is very different from $p_X$, as illustrated in Fig. 3. This demonstrates that minimizing the MSE distortion *does not* generally lead to $p_{\hat{X}} \approx p_X$.
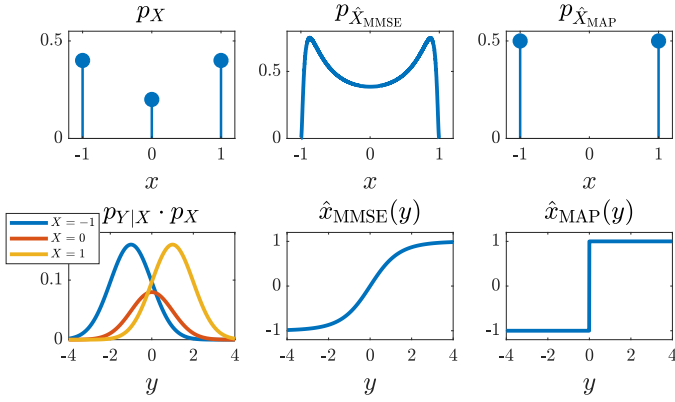
Fig. 3. **The distribution of the MMSE and MAP estimates.** In this example, $Y = X + N$, where $X \sim p_X$ and $N \sim \mathcal{N}(0, 1)$. The distributions of both the MMSE and the MAP estimates deviate significantly from the distribution $p_X$.

The same intuition holds for images. The MMSE estimate is an average over all possible explanations to the measured data, weighted by their likelihoods. However the average of valid images is not necessarily a valid image, so that the MMSE estimate frequently "falls off" the natural image manifold [3]. This leads to unnatural blurry reconstructions, as illustrated in Fig. 4. In this experiment, $x$ is a $280 \times 280$ image comprising 100 smaller $28 \times 28$ digit images. Each digit is chosen uniformly at random from a dataset comprising 54K images from the MNIST dataset [45] and an additional 5.4K blank images. The degraded image $y$ is a noisy version of $x$. As can be seen, the MMSE estimator produces blurry reconstructions, which do not follow the statistics of the (binary) images in the dataset.

### 3.2 The $0 - 1$ distortion

The discussion above may give the impression that unnatural estimates are mainly a problem of the square-error distortion, which causes averaging. One way to avoid averaging, is to minimize the binary $0-1$ loss, which restricts the estimator to choose $\hat{x}$ only from the set of values that $x$ can take. In fact, the minimum mean $0 - 1$ distortion is attained by the maximum-a-posteriori (MAP) rule, which is very popular in image restoration. However, as we exemplify next, the distribution of the MAP estimator also deviates from $p_X$. This behavior has also been studied in [46].

Consider again the setting of (5). In this case, the MAP estimate is given by

$$\hat{x}_{\text{MAP}}(y) = \underset{n \in \{-1,0,1\}}{\arg\max} \ p(X = n | y), \tag{8}$$

where $p(X = n | y)$ is as in (7). Now, it can be easily verified that when $\log(p_1/p_0) > 1/2$, we have $\hat{x}_{\text{MAP}}(y) = \text{sign}(y)$. Namely, the MAP estimator never predicts the value 0. Therefore, in this case, the distribution of the estimate is

$$p_{\hat{X}_{\text{MAP}}}(\hat{x}) = \begin{cases} 0.5 & \hat{x} = +1, \\ 0.5 & \hat{x} = -1, \end{cases} \tag{9}$$

which is obviously different from $p_X$ of (5) (see Fig. 3).

This effect can also be seen in the experiment of Fig. 4. Here, the MAP estimates become increasingly dominated
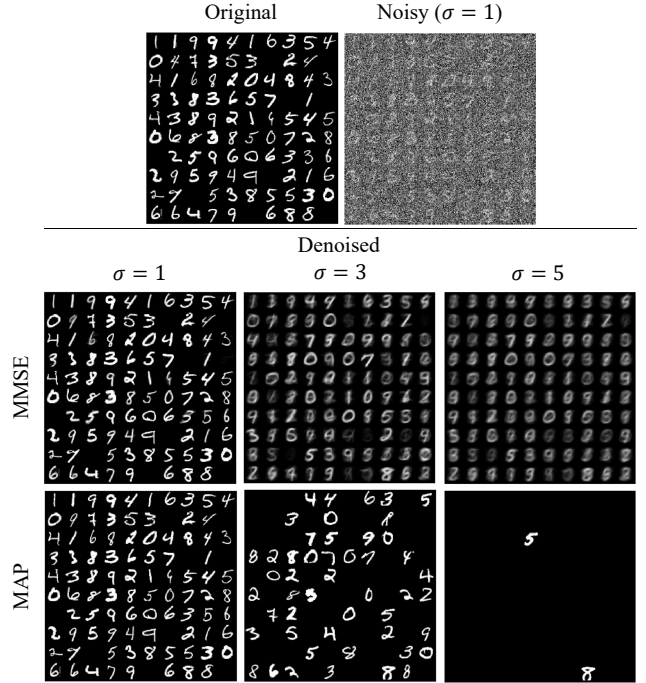


Fig. 4. **MMSE and MAP denoising.** Here, the original image consists of 100 smaller images, chosen uniformly at random from the MNIST dataset enriched with blank images. After adding Gaussian noise ($\sigma = 1, 3, 5$), the image is denoised using the MMSE and MAP estimators. In both cases, the estimates significantly deviate from the distribution of images in the dataset.

by blank images as the noise level rises, and thus clearly deviates from the underlying prior distribution.

### 3.3 Arbitrary distortion measures

We saw that neither the square-error nor the $0 - 1$ loss are distribution preserving. That is, their minimization does not generally lead to $p_{\hat{X}} = p_X$ (i.e. perfect perceptual quality). However these two examples do not yet preclude the existence of a distribution preserving distortion measure. Does there exist a measure whose minimization is guaranteed to lead to $p_{\hat{X}} = p_X$? If we limit ourselves to one single setting, then the answer may be positive. For example, in the setting of Fig. 3, if $p_0$ of (5) equals 0, then the $0 - 1$ loss is distribution preserving as its minimization leads to an estimate satisfying $p_{\hat{X}} = p_X$. This illustrates that a distortion measure may be distribution preserving for certain underlying distributions $p_{X,Y}$ but not for others.

However, from a practical standpoint, we typically want our distortion measure to be adequate in more than one single setting. For example, if our goal is to train a neural network to perform denoising, then it is reasonable to expect that the same distortion measure be equally adequate as a loss function for different noise levels. In fact, we may also want to use the same distortion measure across different tasks (e.g. super-resolution, deblurring, inpainting). The interesting question is, therefore, whether there exists a *stably* distribution preserving distortion measure.

**Definition 2.** *We say that a distortion measure* $\Delta(\cdot, \cdot)$ *is* stably distribution preserving *at* $p_{X,Y}$ *if it is distribution preserving at all* $\tilde{p}_{X,Y}$ *in a TV* $\varepsilon$-*ball around* $p_{X,Y}$ *for some* $\varepsilon > 0$.
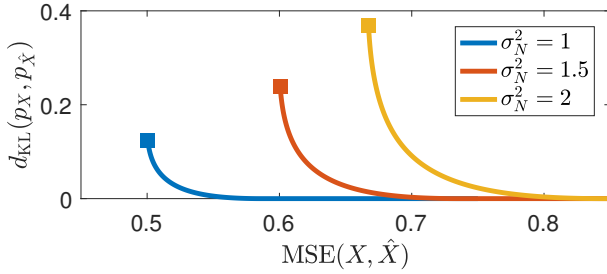
Fig. 5. **Plot of Eq.** (10) **for the setting of Example 1.** The minimal attainable KL distance between $p_X$ and $p_{\hat{X}}$ subject to a constraint on the maximal allowable MSE between $X$ and $\hat{X}$. Here, $Y = X + N$, where $X \sim \mathcal{N}(0,1)$ and $N \sim \mathcal{N}(0, \sigma_N)$, and the estimator is linear, $\hat{X} = aY$. Notice the clear trade-off: The perceptual index ($d_{\mathsf{KL}}$) drops as the allowable distortion (MSE) increases. The graphs cut-off at the MMSE (marked by a square).

As we show next, if the degradation is non-invertible, then no distortion metric can be stably distribution preserving (see proof in Appendix C).

**Theorem 1.** *If $p_{X,Y}$ defines a non-invertible degradation, then $\Delta(\cdot, \cdot)$ is not a stably distribution preserving distortion at $p_{X,Y}$.*

## 4 THE PERCEPTION-DISTORTION TRADEOFF

We saw that for any distortion measure, a low distortion does not generally imply good perceptual-quality. An interesting question, then, is: What is the best perceptual quality that can be attained by an estimator with a prescribed distortion level?

**Definition 3.** *The perception-distortion function of a signal restoration task is given by*

$$P(D) = \min_{p_{\hat{X}|Y}} d(p_X, p_{\hat{X}}) \quad s.t. \quad \mathbb{E}[\Delta(X, \hat{X})] \leq D, \quad (10)$$

*where $\Delta(\cdot, \cdot)$ is a distortion measure and $d(\cdot, \cdot)$ is a divergence between distributions.*

In words, $P(D)$ is the minimal deviation between the distributions $p_X$ and $p_{\hat{X}}$ that can be attained by an estimator with distortion $D$. To gain intuition into the typical behavior of this function, consider the following example.

**Example 1.** *Suppose that $Y = X + N$, where $X \sim \mathcal{N}(0,1)$ and $N \sim \mathcal{N}(0, \sigma_N^2)$ are independent. Take $\Delta(\cdot, \cdot)$ to be the square-error distortion and $d(\cdot, \cdot)$ to be the KL divergence. For simplicity, let us focus on estimators of the form $\hat{X} = aY$. In this case, we can derive a closed form solution to Eq.* (10) *(see Appendix D), which is plotted for several noise levels $\sigma_N$ in Fig. 5. As can be seen, the minimal attainable $d_{\mathsf{KL}}(p_X, p_{\hat{X}})$ drops as the maximal allowable distortion (MSE) increases. Furthermore, the tradeoff is convex and becomes more severe at higher noise levels $\sigma_N$.*

In general settings, it is impossible to solve (10) analytically. However, it turns out that the behavior seen in Fig. 5 is typical, as we show next (see proof in Appendix E).

**Theorem 2** (The perception-distortion tradeoff). *Assume the problem setting of Section 3. If $d(p, q)$ of* (4) *is convex in its*

second argument[2], *then the perception-distortion function $P(D)$ of* (10) *is*
1) *monotonically non-increasing;*
2) *convex.*

Note that Theorem 2 requires no assumptions on the distortion measure $\Delta(\cdot, \cdot)$. This implies that a tradeoff between perceptual quality and distortion exists for *any distortion measure*, including e.g. MSE, SSIM, square error between VGG features [3], [4], etc. Yet, this does not imply that all distortion measures have the same perception-distortion function. Indeed, as we demonstrate in Sec. 6, the tradeoff tends to be less severe for distortion measures that capture semantic similarities between images.

The convexity of $P(D)$ implies that the tradeoff is more severe at the low-distortion and at the high-perceptual-quality extremes. This is particularly important when considering the TV divergence which is associated with the ability to distinguish between real vs. fake images (see Sec. 2.2). Since $P(D)$ is steeper at the low-distortion regime, any *small* improvement in distortion for an algorithm whose distortion is already low, must be accompanied by a *large* degradation in the ability to fool a discriminator. Similarly, any *small* improvement in the perceptual quality of an algorithm whose perceptual index is already low, must be accompanied by a *large* increase in distortion. Let us comment that the assumption that $d(p, q)$ is convex, is not very limiting. For instance, any $f$-divergence (e.g. KL, TV, Hellinger, $\mathcal{X}^2$) as well as the Renyi divergence, satisfy this assumption [47], [48]. In any case, the function $P(D)$ is monotonically non-increasing even without this assumption.

### 4.1 Bounding the Perception-Distortion function

Several past works attempted to answer the question: What is the minimal attainable distortion $D_{\min}$ in various restoration tasks? [49], [50], [51], [52], [53]. This corresponds to the value

$$D_{\min} = \min_{p_{\hat{X}|Y}} \mathbb{E}[\Delta(X, \hat{X})], \quad (11)$$

which is the horizontal coordinate of the leftmost point on the perception-distortion function. However, as the minimum distortion estimator is generally not distribution preserving (Sec. 3.3), an important complementary question is: What is the minimal distortion that can be attained by an estimator *having perfect perceptual quality*? This corresponds to the value

$$D_{\max} = \min_{p_{\hat{X}|Y}} \mathbb{E}[\Delta(X, \hat{X})] \quad s.t. \quad p_{\hat{X}} = p_X, \quad (12)$$

which is the horizontal coordinate of the point where the perception-distortion function first touches the horizontal axis (see Fig. 6).

Observe that perfect perceptual quality ($p_{\hat{X}} = p_X$) is always attainable, for example by drawing $\hat{x}$ from $p_X$ independently of the input $y$. This method, however, ignores the input and is thus not good in terms of distortion. It turns out that perfect perceptual quality can generally be achieved with a significantly lower MSE distortion, as we show next (see proof in Appendix F).

---

2. That is, $d(p, \lambda q_1 + (1-\lambda)q_2) \leq \lambda d(p, q_1) + (1-\lambda)d(p, q_2)$ for any three distributions $p, q_1, q_2$ and any $\lambda \in [0, 1]$.
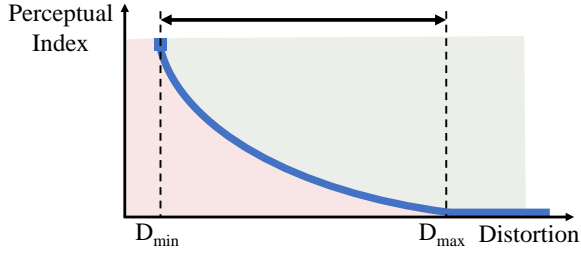
Fig. 6. **Bounding the perception-distortion function.** The distance between $D_{\min}$ and $D_{\max}$ is the increase in distortion which is needed to obtain perfect perceptual quality. For the MSE, Theorem 3 proves this will never be more than a factor of $2$ (which is $3$dB in terms of PSNR).

**Theorem 3.** *For the square error distortion* $\Delta(x, \hat{x}) = \|\hat{x} - x\|^2$,

$$D_{\max} \le 2 D_{\min}, \tag{13}$$

*where $D_{\min}$ and $D_{\max}$ are defined by* (11) *and* (12), *respectively. This bound is attained by the estimator $\hat{X}$ defined through*

$$p_{\hat{X}|Y}(x|y) = p_{X|Y}(x|y), \tag{14}$$

*which achieves $p_{\hat{X}} = p_X$ and has an MSE of $2 D_{\min}$.*

In simple words, Theorem 3 states that one would never need to sacrifice more than 3dB in PSNR to obtain perfect perceptual quality. This can be achieved by drawing $\hat{x}$ from the posterior distribution $p_{X|Y}$. Interestingly, such a degradation was indeed incurred by all super-resolution methods that achieved state-of-the-art perceptual quality to date. This can be seen in Fig. 9, where the RMSE of the algorithms with the lowest perceptual index is nearly a factor of $\sqrt{2}$ larger than the RMSE of the methods with the lowest RMSE (see also [3], [6]). However, note that this bound is generally not tight. For example, in the scalar Gaussian toy example of Fig. 5, $D_{\max}$ can be quite smaller than $2 D_{\min}$, depending on the noise level.

### 4.2 Connection to rate-distortion theory

The perception-distortion tradeoff is closely related to the well-established rate-distortion theory [19]. This theory characterizes the tradeoff between the bit-rate required to communicate a signal, and the distortion incurred in the signal's reconstruction at the receiver. More formally, the rate-distortion function of a signal $X$ is defined by

$$R(D) = \min_{p_{\hat{X}|X}} I(X; \hat{X}) \quad \text{s.t.} \quad \mathbb{E}[\Delta(X, \hat{X})] \le D, \tag{15}$$

where $I(X; \hat{X})$ is the mutual information between $X$ and $\hat{X}$.

There are, however, several key differences between the two tradeoffs. First, in rate-distortion the optimization is over all conditional distributions $p_{\hat{X}|X}$, i.e. given the *original* signal. In the perception-distortion case, the estimator has access only to the degraded signal $Y$, so that the optimization is over the conditional distributions $p_{\hat{X}|Y}$, which is more restrictive. In other words, the perception-distortion tradeoff depends on the degradation $p_{Y|X}$, and not only on the signal's distribution $p_X$ (see Example 1). Second, in rate-distortion the rate is quantified by the mutual information $I(X; \hat{X})$, which depends on the joint distribution $p_{X, \hat{X}}$. In our case, perception is quantified by the similarity between

$p_X$ and $p_{\hat{X}}$, which does not depend on their joint distribution. Lastly, mutual information is inherently convex, while the convexity of the perception-distortion curve is guaranteed only when $d(\cdot, \cdot)$ is convex.

While the two tradeoffs are different, it is important to note that perceptual quality does play a role in lossy compression, as evident from the success of recent GAN based compression schemes [39], [40], [54]. Theoretically, its effect can be studied through the rate-distortion-perception function [55], [56], [57], which is an extension of the rate-distortion function (15) and the perception-distortion function (10), characterizing the triple tradeoff between rate, distortion, and perceptual quality.

## 5 TRAVERSING THE TRADEOFF WITH A GAN

There exists a systematic way to design estimators that approach the perception-distortion curve: Using GANs. Specifically, motivated by [3], [6], [7], [11], [20], [21], restoration problems can be approached by modifying the loss of the generator of a GAN to be

$$\ell_{\text{gen}} = \ell_{\text{distortion}} + \lambda \, \ell_{\text{adv}}, \tag{16}$$

where $\ell_{\text{distortion}}$ is the distortion between the original and reconstructed images, and $\ell_{\text{adv}}$ is the standard GAN adversarial loss. It is well known that $\ell_{\text{adv}}$ is proportional to some divergence $d(p_X, p_{\hat{X}})$ between the generator and data distributions [18], [43], [44] (the type of divergence depends on the loss). Thus, (16) in fact approximates the objective

$$\ell_{\text{gen}} \approx \mathbb{E}[\Delta(x, \hat{x})] + \lambda \, d(p_X, p_{\hat{X}}). \tag{17}$$

Viewing $\lambda$ as a Lagrange multiplier, it is clear that minimizing $\ell_{\text{gen}}$ is equivalent to minimizing (10) for some $D$. Varying $\lambda$ corresponds to varying $D$, thus producing estimators along the perception-distortion function.

Let us use this approach to explore the perception-distortion tradeoff for the digit denoising example of Fig. 4 with $\sigma = 3$. We train a Wasserstein GAN (WGAN) based denoiser [43], [58] with an MSE distortion loss $\ell_{\text{distortion}}$. Here, $\ell_{\text{adv}}$ is proportional to the Wasserstein distance $d_W(p_X, p_{\hat{X}})$ between the generator and data distributions. The WGAN has the valuable property that its discriminator (critic) loss is an accurate estimate (up to a constant factor) of $d_W(p_X, p_{\hat{X}})$ [43]. This allows us to easily compute the perceptual quality index of the trained denoiser. We obtain a set of estimators with several values of $\lambda \in [0, 0.3]$. For each denoiser, we evaluate the perceptual quality by the final discriminator loss. As seen in Fig. 7, the curve connecting the estimators on the perception-distortion plane is monotonically decreasing. Moreover, it is associated with estimates that gradually transition from blurry and accurate to sharp and inaccurate. This curve obviously does not coincide with the analytic bound (10) (illustrated by a dashed line). However, it seems to be adjacent to it. This is indicated by the fact that the left-most point of the WGAN curve is very close to the left-most point of the theoretical bound, which corresponds to the MMSE estimator. See Appendix G for the WGAN training details and architecture.

Besides the MMSE estimator, Figure 7 also includes the MAP estimator, the random draw estimator $\hat{x} \sim p_X$ (which ignores the noisy image $y$), and the conditional
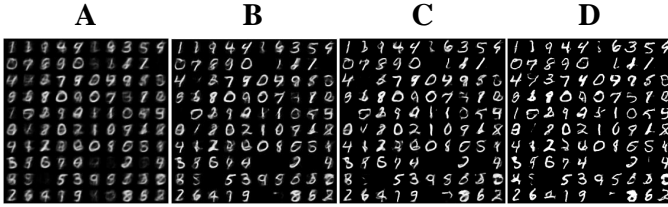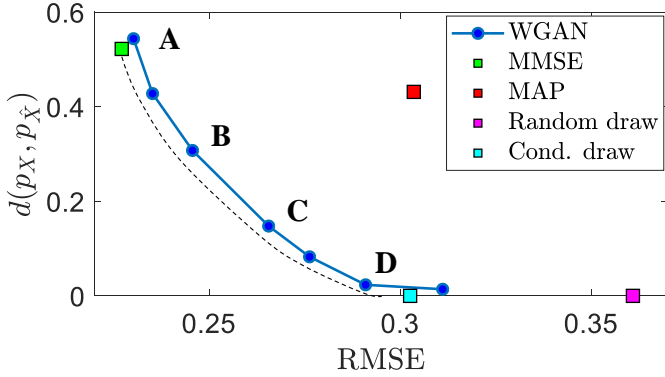
Fig. 7. **Image denoising utilizing a GAN.** A Wasserstein GAN was trained to denoise the images of the experiment in Fig. 4. The generator loss $l_{gen} = l_{MSE} + \lambda l_{adv}$ consists of a perceptual quality (adversarial) loss and a distortion (MSE) loss, where $\lambda$ controls the trade-off between the two. For each $\lambda \in [0, 0.3]$, the graph depicts the distortion (MSE) and perceptual quality (Wasserstein distance between $p_X$ and $p_{\hat{X}}$). The curve connecting the estimators is a good approximation to the theoretical perception-distortion tradeoff (illustrated by a dashed line).

draw estimator of (14). The perceptual quality of these estimators is evaluated, as above, by the final loss of the WGAN discriminator [43], trained (without a generator) to distinguish between the estimators' outputs and images from the dataset. Note that the denoising WGAN estimator (D) achieves the same distortion as the MAP estimator, but with far better perceptual quality. Furthermore, it achieves nearly the same perceptual quality as the random draw estimator, but with a significantly lower distortion.

## 6 PRACTICAL METHOD FOR EVALUATING ALGORITHMS

Certain applications may require low-distortion (e.g. in medical imaging), while others may prefer superior perceptual quality. How should image restoration algorithms be evaluated, then?

**Definition 4.** *We say that Algorithm A* dominates *Algorithm B if it has better perceptual quality* and *less distortion.*

Note that if Algorithm A is better than B in only one of the two criteria, then neither $A$ dominates $B$ nor $B$ dominates $A$. Therefore, among a group of algorithms, there may be a large subset which can be considered equally good.

**Definition 5.** *We say that an algorithm is* admissible *among a group of algorithms, if it is not dominated by any other algorithm in the group.*

As shown in Figure 8, these definitions have very simple interpretations when plotting algorithms on the perception-distortion plane. In particular, the admissible algorithms in
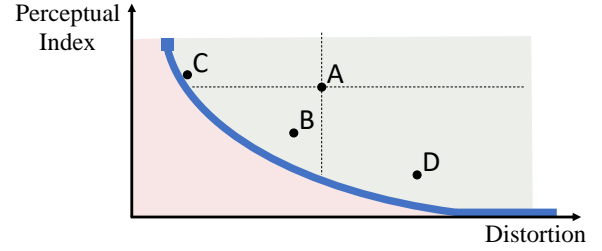


Fig. 8. **Dominance and admissibility.** Algorithm A is dominated by Algorithm B, and is thus inadmissible. Algorithms B, C and D are all admissible, as they are not dominated by any algorithm.

the group, are those which lie closest to the perception-distortion bound.

As discussed in Sec. 2, distortion is measured by *full-reference* (FR) metrics, e.g. [2], [4], [23], [24], [25], [26], [27]. The choice of the FR metric, depends on the type of similarities we want to measure (per-pixel, semantic, etc.). Perceptual quality, on the other hand, is ideally quantified by collecting human opinion scores, which is time consuming and costly [29], [35]. Instead, the divergence $d(p_X, p_{\hat{X}})$ can be computed, for instance by training a discriminator net (see Sec. 5). However, this requires *many* training images and is thus also time consuming. A practical alternative is to utilize *no*-reference (NR) metrics, e.g. [29], [30], [35], [36], [59], [60], [61], which quantify the perceptual quality of an image *without* a corresponding original image. In scenarios where NR metrics are highly correlated with human mean-opinion-scores (e.g. $4\times$ super-resolution [61]), they can be used as a fast and simple method for approximating the perceptual quality of an algorithm[3].

We use this approach to evaluate 16 SR algorithms in a $4\times$ magnification task, by plotting them on the perception-distortion plane (Fig. 9). We measure perceptual quality using the NR metric NIQE [36], which was shown to correlate well with human opinion scores in a recent SR challenge [64] (see Appendix H for experiments with the NR metrics BRISQUE [30], BLIINDS-II [35] and the recent NR metric by Ma et al. [61]). We measure distortion by the five common FR metrics RMSE, SSIM [2], MS-SSIM [23], IFC [24] and VIF [25], and additionally by the recent $VGG_{2,2}$ metric (the distance in the feature space of a VGG net) [3], [4]. To conform to previous evaluations, we compute all metrics on the y-channel after discarding a 4-pixel border (except for $VGG_{2,2}$, which is computed on RGB images). Comparisons on color images can be found in Appendix H. The algorithms are evaluated on the BSD100 dataset [65]. The evaluated algorithms include: A+ [66], SRCNN [67], SelfEx [68], VDSR [69], Johnson et al. [4], LapSRN [70], Bae et al. [71] ("primary" variant), EDSR [72], SRResNet variants which optimize MSE and $VGG_{2,2}$ [3], SRGAN variants which optimize MSE, $VGG_{2,2}$, and $VGG_{5,4}$, in addition to an adversarial loss [3], ENet [6] ("PAT" variant), Deng [73] ($\gamma = 0.55$), and Mechrez et al. [74].

---

3. In scenarios where NR metrics are inaccurate (e.g. blind deblurring with large blurs [62], [63]), the perceptual metric should be human-opinion-scores or the loss of a discriminator trained to distinguish the algorithms' outputs from natural images.
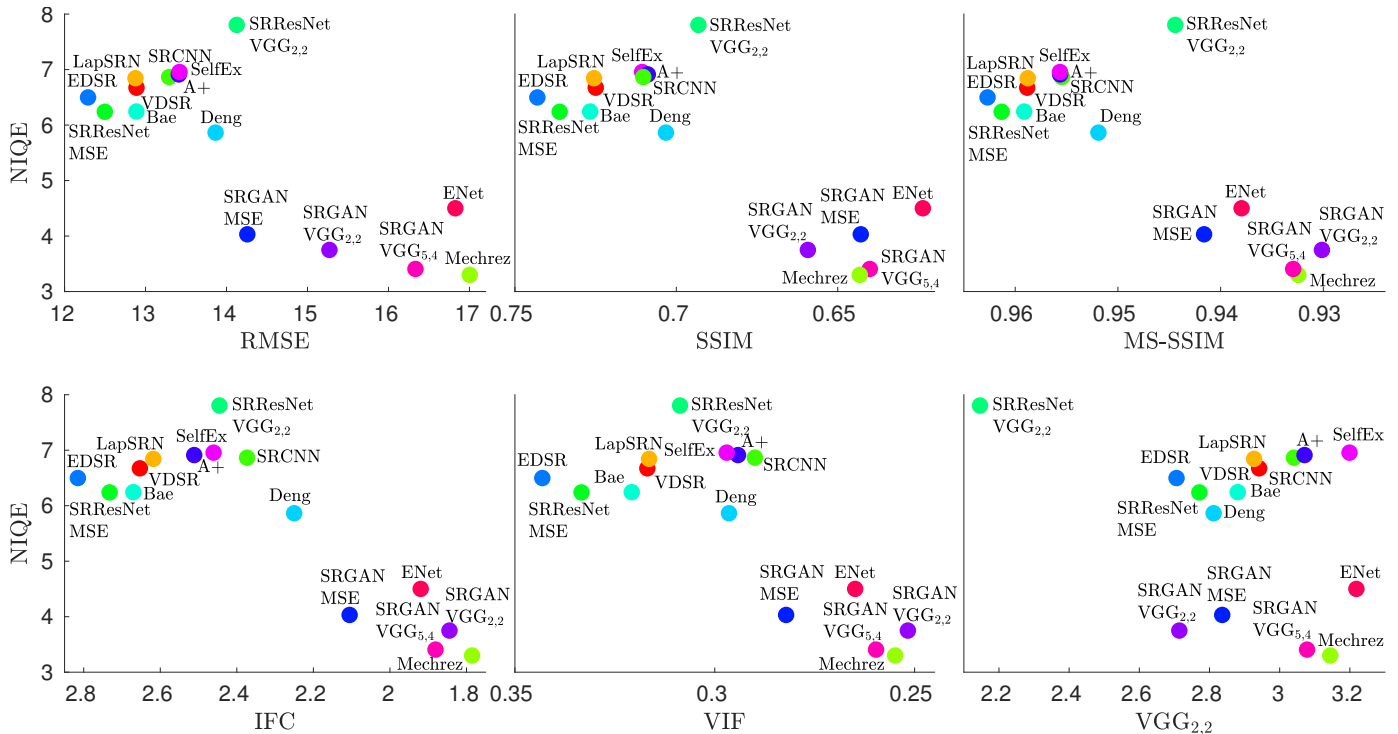
Fig. 9. **Perception-distortion evaluation of SR algorithms.** We plot 16 algorithms on the perception-distortion plane. Perception is measured by the NR metric NIQE [36]. Distortion is measured by the common full-reference metrics RMSE, SSIM, MS-SSIM, IFC, VIF and $VGG_{2,2}$. In all plots, the lower left corner is blank, revealing an unattainable region in the perception-distortion plane. In proximity of the unattainable region, an improvement in perceptual quality comes at the expense of higher distortion.



Fig. 10. **Visual comparison of algorithms closest to the perception-distortion bound.** The algorithms are ordered from low to high distortion (as evaluated by RMSE, MS-SSIM, IFC, VIF). Notice the co-occurring increase in perceptual quality.

Interestingly, the same pattern is observed in all plots: (i) The lower left corner is blank, revealing an unattainable region in the perception-distortion plane. (ii) In proximity of this blank region, NR and FR metrics are *anti-correlated*, indicating a tradeoff between perception and distortion. Notice that the tradeoff exists even for the IFC, VIF and $VGG_{2,2}$ measures, which are considered to capture visual quality better than MSE and SSIM.

Figure 10 depicts the outputs of several algorithms lying closest to the perception-distortion bound in the IFC graph in Fig. 9. While the images are ordered from low to high distortion (according to IFC), their perceptual quality clearly improves from left to right.

Both FR and NR measures are commonly validated by calculating their correlation with human opinion scores, based on the assumption that both should be correlated with perceptual quality. However, as Fig. 11 shows, while FR measures can be well-correlated with perceptual quality when distant from the unattainable region, this is clearly

not the case when approaching the perception-distortion bound. In particular, all tested FR methods are inconsistent with human opinion scores which found the SRGAN to be superb in terms of perceptual quality [3], while NR methods successfully determine this. We conclude that image restoration algorithms should always be evaluated by a pair of NR and FR metrics, constituting a reliable, reproducible and simple method for comparison, which accounts for both perceptual quality and distortion. This evaluation method was demonstrated and validated by a human opinion study in the 2018 PIRM super-resolution challenge [64].

Up until 2016, SR algorithms occupied only the upper-left section of the perception-distortion plane. Nowadays, emerging techniques are exploring new regions in this plane. The SRGAN, ENet, Deng, Johnson et al. and Mechrez et al. methods are the first (to our knowledge) to populate the high perceptual quality region. In the near future we will most likely witness continued efforts to approach the perception-distortion bound, not only in the low-distortion
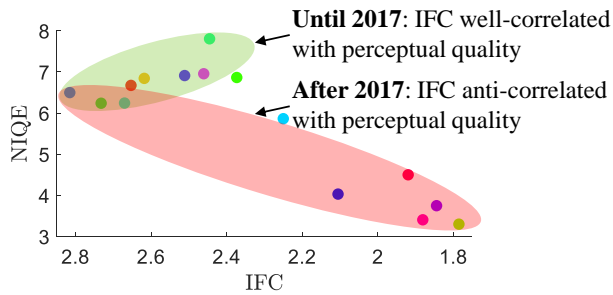
Fig. 11. **Correlation between distortion and perceptual quality.** In proximity of the perception-distortion bound, distortion and perceptual quality are *anti-correlated*. However, correlation is possible at distance from the bound.

region, but throughout the entire plane.

## 7 CONCLUSION

We proved and demonstrated the counter-intuitive phenomenon that distortion and perceptual quality are at odds with each other. Namely, the lower the distortion of an algorithm, the more its distribution must deviate from the statistics of natural scenes. We showed empirically that this tradeoff exists for many popular distortion measures, including those considered to be well-correlated with human perception. Therefore, any distortion measure alone, is unsuitable for assessing image restoration methods. Our novel methodology utilizes a pair of NR and FR metrics to place each algorithm on the perception-distortion plane, facilitating a more informative comparison of image restoration methods.

## REFERENCES

[1] Y. Blau and T. Michaeli, "The perception-distortion tradeoff," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 6228–6237, doi: 10.1109/CVPR.2018.00652.

[2] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: From error visibility to structural similarity," *IEEE Transactions on Image Processing*, vol. 13, no. 4, pp. 600–612, 2004.

[3] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, and W. Shi, "Photo-realistic single image super-resolution using a generative adversarial network," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 4681–4690.

[4] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *European Conference on Computer Vision (ECCV)*, 2016, pp. 694–711.

[5] R. Dahl, M. Norouzi, and J. Shlens, "Pixel recursive super resolution," in *International Conference on Computer Vision (ICCV)*, 2017, pp. 5439–5448.

[6] M. S. M. Sajjadi, B. Scholkopf, and M. Hirsch, "EnhanceNet: Single image super-resolution through automated texture synthesis," in *International Conference on Computer Vision (ICCV)*, 2017, pp. 4491–4500.

[7] R. A. Yeh, C. Chen, T. Y. Lim, A. G. Schwing, M. Hasegawa-Johnson, and M. N. Do, "Semantic image inpainting with deep generative models," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 5485–5493.

[8] C.-Y. Yang, C. Ma, and M.-H. Yang, "Single-image super-resolution: A benchmark," in *European Conference on Computer Vision (ECCV)*, 2014, pp. 372–386.

[9] T. R. Shaham, T. Dekel, and T. Michaeli, "SinGAN: Learning a generative model from a single natural image," in *International Conference on Computer Vision (ICCV)*, 2019.

[10] Z. Wang and A. C. Bovik, "Mean squared error: Love it or leave it? A new look at signal fidelity measures," *IEEE Signal Processing Magazine*, vol. 26, no. 1, pp. 98–117, 2009.

[11] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 1125–1134.

[12] R. Zhang, P. Isola, and A. A. Efros, "Colorful image colorization," in *European Conference on Computer Vision (ECCV)*, 2016, pp. 649–666.

[13] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen, "Improved techniques for training gans," in *Advances in Neural Information Processing Systems (NIPS)*, 2016, pp. 2234–2242.

[14] E. L. Denton, S. Chintala, R. Fergus *et al.*, "Deep generative image models using a Laplacian pyramid of adversarial networks," in *Advances in Neural Information Processing Systems (NIPS)*, 2015, pp. 1486–1494.

[15] S. Iizuka, E. Simo-Serra, and H. Ishikawa, "Let there be color!: Joint end-to-end learning of global and local image priors for automatic image colorization with simultaneous classification," *ACM Transactions on Graphics (TOG)*, vol. 35, no. 4, p. 110, 2016.

[16] R. Zhang, J.-Y. Zhu, P. Isola, X. Geng, A. S. Lin, T. Yu, and A. A. Efros, "Real-time user-guided image colorization with learned deep priors," *ACM Transactions on Graphics (TOG)*, vol. 9, no. 4, 2017.

[17] S. Guadarrama, R. Dahl, D. Bieber, M. Norouzi, J. Shlens, and K. Murphy, "PixColor: Pixel recursive colorization," *British Machine Vision Conference (BMVC)*, 2017.

[18] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in Neural Information Processing Systems (NIPS)*, 2014, pp. 2672–2680.

[19] T. M. Cover and J. A. Thomas, *Elements of information theory*. John Wiley & Sons, 2012.

[20] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros, "Context encoders: Feature learning by inpainting," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 2536–2544.

[21] O. Rippel and L. Bourdev, "Real-time adaptive image compression," in *International Conference on Machine Learning (ICML)*, 2017, pp. 2922–2930.

[22] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *International Conference on Computer Vision (ICCV)*, 2017, pp. 2223–2232.

[23] Z. Wang, E. P. Simoncelli, and A. C. Bovik, "Multiscale structural similarity for image quality assessment," in *Conference on Signals, Systems and Computers*, vol. 2, 2003, pp. 1398–1402.

[24] H. R. Sheikh, A. C. Bovik, and G. De Veciana, "An information fidelity criterion for image quality assessment using natural scene statistics," *IEEE Transactions on Image Processing*, vol. 14, no. 12, pp. 2117–2128, 2005.

[25] H. R. Sheikh and A. C. Bovik, "Image information and visual quality," *IEEE Transactions on Image Processing*, vol. 15, no. 2, pp. 430–444, 2006.

[26] D. M. Chandler and S. S. Hemami, "VSNR: A wavelet-based visual signal-to-noise ratio for natural images," *IEEE Transactions on Image Processing*, vol. 16, no. 9, pp. 2284–2298, 2007.

[27] L. Zhang, L. Zhang, X. Mou, and D. Zhang, "FSIM: A feature similarity index for image quality assessment," *IEEE Transactions on Image Processing*, vol. 20, no. 8, pp. 2378–2386, 2011.

[28] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 586–595.

[29] A. K. Moorthy and A. C. Bovik, "Blind image quality assessment: From natural scene statistics to perceptual quality," *IEEE transactions on Image Processing*, vol. 20, no. 12, pp. 3350–3364, 2011.

[30] A. Mittal, A. K. Moorthy, and A. C. Bovik, "No-reference image quality assessment in the spatial domain," *IEEE Transactions on Image Processing*, vol. 21, no. 12, pp. 4695–4708, 2012.

[31] F. Nielsen, "Hypothesis testing, information divergence and computational geometry," in *Geometric Science of Information*, 2013, pp. 241–248.

[32] Z. Wang and E. P. Simoncelli, "Reduced-reference image quality assessment using a wavelet-domain natural image statistic model." in *Human Vision and Electronic Imaging*, vol. 5666, 2005, pp. 149–159.

[33] Z. Wang, G. Wu, H. R. Sheikh, E. P. Simoncelli, E.-H. Yang, and A. C. Bovik, "Quality-aware images," *IEEE Transactions on Image Processing*, vol. 15, no. 6, pp. 1680–1689, 2006.

[34] Q. Li and Z. Wang, "Reduced-reference image quality assessment using divisive normalization-based image representation," *IEEE Journal of Selected Topics in Signal Processing*, vol. 3, no. 2, pp. 202–211, 2009.

[35] M. A. Saad, A. C. Bovik, and C. Charrier, "Blind image quality assessment: A natural scene statistics approach in the DCT domain," *IEEE transactions on Image Processing*, vol. 21, no. 8, pp. 3339–3352, 2012.

[36] A. Mittal, R. Soundararajan, and A. C. Bovik, "Making a "completely blind" image quality analyzer," *IEEE Signal Processing Letters*, vol. 20, no. 3, pp. 209–212, 2013.

[37] X. Wang, K. Yu, S. Wu, J. Gu, Y. Liu, C. Dong, Y. Qiao, and C. Change Loy, "Esrgan: Enhanced super-resolution generative adversarial networks," in *Proceedings of the European Conference on Computer Vision (ECCV) workshops*, 2018.

[38] J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, and T. S. Huang, "Generative image inpainting with contextual attention," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 5505–5514.

[39] E. Agustsson, M. Tschannen, F. Mentzer, R. Timofte, and L. Van Gool, "Generative adversarial networks for extreme learned image compression," *arXiv preprint arXiv:1804.02958*, 2018.

[40] M. Tschannen, E. Agustsson, and M. Lucic, "Deep generative models for distribution-preserving lossy compression," in *Proc. Conference on Neural Information Processing Systems (NeurIPS)*, 2018.

[41] O. Kupyn, V. Budzan, M. Mykhailych, D. Mishkin, and J. Matas, "Deblurgan: Blind motion deblurring using conditional adversarial networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 8183–8192.

[42] M.-Y. Liu, T. Breuel, and J. Kautz, "Unsupervised image-to-image translation networks," in *Advances in neural information processing systems (NIPS)*, 2017, pp. 700–708.

[43] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein generative adversarial networks," in *International Conference on Machine Learning (ICML)*, 2017, pp. 214–223.

[44] S. Nowozin, B. Cseke, and R. Tomioka, "f-gan: Training generative neural samplers using variational divergence minimization," in *Advances in Neural Information Processing Systems (NIPS)*, 2016, pp. 271–279.

[45] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.

[46] M. Nikolova, "Model distortions in Bayesian MAP reconstruction," *Inverse Problems and Imaging*, vol. 1, no. 2, p. 399, 2007.

[47] I. Csiszár, P. C. Shields *et al.*, "Information theory and statistics: A tutorial," *Foundations and Trends® in Communications and Information Theory*, vol. 1, no. 4, pp. 417–528, 2004.

[48] T. Van Erven and P. Harremos, "Rényi divergence and Kullback-Leibler divergence," *IEEE Transactions on Information Theory*, vol. 60, no. 7, pp. 3797–3820, 2014.

[49] A. Levin and B. Nadler, "Natural image denoising: Optimality and inherent bounds," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2011, pp. 2833–2840.

[50] A. Levin, B. Nadler, F. Durand, and W. T. Freeman, "Patch complexity, finite pixel correlations and optimal denoising," in *European Conference on Computer Vision (ECCV)*, 2012, pp. 73–86.

[51] P. Chatterjee and P. Milanfar, "Is denoising dead?" *IEEE Transactions on Image Processing*, vol. 19, no. 4, pp. 895–911, 2010.

[52] P. Chatterjee and P. Milanfar, "Practical bounds on image denoising: From estimation to information," *IEEE Transactions on Image Processing (TIP)*, vol. 20, no. 5, pp. 1221–1233, 2011.

[53] S. Baker and T. Kanade, "Limits on super-resolution and how to break them," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 24, no. 9, pp. 1167–1183, 2002.

[54] S. Santurkar, D. Budden, and N. Shavit, "Generative compression," in *Proc. Picture Coding Symposium (PCS)*, 2018.

[55] Y. Blau and T. Michaeli, "Rethinking lossy compression: The Rate-distortion-perception tradeoff," in *International Conference on Machine Learning (ICML)*, 2019.

[56] R. Matsumoto, "Introducing the perception-distortion tradeoff into the rate-distortion theory of general information sources," *IEICE Communications Express*, vol. 7, no. 11, pp. 427–431, 2018.

[57] R. Matsumoto, "Rate-distortion-perception tradeoff of variable-length source coding for general information sources," *IEICE Communications Express*, vol. 8, no. 2, pp. 38–42, 2019.

[58] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville, "Improved training of wasserstein gans," in *Advances in Neural Information Processing Systems (NIPS)*, 2017, pp. 5769–5779.

[59] P. Ye, J. Kumar, L. Kang, and D. Doermann, "Unsupervised feature learning framework for no-reference image quality assessment," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012, pp. 1098–1105.

[60] L. Kang, P. Ye, Y. Li, and D. Doermann, "Convolutional neural networks for no-reference image quality assessment," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014, pp. 1733–1740.

[61] C. Ma, C.-Y. Yang, X. Yang, and M.-H. Yang, "Learning a no-reference quality metric for single-image super-resolution," *Computer Vision and Image Understanding*, vol. 158, pp. 1–16, 2017.

[62] W.-S. Lai, J.-B. Huang, Z. Hu, N. Ahuja, and M.-H. Yang, "A comparative study for single image blind deblurring," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 1701–1709.

[63] Y. Liu, J. Wang, S. Cho, A. Finkelstein, and S. Rusinkiewicz, "A no-reference metric for evaluating the quality of motion deblurring." *ACM Transactions on Graphics (TOG)*, vol. 32, no. 6, pp. 175–1, 2013.

[64] Y. Blau, R. Mechrez, R. Timofte, T. Michaeli, and L. Zelnik-Manor, "The 2018 pirm challenge on perceptual image super-resolution," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.

[65] D. Martin, C. Fowlkes, D. Tal, and J. Malik, "A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics," in *International Conference on Computer Vision (ICCV)*, vol. 2, 2001, pp. 416–423.

[66] R. Timofte, V. De Smet, and L. Van Gool, "A+: Adjusted anchored neighborhood regression for fast super-resolution," in *Asian Conference on Computer Vision*, 2014, pp. 111–126.

[67] C. Dong, C. C. Loy, K. He, and X. Tang, "Learning a deep convolutional network for image super-resolution," in *European Conference on Computer Vision (ECCV)*, 2014, pp. 184–199.

[68] J.-B. Huang, A. Singh, and N. Ahuja, "Single image super-resolution from transformed self-exemplars," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 5197–5206.

[69] J. Kim, J. Kwon Lee, and K. Mu Lee, "Accurate image super-resolution using very deep convolutional networks," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 1646–1654.

[70] W.-S. Lai, J.-B. Huang, N. Ahuja, and M.-H. Yang, "Deep Laplacian pyramid networks for fast and accurate super-resolution," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 624–632.

[71] W. Bae, J. Yoo, and J. Chul Ye, "Beyond deep residual learning for image restoration: Persistent homology-guided manifold simplification," in *Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2017, pp. 145–153.

[72] B. Lim, S. Son, H. Kim, S. Nah, and K. M. Lee, "Enhanced deep residual networks for single image super-resolution," in *Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2017.

[73] X. Deng, "Enhancing image quality via style transfer for single image super-resolution," *IEEE Signal Processing Letters*, 2018.

[74] R. Mechrez, I. Talmi, F. Shama, and L. Zelnik-Manor, "Maintaining natural image statistics with the contextual loss," in *Asian Conference on Computer Vision*. Springer, 2018, pp. 427–443.

[75] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *International Conference on Learning Representations (ICLR)*, 2014.

# APPENDIX A
## REAL-VS.-FAKE USER STUDIES AND HYPOTHESIS TESTING

We assume the setting where an observer is shown a real image (a draw from $p_X$) or an algorithm output (a draw from $p_{\hat{X}}$), with a prior probability of 0.5 each. The task is to identify which distribution the image was drawn from ($p_X$ or $p_{\hat{X}}$) with maximal probability of success. This is the setting of the Bayesian hypothesis testing problem, for which the maximum a-posteriori (MAP) decision rule minimizes the probability of error (see Section 1 in [31]). When there are two possible hypotheses with equal probabilities (as in our setting), the relation between the probability of error and the total-variation distance between $p_X$ and $p_{\hat{X}}$ in (1) can be easily derived (see Section 2 in [31]).

# APPENDIX B
## THE MMSE AND MAP EXAMPLES OF SEC. 3

Sections 3.1 and 3.2 exemplify that the MSE and the $0-1$ loss are not distribution preserving in the setting of estimating a discrete random variable (vector) $X$ from its noisy version $Y = X + N$, where $N \sim \mathcal{N}(0, \sigma^2 I)$ is independent of $X$. Since the conditional distribution of $Y$ given $X = x$ is $\mathcal{N}(x, \sigma^2 I)$, the MMSE estimator is given by

$$
\begin{aligned}
\hat{x}_{\text{MMSE}}(y) &= \mathbb{E}[X|Y = y] \\
&= \sum_x x p(x|y) \\
&= \sum_x x \frac{p(y|x)p(x)}{\sum_{x'} p(y|x')p(x')} \\
&= \sum_x x \frac{\exp(-\frac{1}{2\sigma^2}\|y-x\|^2)p(x)}{\sum_{x'} \exp(-\frac{1}{2\sigma^2}\|y-x'\|^2)p(x')}, \quad (18)
\end{aligned}
$$

and the MAP estimator is given by

$$
\begin{aligned}
\hat{x}_{\text{MAP}}(y) &= \arg\max_x p(x|y) \\
&= \arg\min_x -\log(p(y|x)p(x)) \\
&= \arg\min_x \frac{1}{2\sigma^2}\|y-x\|^2 - \log(p(x)). \quad (19)
\end{aligned}
$$

In the example of Fig. 4, $x$ is a $280 \times 280$ binary image comprising $28 \times 28$ blocks chosen uniformly at random from a finite database. Since the noise $N$ is i.i.d., each $28 \times 28$ block of $y$ can be denoised separately, both in the case of the MSE criterion and in the case of MAP. For each block, we have $p(x) = 1/59400$ for the non-blank images and $p(x) = 1/11$ for the blank image.

In the trinary example (5), we calculate the distribution of the MMSE estimate (Fig. 3) by

$$
p_{\hat{X}_{\text{MMSE}}}(\hat{x}) = p_Y(\hat{x}_{\text{MMSE}}^{-1}(\hat{x})) \left| \frac{d}{d\hat{x}} \hat{x}_{\text{MMSE}}^{-1}(\hat{x}) \right| \quad (20)
$$

where the inverse of $\hat{x}_{\text{MMSE}}(y)$ (see (6)) and its derivative are calculated numerically, and $p_Y(y) = \sum_x p(y|x)p(x)$ with $p(y|x) \sim \mathcal{N}(x, 1)$ and $p(x)$ of (5).

# APPENDIX C
## PROOF OF THEOREM 1

We will show that a stably distribution preserving optimal estimator is necessarily unique. At the same time, we will show that a non-invertible degradation implies that this optimal estimator is non-unique. Specifically, we use the following definitions.

**Definition 6.** *We say that a degradation is* not invertible *if $p_{X|Y}(x|y) > 0$ for all $(x, y) \in \mathcal{S}_x \times \mathcal{S}_y$, where $\mathcal{S}_x$ is a non-singleton set and $\mathcal{S}_y$ satisfies $\mathbb{P}(Y \in \mathcal{S}_y) > 0$.*

**Definition 7.** *We say that the optimal estimator is* not unique *if there exist two estimators, $p_{\hat{X}_1|Y}$ and $p_{\hat{X}_2|Y}$ that minimize the mean distortion (3) and differ from one another in the sense that*

$$
d_{TV}\left(p_{\hat{X}_1|Y}(\cdot|y), p_{\hat{X}_2|Y}(\cdot|y)\right) > 0 \quad \forall y \in \mathcal{S}_y \quad (21)
$$

*where $\mathcal{S}_y$ is a set that satisfies $\mathbb{P}(Y \in \mathcal{S}_y) > 0$.*

The outline of the proof of Theorem 1 will be as follows:

1) In Lemma 1 we will show that if the distortion measure is stably distribution preserving, then the optimal estimator $\hat{X}^*$ is uniquely defined by $p_{\hat{X}^*|Y} = p_{X|Y}$.
2) In Lemma 2 we will show that if the estimator $\hat{X}^*$ defined by $p_{\hat{X}^*|Y} = p_{X|Y}$ is an optimal estimator and the degradation is non-invertible, then the optimal estimator is non-unique.
3) This leads to a contradiction, proving that there does not exist a stably distribution preserving distortion metric if the degradation in non-invertible.

**Lemma 1.** *If the distortion measure $\Delta(\cdot, \cdot)$ is stably distribution preserving at $p_{X,Y}$, then the optimal estimator $\hat{X}^*$ that minimizes the mean distortion (3) is uniquely defined by $p_{\hat{X}^*|Y} = p_{X|Y}$.*

*Proof.* We start by noting that the optimal estimator $p_{\hat{X}|Y}$ depends only on $p_{X|Y}$ and not on $p_Y$. Indeed, since $X$ and $\hat{X}$ are independent given $Y$, the mean distortion can be written as

$$
\begin{aligned}
\mathbb{E}[\Delta(X, \hat{X})] &= \iiint \Delta(x, \hat{x}) p_{X|Y}(x|y) p_{\hat{X}|Y}(\hat{x}|y) p_Y(y) dx d\hat{x} dy \\
&= \int \left( \int f(\hat{x}, y) p_{\hat{X}|Y}(\hat{x}|y) d\hat{x} \right) p_Y(y) dy, \quad (22)
\end{aligned}
$$

where we defined

$$
f(\hat{x}, y) = \int \Delta(x, \hat{x}) p_{X|Y}(x|y) dx. \quad (23)
$$

Therefore, the optimal $p_{\hat{X}|Y}$ is that which minimizes $\int f(\hat{x}, y) p_{\hat{X}|Y}(\hat{x}|y) d\hat{x}$ for each $y$. Since $f(\hat{x}, y)$ depends only on $p_{X|Y}$, the optimal estimator depends only on $p_{X|Y}$.

Next, we observe that if a distortion measure is stably distribution preserving at $p_{X,Y}$, then there exists an $\alpha \in (0, 1)$ such that the measure is distribution preserving at any perturbed joint distribution of the form $\tilde{p}_{X,Y} = p_{X|Y}\tilde{p}_Y$, where

$$
\tilde{p}_Y = \alpha p_Y + (1 - \alpha)q \quad (24)
$$

and $q$ is any distribution. That is, we take a perturbed joint distribution having the same posterior $p_{X|Y}$ as $p_{X,Y}$, but

a perturbed marginal. Indeed, taking $\alpha \geq 1 - \varepsilon$, any such $\tilde{p}_{X,Y}$ is in the TV $\varepsilon$-ball around $p_{X,Y}$, as

$$
\begin{aligned}
d_{TV}(p_{X,Y}, \tilde{p}_{X,Y}) &= \tfrac{1}{2} \iint |p_{X,Y}(x,y) - \tilde{p}_{X,Y}(x,y)| dx dy \\
&= \tfrac{1}{2} \iint |p_{X|Y}(x|y)p_Y(y) - p_{X|Y}(x|y)\tilde{p}_Y(y)| dx dy \\
&= \tfrac{1}{2}(1-\alpha) \iint |p_{X|Y}(x|y)p_Y(y) - p_{X|Y}(x|y)q(y)| dx dy \\
&\leq 1 - \alpha \\
&\leq \varepsilon. \quad (25)
\end{aligned}
$$

By our assumption that the optimal estimator is stably distribution preserving, it must satisfy $p_{\hat{X}^*} = p_X$ for any perturbation of $p_{X,Y}$ of the form (24). Since the posterior has not changed, the optimal estimator $p_{\hat{X}^*|Y}$ remains the same. Its marginal $\tilde{p}_{\hat{X}^*}$, however, is modified to

$$
\begin{aligned}
\tilde{p}_{\hat{X}^*}(x) &= \int p_{\hat{X}^*|Y}(x|y)\tilde{p}_Y(y) dy \\
&= \alpha \int p_{\hat{X}^*|Y}(x|y)p_Y(y) dy + (1-\alpha) \int p_{\hat{X}^*|Y}(x|y)q(y) dy \\
&= \alpha p_X(x) + (1-\alpha) \int p_{\hat{X}^*|Y}(x|y)q(y) dy, \quad (26)
\end{aligned}
$$

where we used the assumption that $p_{\hat{X}^*} = p_X$. Similarly, the distribution of $X$ has changed to

$$
\begin{aligned}
\tilde{p}_X(x) &= \int p_{X|Y}(x|y)\tilde{p}_Y(y) dy \\
&= \alpha \int p_{X|Y}(x|y)p_Y(y) dy + (1-\alpha) \int p_{X|Y}(x|y)q(y) dy \\
&= \alpha p_X(x) + (1-\alpha) \int p_{X|Y}(x|y)q(y) dy. \quad (27)
\end{aligned}
$$

Thus, equality between $\tilde{p}_{\hat{X}^*}$ and $\tilde{p}_X$ is kept only if

$$
\int p_{\hat{X}^*|Y}(x|y)q(y) dy = \int p_{X|Y}(x|y)q(y) dy. \quad (28)
$$

This equality can hold for *every* perturbation $q$ only if $p_{\hat{X}^*|Y} = p_{X|Y}$, completing the proof.

Notice that this also proves that the optimal estimator is unique (under the stably distribution preserving assumption), as we demonstrated that only $p_{\hat{X}^*|Y} = p_{X|Y}$ minimizes the mean distortion. $\square$

**Lemma 2.** *If the degradation is non-invertible, and the estimator $\hat{X}^*$ defined by $p_{\hat{X}^*|Y} = p_{X|Y}$ is an optimal estimator, then the optimal estimator is non-unique.*

*Proof.* Since the degradation is non-invertible, $p_{X|Y}(x|y) > 0$ for all $(x,y) \in \mathcal{S}_x \times \mathcal{S}_y$, where $\mathcal{S}_x$ is a non-singleton set and $\mathcal{S}_y$ is a set that satisfies $\mathbb{P}(Y \in \mathcal{S}_y) > 0$ (Definition 6). As $p_{\hat{X}^*|Y} = p_{X|Y}$, we also have that $p_{\hat{X}^*|Y}(x|y) > 0$ for all $(x,y) \in \mathcal{S}_x \times \mathcal{S}_y$.

Now, since $\hat{X}^*$ is an optimal estimator, $p_{\hat{X}^*|Y}$ must minimize $\int f(\hat{x},y)p_{\hat{X}^*|Y}(\hat{x}|y) d\hat{x}$ for each $y$ (see proof of Lemma 1). This means that for any $y$, the conditional $p_{\hat{X}^*|Y}(\hat{x}|y)$ must assign positive probability only to $\hat{x}$ in the set of minima $\mathcal{S}_{\min}(y) = \arg\min_{\hat{x}} f(\hat{x},y)$. We conclude that $\mathcal{S}_x \subseteq \mathcal{S}_{\min}(y)$ for every $y \in \mathcal{S}_y$. This implies that any other estimator that assigns zero probability to $\hat{x} \notin \mathcal{S}_x$ for every $y \in \mathcal{S}_y$, is also optimal.

Let $\mathcal{S}_x^1, \mathcal{S}_x^2$ be non-empty disjoint sets such that $\mathcal{S}_x^1 \cup \mathcal{S}_x^2 = \mathcal{S}_x$. Now, define two estimators, such that $p_{\hat{X}_1|Y}(\hat{x}|y) > 0$ only for $\hat{x} \in \mathcal{S}_x^1$, and $p_{\hat{X}_2|Y}(\hat{x}|y) > 0$ only for $\hat{x} \in \mathcal{S}_x^2$, for every $y \in \mathcal{S}_y$. Both are optimal estimators (as they only assign positive probability to $\hat{x} \in \mathcal{S}_x$). Yet, these two estimators have conditional distributions with disjoint supports for every $y$, and thus $d_{TV}(p_{\hat{X}_1|Y}(\cdot|y), p_{\hat{X}_2|Y}(\cdot|y)) > 0 \quad \forall y \in \mathcal{S}_y$. Therefore, by Definition 7, the optimal estimator is non-unique. $\square$

Now, let us assume to the contrary that $p_{X,Y}$ defines a non-invertible degradation, and that the distortion function $\Delta(\cdot,\cdot)$ is stably distribution preserving at $p_{X,Y}$. By Lemma 1, the optimal estimator $\hat{X}^*$ is uniquely defined by $p_{\hat{X}^*|Y} = p_{X|Y}$. But now according to Lemma 2, since the degradation is non-invertible, the optimal estimator is non-unique, leading to a contradiction.

## APPENDIX D
## DERIVATION OF EXAMPLE 1

Since $\hat{X} = aY = a(X + N)$, it is a zero-mean Gaussian random variable. Now, the Kullback-Leibler distance between two zero-mean normal distributions is given by

$$
d_{KL}(p_X \| p_{\hat{X}}) = \ln\left(\frac{\sigma_{\hat{X}}}{\sigma_X}\right) + \frac{\sigma_X^2}{2\sigma_{\hat{X}}^2} - \frac{1}{2}, \quad (29)
$$

and the MSE between $X$ and $\hat{X}$ is given by

$$
\text{MSE}(X, \hat{X}) = E[(X - \hat{X})^2] = \sigma_X^2 - 2\sigma_{X\hat{X}} + \sigma_{\hat{X}}^2. \quad (30)
$$

Substituting $\hat{X} = aY$ and $\sigma_X^2 = 1$, we obtain that $\sigma_{\hat{X}} = |a|\sqrt{1 + \sigma_N^2}$ and $\sigma_{X\hat{X}} = a$, so that

$$
d_{KL}(a) = \ln\left(|a|\sqrt{1 + \sigma_N^2}\right) + \frac{1}{2a^2(1 + \sigma_N^2)} - \frac{1}{2}, \quad (31)
$$

$$
\text{MSE}(a) = 1 + a^2(1 + \sigma_N^2) - 2a, \quad (32)
$$

and

$$
P(D) = \min_a d_{KL}(a) \quad \text{s.t.} \quad \text{MSE}(a) \leq D. \quad (33)
$$

Notice that $d_{KL}$ is symmetric, and $\text{MSE}(|a|) \leq \text{MSE}(a)$ (see Fig. 12). Thus, for any negative $a$, there always exists a positive $a$ with which $d_{KL}$ is the same and the MSE is not larger. Therefore, without loss of generality, we focus on the range $a \geq 0$.

For $D < D_{\min} = \frac{\sigma_N^2}{1 + \sigma_N^2}$ the constraint set of $\text{MSE}(a) < D$ is empty, and there is no solution to (33). For $D \geq D_{\min}$, the constraint is satisfied for $a_- \leq a \leq a_+$, where

$$
a_\pm(D) = \frac{1}{(1 + \sigma_N^2)}\left(1 \pm \sqrt{D(1 + \sigma_N^2) - \sigma_N^2}\right). \quad (34)
$$

For $D = D_{\min}$, the optimal (and only possible) $a$ is

$$
a = a_+(D_{\min}) = a_-(D_{\min}) = \frac{1}{(1 + \sigma_N^2)}. \quad (35)
$$

For $D > D_{\min}$, $a_+$ monotonically increases with $D$, broadening the constraint set. The objective $d_{KL}(a)$ monotonically decreases with $a$ in the range $a \in (0, 1/\sqrt{(1 + \sigma_N^2)})$ (see Fig. 12 and the mathematical justification below). Thus, for $D_{\min} < D \leq D_0$, the optimal $a$ is always the largest
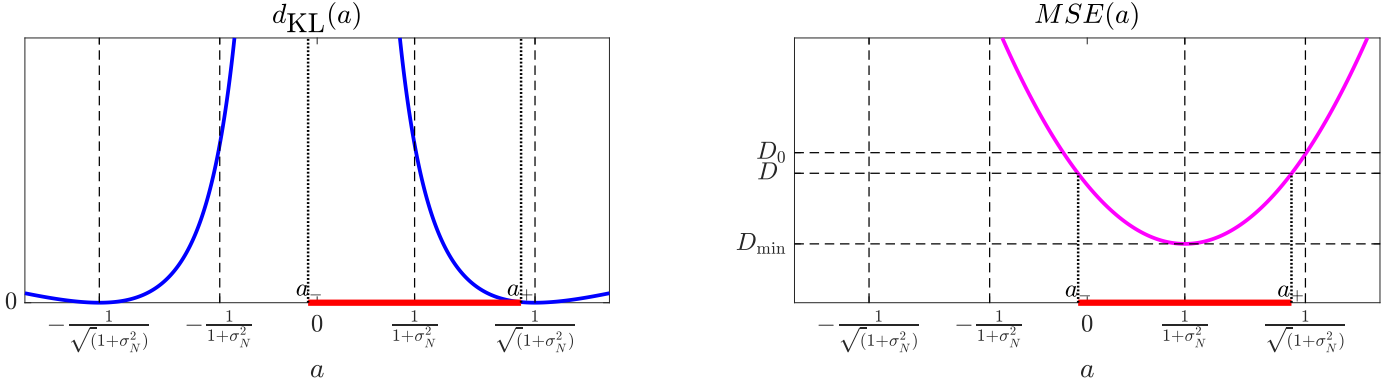
Fig. 12. Plots of (31) and (32). $D$ defines the range $(a_-, a_+)$ of $a$ values complying with the MSE constraint (marked in red). The objective $d_{\mathsf{KL}}$ is minimized over this range of possible $a$ values.
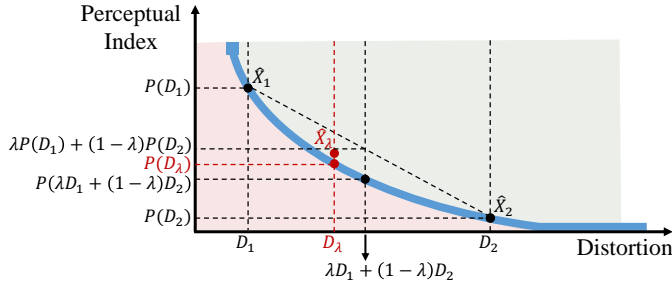


Fig. 13. Illustration of the proof of Theorem 2.

possible $a$, which is $a = a_+(D)$, where $D_0$ is defined by $a_+(D_0) = 1/\sqrt{(1+\sigma_N^2)}$ (see Fig. 12). For $D > D_0$, the optimal $a$ is $a = 1/\sqrt{(1+\sigma_N^2)}$, which achieves the global minimum $d_{\mathsf{KL}}(a) = 0$. The closed form solution is therefore given by

$$P(D) = \begin{cases} d_{\mathsf{KL}}(a_+(D)) & D_{\min} \leq D < D_0 \\ 0 & D_0 \leq D \end{cases} \tag{36}$$

To justify the monotonicity of $d_{\mathsf{KL}}(a)$ in the range $a \in (0, 1/\sqrt{(1+\sigma_N^2)})$, notice that for $a > 0$,

$$\frac{d}{da} d_{\mathsf{KL}}(a) = \frac{1}{a} - \frac{1}{(1+\sigma_N^2)} \frac{1}{a^3}, \tag{37}$$

which is negative for $a \in (0, 1/\sqrt{(1+\sigma_N^2)})$.

## APPENDIX E
## PROOF OF THEOREM 2

The proof of Theorem 2 follows closely that of the rate-distortion theorem from information theory [19]. The value $P(D)$ is the minimal distance $d(p_X, p_{\hat{X}})$ over a constraint set whose size does not decrease with $D$. This implies that the function $P(D)$ is non-increasing in $D$. Now, to prove the convexity of $P(D)$, we will show that

$$\lambda P(D_1) + (1-\lambda)P(D_2) \geq P(\lambda D_1 + (1-\lambda)D_2), \tag{38}$$

for all $\lambda \in [0, 1]$ (see Fig. 13). First, by definition, the left hand side of (38) can be written as

$$\lambda d(p_X, p_{\hat{X}_1}) + (1-\lambda)d(p_X, p_{\hat{X}_2}), \tag{39}$$

where $\hat{X}_1$ and $\hat{X}_2$ are the estimators defined by

$$p_{\hat{X}_1|Y} = \arg\min_{p_{\hat{X}|Y}} d(p_X, p_{\hat{X}}) \text{ s.t. } \mathbb{E}\left[\Delta(X, \hat{X})\right] \leq D_1, \tag{40}$$

$$p_{\hat{X}_2|Y} = \arg\min_{p_{\hat{X}|Y}} d(p_X, p_{\hat{X}}) \text{ s.t. } \mathbb{E}\left[\Delta(X, \hat{X})\right] \leq D_2. \tag{41}$$

Since $d(\cdot, \cdot)$ is convex in its second argument,

$$\lambda d(p_X, p_{\hat{X}_1}) + (1-\lambda)d(p_X, p_{\hat{X}_2}) \geq d(p_X, p_{\hat{X}_\lambda}), \tag{42}$$

where $\hat{X}_\lambda$ is defined by

$$p_{\hat{X}_\lambda|Y} = \lambda p_{\hat{X}_1|Y} + (1-\lambda)p_{\hat{X}_2|Y}. \tag{43}$$

Denoting $D_\lambda = \mathbb{E}[\Delta(X, \hat{X}_\lambda)]$, we have that

$$d(p_X, p_{\hat{X}_\lambda}) \geq \min_{p_{\hat{X}|Y}} \left\{ d(p_X, p_{\hat{X}}) : \mathbb{E}[\Delta(X, \hat{X})] \leq D_\lambda \right\}$$
$$= P(D_\lambda), \tag{44}$$

because $\hat{X}_\lambda$ is in the constraint set. Below, we show that

$$D_\lambda \leq \lambda D_1 + (1-\lambda)D_2. \tag{45}$$

Therefore, since $P(D)$ is non-increasing in $D$, we have that

$$P(D_\lambda) \geq P(\lambda D_1 + (1-\lambda)D_2). \tag{46}$$

Combining (39), (42), (44) and (46) proves (38), thus demonstrating that $P(D)$ is convex.

To justify (45), note that

$$D_\lambda = \mathbb{E}\left[\Delta(X, \hat{X}_\lambda)\right]$$
$$= \mathbb{E}\left[\mathbb{E}\left[\Delta(X, \hat{X}_\lambda)|Y\right]\right]$$
$$= \mathbb{E}\left[\lambda\mathbb{E}\left[\Delta(X, \hat{X}_1)|Y\right] + (1-\lambda)\mathbb{E}\left[\Delta(X, \hat{X}_2)|Y\right]\right]$$
$$= \lambda\mathbb{E}\left[\Delta(X, \hat{X}_1)\right] + (1-\lambda)\mathbb{E}\left[\Delta(X, \hat{X}_2)\right]$$
$$\leq \lambda D_1 + (1-\lambda)D_2, \tag{47}$$

where the second and fourth transitions are according to the law of total expectation and the third transition is justified by

$$p(x, \hat{x}_\lambda | y) = p(\hat{x}_\lambda | x, y) p(x | y)$$
$$= p(\hat{x}_\lambda | y) p(x | y)$$
$$= (\lambda p(\hat{x}_1 | y) + (1 - \lambda) p(\hat{x}_2 | y)) p(x | y)$$
$$= \lambda p(\hat{x}_1 | y) p(x | y) + (1 - \lambda) p(\hat{x}_2 | y)) p(x | y)$$
$$= \lambda p(x, \hat{x}_1 | y) + (1 - \lambda) p(x, \hat{x}_2 | y)). \tag{48}$$

Here we used (43) and the fact that $X$ and $\hat{X}_\lambda$ are independent given $Y$, and similarly for the pairs $(X, \hat{X}_1)$ and $(X, \hat{X}_2)$.

## APPENDIX F
## PROOF OF THEOREM 3

The estimator $\hat{X}$ of (14) attains perfect perceptual quality since

$$p_{\hat{X}}(x) = \int p_{\hat{X}|Y}(x|y) p_Y(y) dy$$
$$= \int p_{X|Y}(x|y) p_Y(y) dy$$
$$= p_X(x). \tag{49}$$

Furthermore, note that

$$\mathbb{E}[X^T \hat{X}] = \mathbb{E}[\mathbb{E}[X^T \hat{X} | Y]]$$
$$= \mathbb{E}[\mathbb{E}[X | Y]^T \mathbb{E}[\hat{X} | Y]]$$
$$= \mathbb{E}[\|\mathbb{E}[X | Y]\|^2], \tag{50}$$

and

$$\mathbb{E}[\|\hat{X}\|^2] = \mathbb{E}[\mathbb{E}[\|\hat{X}\|^2 | Y] = \mathbb{E}[\mathbb{E}[\|X\|^2 | Y] = \mathbb{E}[\|X\|^2], \tag{51}$$

where we used the law of total expectation and the fact that given $Y$, $X$ and $\hat{X}$ are independent and identically distributed. The MSE of $\hat{X}$ is therefore

$$\mathbb{E}[\|X - \hat{X}\|^2] = \mathbb{E}[\|X\|^2] - 2\mathbb{E}[X^T \hat{X}] + \mathbb{E}[\|\hat{X}\|^2]$$
$$= 2(\mathbb{E}[\|X\|^2] - \mathbb{E}[\|\mathbb{E}[X | Y]\|^2])$$
$$= 2\mathbb{E}[\|X - \mathbb{E}[X | Y]\|^2]$$
$$= 2\mathbb{E}[\|X - \hat{X}_{\text{MMSE}}\|^2], \tag{52}$$

where the second equality is due to (50) and (51), and the third equality is due to the orthogonality principle. We thus established that $\hat{X}$ is a distribution preserving estimator whose MSE is precisely twice the MSE of the MMSE estimator. This implies that

$$D_{\max} \le \mathbb{E}[\|X - \hat{X}\|^2] = 2D_{\min}, \tag{53}$$

completing the proof.

## APPENDIX G
## WGAN ARCHITECTURE AND TRAINING DETAILS (SEC. 5)

The architecture of the WGAN trained for denoising the MNIST images is detailed in Table 1. The training algorithm and adversarial losses are as proposed in [58]. The generator loss was modified to include a content loss term,

i.e. $\ell_{\text{gen}} = \ell_{\text{MSE}} + \lambda \ell_{\text{adv}}$, where $\ell_{\text{MSE}}$ is the standard MSE loss. For each $\lambda$ the WGAN was trained for 35 epochs, with a batch size of 64 images. The ADAM optimizer [75] was used, with $\beta_1 = 0.5$, $\beta_2 = 0.9$. The generator/discriminator initial learning rate is $10^{-3}/10^{-4}$ respectively, where learning rate of both decreases by half every 10 epochs. The filter size of the discriminator convolutional layers is $5 \times 5$, and these are performed without padding. The filter size in the generator transposed-convolutional layers is $5 \times 5/4 \times 4$, and these are performed with $2/1$ pixel padding for the first/ second and third transposed-convolutional layers, respectively. The stride of each convolutional layer and the slope for the leaky-ReLU layers appear in Table 1. Note that the perception-distortion curve in Fig. 7 is generated by training on single digit images, which in general may deviate from the perception-distortion curve of whole images containing i.i.d. sub-blocks of digits.

## APPENDIX H
## SUPER-RESOLUTION EVALUATION DETAILS (SEC. 6) AND ADDITIONAL COMPARISONS

The no-reference (NR) and full-reference (FR) methods BRISQUE, BLIINDS-II, NIQE, SSIM, MS-SSIM, IFC and VIF were obtained from the LIVE laboratory website[4], the NR method of Ma et al. was obtained from the project webpage[5], and the pretrained VGG-19 network was obtained through the PyTorch torchvision package[6]. The low-resolution images were obtained by factor 4 downsampling with a bicubic kernel. The super-resolution results on the BSD100 dataset of the SRGAN and SRResNet variants were obtained online[7], and the results of EDSR, Deng, Johnson et al. and Mechrez et al. were kindly provided by the authors. The algorithms for testing the other SR methods were obtained online: A+[8], SRCNN[9], SelfEx[10], VDSR[11], LapSRN[12], Bae et al. [13] and ENet[14]. All NR and FR metrics and all SR algorithms were used with the default parameters and models. In the paper, we reported comparisons on the y-channel (except for the $\text{VGG}_{2,2}$ measure). In the supplementary material, we report results with additional NR metrics on the y-channel, as well as results on color images. When comparing color images, for SR algorithms which treat the y-channel alone, the Cb and Cr channels are upsampled by bicubic interpolation.

The general pattern appearing in Fig. 9 will appear for any NR method which accurately predicts the perceptual quality of images. We show here three additional popular NR methods: BRISQUE [30], BLIINDS-II [35] and the recent measure by Ma et al. [61] in Figs. 14, 15, 16, where the same conclusions as for NIQE [36] (see Sec. 6) are apparent. The

4. http://live.ece.utexas.edu/research/Quality/index.htm
5. https://github.com/chaoma99/sr-metric
6. http://pytorch.org/docs/master/torchvision/index.html
7. https://twitter.box.com/s/lcue6vlrd01ljkdtdkhmfvk7vtjhetog
8. http://www.vision.ee.ethz.ch/~timofter/ACCV2014_ID820_SUPPLEMENTARY/
9. http://mmlab.ie.cuhk.edu.hk/projects/SRCNN.html
10. https://github.com/jbhuang0604/SelfExSR
11. http://cv.snu.ac.kr/research/VDSR/
12. https://github.com/phoenix104104/LapSRN
13. https://github.com/iorism/CNN
14. https://webdav.tue.mpg.de/pixel/enhancenet/

TABLE 1
Generator and discriminator architecture. FC is a fully-connected layer, BN is a batch-norm layer, and l-ReLU is a leaky-ReLU layer.

| Discriminator | |
|---|---|
| Size | Layer |
| $28 \times 28 \times 1$ | Input |
| $12 \times 12 \times 32$ | Conv (stride=2), l-ReLU (slope=0.2) |
| $4 \times 4 \times 64$ | Conv (stride=2), l-ReLU (slope=0.2) |
| 1024 | Flatten |
| 1 | FC |
| 1 | Output |

| Generator | |
|---|---|
| Size | Layer |
| $28 \times 28 \times 1$ | Input |
| 784 | Flatten |
| $4 \times 4 \times 128$ | FC, unflatten, BN, ReLU |
| $7 \times 7 \times 64$ | transposed-Conv (stride=2), BN, ReLU |
| $14 \times 14 \times 32$ | transposed-Conv (stride=2), BN, ReLU |
| $28 \times 28 \times 1$ | transposed-Conv (stride=2), sigmoid |
| $28 \times 28 \times 1$ | Output |

same pattern appears for RGB images as well, as shown in Figs. 17, 18. Note that the perceptual quality of Johnson et al. and SRResNet-VGG$_{2,2}$ is inconsistent between NR metrics, likely due to varying sensitivity to the cross-hatch pattern artifacts which are present in these method's outputs. For this reason, Johnson et al. does not appear in the NIQE plots, as its NIQE score is 13.55 (far off the plots).

Fig. 14. Plot of 15 algorithms on the perception-distortion plane, where perception is measured by the NR metric by Ma et al. [61], and distortion is measured by the common full-reference metrics RMSE, SSIM, MS-SSIM, IFC, VIF and VGG$_{2,2}$. All metrics were calculated on the **y-channel** alone.



Fig. 15. Plot of 16 algorithms on the perception-distortion plane, where perception is measured by the NR metric BRISQUE, and distortion is measured by the common full-reference metrics RMSE, SSIM, MS-SSIM, IFC, VIF and VGG$_{2,2}$. All metrics were calculated on the **y-channel** alone.
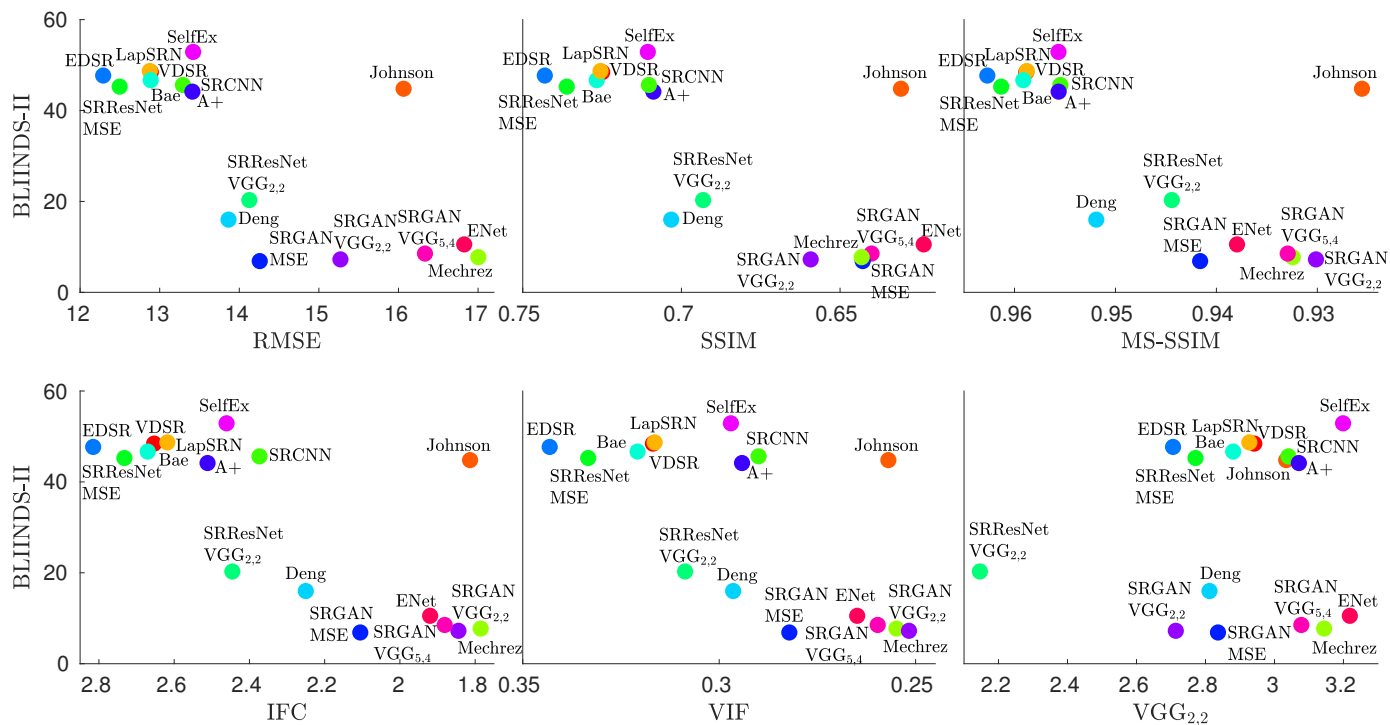
Fig. 16. Plot of 16 algorithms on the perception-distortion plane, where perception is measured by the NR metric BLIINDS-II, and distortion is measured by the common full-reference metrics RMSE, SSIM, MS-SSIM, IFC, VIF and VGG$_{2,2}$. All metrics were calculated on the **y-channel** alone.
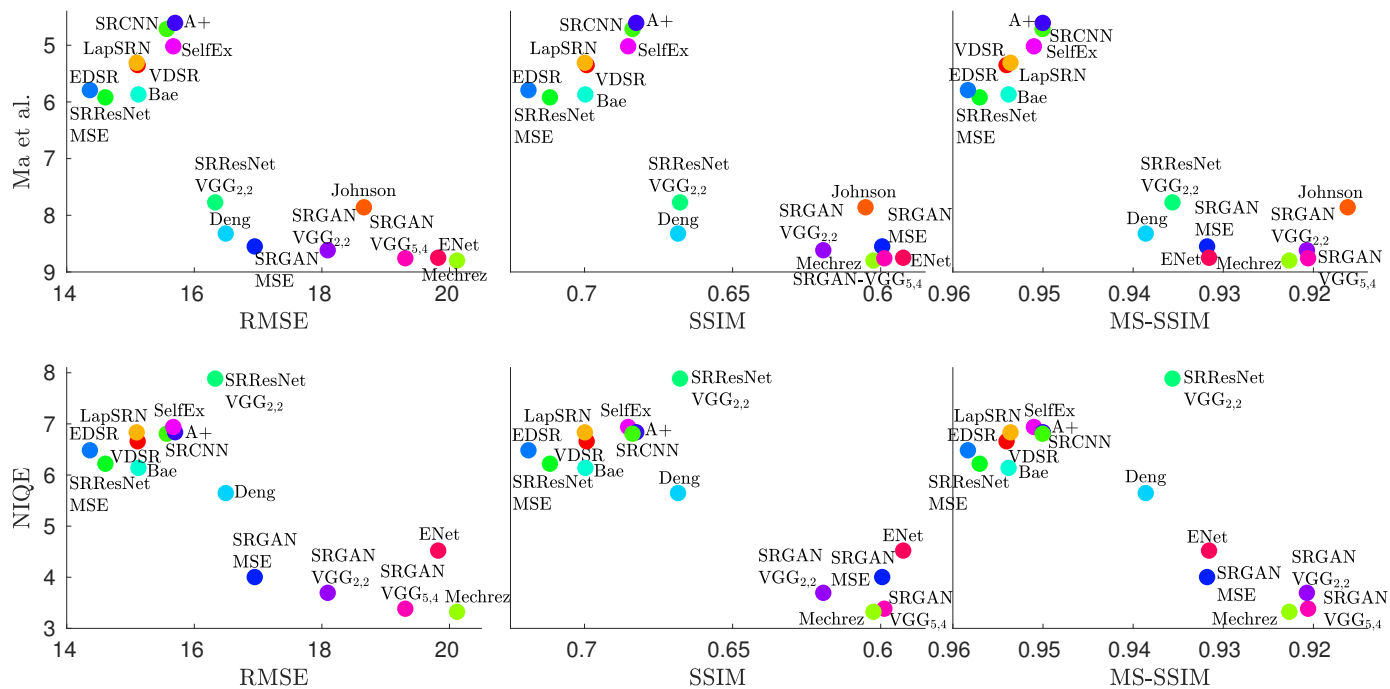


Fig. 17. Plot of 16 algorithms on the perception-distortion plane. Perception is measured by the the NR metrics of Ma et al. and NIQE, and distortion is measured by the common full-reference metrics RMSE, SSIM and MS-SSIM. All metrics were calculated on **three channel RGB** images.
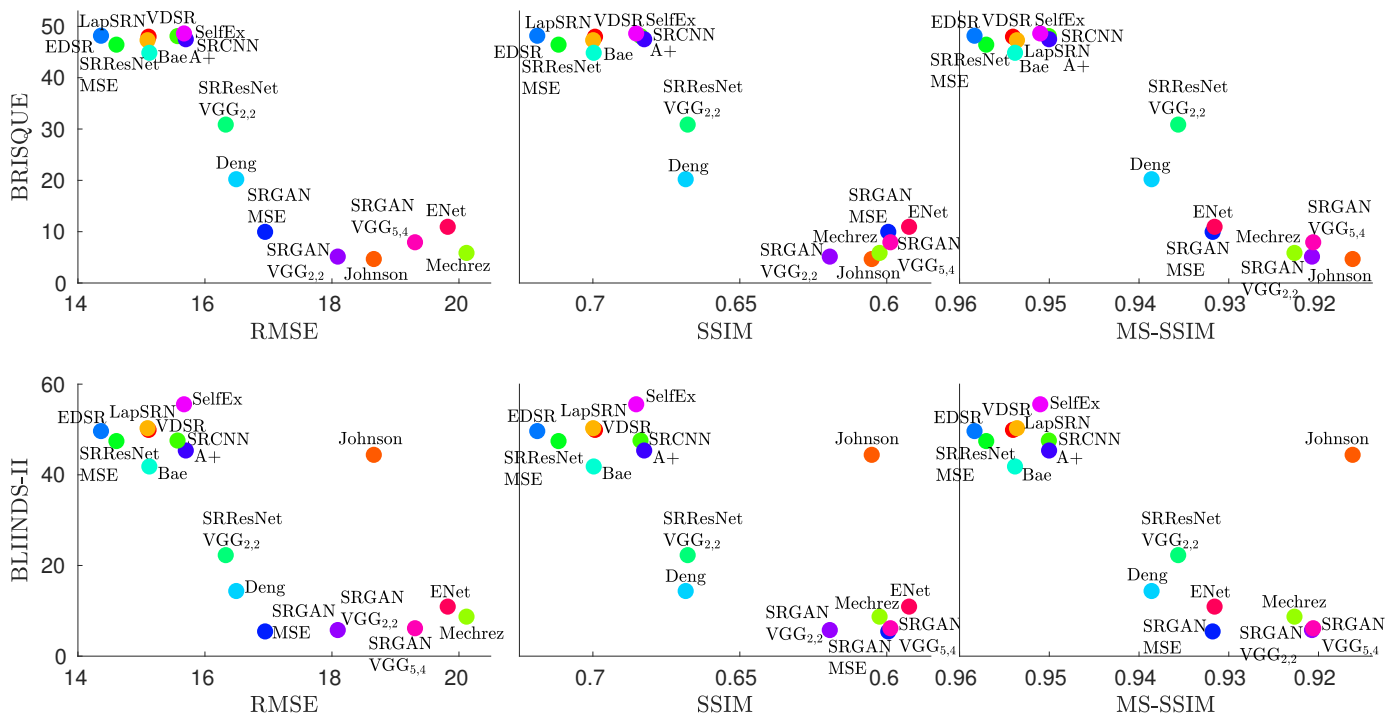
Fig. 18. Plot of 16 algorithms on the perception-distortion plane. Perception is measured by the the NR metrics BRISQUE and BLIINDS-II, and distortion is measured by the common full-reference metrics RMSE, SSIM and MS-SSIM. All metrics were calculated on **three channel RGB** images.