# Understanding SSIM

Jim Nilsson
NVIDIA

Tomas Akenine-Möller
NVIDIA
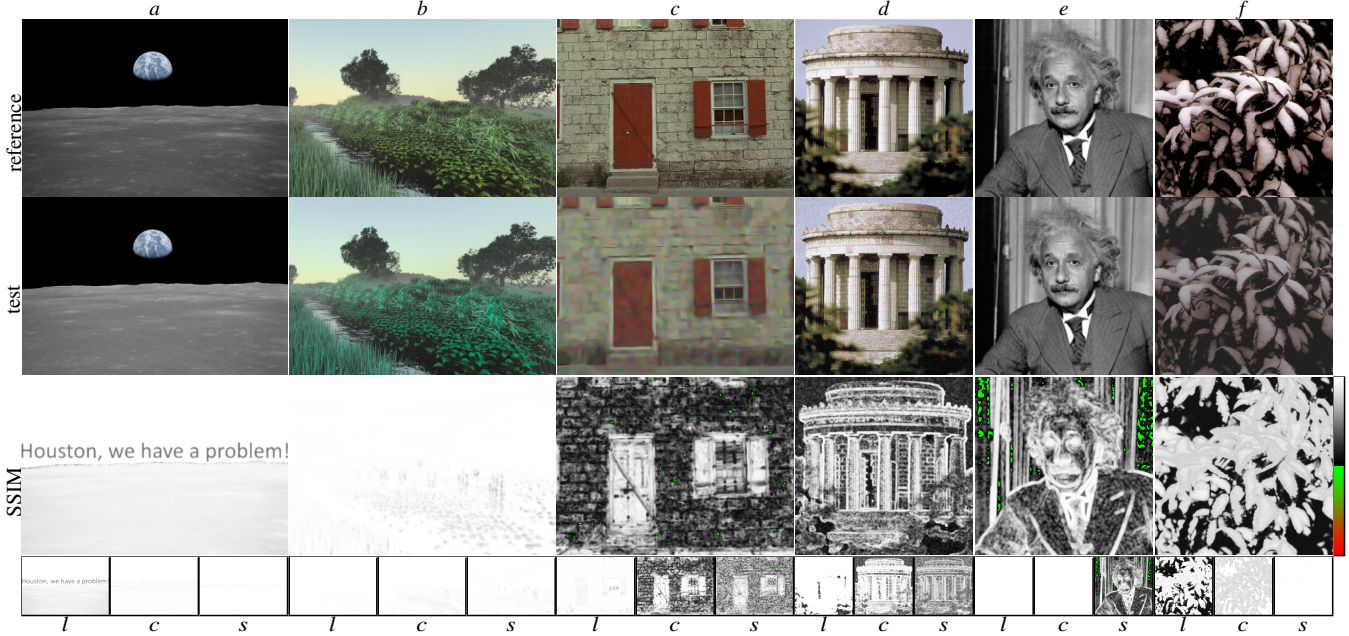
Figure 1: Images should ideally be viewed on a display at 100% scale, so we urge the reader to look at the images in our supplemental material.* The SSIM images (3rd row) are generated from the references (1st row) and the test images (2nd row), and SSIM is visualized using the heatmap to the lower right, where white is SSIM = 1 (identical images), black is SSIM = 0, and SSIM values in $(-1, 0]$ map to (red, green). The 4th row contains the $l$, $c$, and $s$ components of SSIM. Comments: *a)* the dark text is detected, *b)* chrominance differences are hardly noticed, *c)* large areas on the red door and window shutter are not detected, *d)* edges in images often generate an SSIM value close to one (white), even though the error is large all over this test image, *e)* the Einstein images are similar unless highly zoomed-in, but SSIM reports large errors, and *f)* the final pair where the SSIM image is leaning toward being almost inverted compared to what one may expect. Note that all reference and test images, but not SSIM images, are reduced in size for space considerations.

## Abstract

The use of the structural similarity index (SSIM) is widespread. For almost two decades, it has played a major role in image quality assessment in many different research disciplines. Clearly, its merits are indisputable in the research community. However, little deep scrutiny of this index has been performed. Contrary to popular belief, there are some interesting properties of SSIM that merit such scrutiny. In this paper, we analyze the mathematical factors of SSIM and show that it can generate results, in both synthetic and realistic use cases, that are unexpected, sometimes undefined, and nonintuitive. As a consequence, assessing image quality based on SSIM can lead to incorrect conclusions and using SSIM as a loss function for deep learning can guide neural network training in the wrong direction.

## 1 Introduction

The original SSIM paper [19] has over 20,000 citations on Google Scholar. Thousands of research papers have used it as a quality index when comparing images, and we are indeed authors of a few of those. In this paper, we provide a review and a deep inspection of SSIM, and we show that SSIM can deliver unexpected or invalid results in both simple use cases and for real image pairs. See Figure 1. We start with an overview of SSIM.

The input color space of SSIM is never defined. As the reference Matlab script performs no color space transformations on inputs, our assumption throughout this paper is that all images are encoded in sRGB color space, i.e., approximately gamma encoded with an exponent $\approx 2.4$. Note that this means that an image that is loaded by the SSIM script is assumed to be viewed directly on screen as is. For two images **A** and **B**, the original formula [19] for per-pixel SSIM is given by

$$\text{SSIM}(x, y) = \left(l(x, y)\right)^{\alpha} \left(c(x, y)\right)^{\beta} \left(s(x, y)\right)^{\gamma}, \quad (1)$$

where **A** and **B** are inputs to all functions, but omitted for clarity. To compute mean, variance, and covariance in a patch around a pixel, they use Gaussian-weighted versions of these formulae with a filter kernel of $11 \times 11$ pixels and $\sigma = 1.5$. The *luminance* component, $l$, is then

$$l(x, y) = \frac{2\mu_{\mathbf{A}}\mu_{\mathbf{B}} + C_1}{\mu_{\mathbf{A}}^2 + \mu_{\mathbf{B}}^2 + C_1}, \quad (2)$$

---

where the mean values $\mu_{\mathbf{A}}$ and $\mu_{\mathbf{B}}$ are functions of $x, y$ as well, e.g., $\mu_{\mathbf{A}}(x, y)$, but we skip the $(x, y)$ in favor of shorter notation, and likewise for the variances $\sigma_{\mathbf{A}}^2$ and $\sigma_{\mathbf{B}}^2$, and for the covariance $\sigma_{\mathbf{AB}}$. The *contrast* component, $c$, is

$$c(x, y) = \frac{2\sigma_{\mathbf{A}}\sigma_{\mathbf{B}} + C_2}{\sigma_{\mathbf{A}}^2 + \sigma_{\mathbf{B}}^2 + C_2}. \tag{3}$$

Finally, the *structure* component, $s$, is

$$s(x, y) = \frac{\sigma_{\mathbf{AB}} + C_3}{\sigma_{\mathbf{A}}\sigma_{\mathbf{B}} + C_3}. \tag{4}$$

Wang et al. propose that $C_1 = (K_1 L)^2$, $C_2 = (K_2 L)^2$, and $C_3 = C_2/2$, where $L = 255$ for 8-bit component images. Furthermore, they chose $K_1 = 0.01$ and $K_2 = 0.03$. If the range for an image is $[0, 1]$, we set $L = 1$ in order to get the same result, i.e., $C_1 = K_1^2$ and $C_2 = K_2^2$. Finally, the mean SSIM (MSSIM) value, which is pooled over the entire image, is

$$\mathrm{MSSIM}(\mathbf{A}, \mathbf{B}) = \frac{1}{wh} \sum_x \sum_y \mathrm{SSIM}(x, y), \tag{5}$$

where $w$ and $h$ are the width and height of the image. As can be seen, the term $\sigma_{\mathbf{A}}\sigma_{\mathbf{B}}$ is in the numerator in Equation 3 and in the denominator in Equation 4. To create a simplified expression, Wang et al. therefore proposed to use $\alpha = \beta = \gamma = 1$ and $C_3 = C_2/2$, which results in

$$\mathrm{SSIM}(x, y) = \frac{(2\mu_{\mathbf{A}}\mu_{\mathbf{B}} + C_1)(2\sigma_{\mathbf{AB}} + C_2)}{(\mu_{\mathbf{A}}^2 + \mu_{\mathbf{B}}^2 + C_1)(\sigma_{\mathbf{A}}^2 + \sigma_{\mathbf{B}}^2 + C_2)}. \tag{6}$$

Next, we review some related work.

## 2 The History of SSIM

The universal quality index (UQI), which was introduced by Wang and Bovik [17], is essentially SSIM as presented above, but without any of the constants $C_i$. These constants were added later to avoid division by zero [19]. Multi-scale SSIM (MS-SSIM)[1] was introduced as a means for including image details at different scales [21]. MS-SSIM adds more components in the expression, where both the contrast and structure expressions are evaluated at five low-pass filtered and downsampled versions of the original images. However, the components have the same form as in SSIM.

Sampat et al. [16] introduce complex wavelet structural similarity (CW-SSIM), where the expression in Equation 6 is used, but where the components are replaced by complex wavelet coefficients. CW-SSIM is more tolerant to small translations and rotations, which may be a desired effect in some contexts. However, for rendered images, which often contain geometrical edges, it is most likely not a desired feature, since, for instance, a game designer usually wants the geometry to be precisely where he/she intends. 3D-SSIM [24] is an extension of SSIM for video, where the formulae are evaluated for three-dimensional blocks of pixel values and multiplied with information content weights and local distortion weights. SSIM is still being adapted for new uses, e.g., spherical SSIM [6], where SSIM was adapted to handle a spherical projection, and for medical images [14].

---

[1]Note the difference between MSSIM, which is the average SSIM value over an image pair, and MS-SSIM, which is the multi-scale variant of SSIM.

While mean square error (MSE) has been criticized [18] for not delivering a truthful value compared to image error, Dosselman and Yang [8] and Horé and Ziou [10] have, at the same time, shown that there is a close relationship between MSE and SSIM. Whittle et al. [22] evaluate image metrics for Monte Carlo rendered images with different levels of noise. Their conclusion is that MS-SSIM performs well for this task. Čadík et al. [5] perform an extensive evaluation of image indices and metrics together with a user study, and find contradictory results for several of the algorithms, including SSIM, for various image distortions. Recently, SSIM has found uses as a loss function for deep learning [25], and is also included in tool kits such as Tensorflow.

Neither the UQI nor SSIM make any claim to be a *metric* in the mathematical sense, for which the triangle inequality, i.e., $d(x, z) \leq d(x, y) + d(y, z)$, must hold. It has, however, been shown that $\sqrt{1 - l(x, y)}$ and $\sqrt{1 - c(x, y)s(x, y)}$ do fulfill the triangle inequality, and thus are metrics [3]. Interestingly, the derivation to transform a modified version of SSIM into a metric, also revealed subtle—while important—properties of the index. For instance, for $l(x, y)$, we have:

$$\sqrt{1 - l(x, y)} = \frac{|\mu_{\mathbf{A}} - \mu_{\mathbf{B}}|}{\sqrt{\mu_{\mathbf{A}}^2 + \mu_{\mathbf{B}}^2 + C_1}}, \tag{7}$$

which can be seen as a normalized version of the root mean square error (RMSE).

This is notable, since MSE is not a perception-based metric [9, 18] and the finding above has the implication that there could exist a direct relationship between SSIM and MSE, which has indeed been independently discovered by Dosselman and Yang [7, 8] and Horé and Ziou [10]. Specifically, Dosselman and Yang show that there is a direct mathematical transform between $\mathrm{SSIM}^*$ and $\mathrm{MSE}^*$. $\mathrm{SSIM}^*$ is SSIM with constants $C_i = 0$, which does not alter the validity of the analysis, while $\mathrm{MSE}^*$ is a local MSE, using the same footprint as SSIM. They further show this empirically by correlating MSE versus SSIM for a range of images, using the coefficient of multiple determination, $R^2$, which is $0.0$ for no association between the variables and values closer to $1.0$ indicate strong degree of correspondence [7]. It was found that $R^2$ was between $0.9322$ and $1.0$, which implies that $\mathrm{SSIM}^*$ and $\mathrm{MSE}^*$ perform similarly.

Furthermore, Horé and Ziou [10] make a similar discovery, i.e., that there is a close relationship between $\mathrm{SSIM}^*$, peak-signal-to-noise-ratio (PSNR), and MSE. They identify that $\mathrm{MSE} = \sigma_{\mathbf{A}}^2 + \sigma_{\mathbf{B}}^2 - 2\sigma_{\mathbf{AB}} + (\mu_{\mathbf{A}} - \mu_{\mathbf{B}})^2$, and can then derive PSNR as a function of $\mathrm{SSIM}^*$. For images of similar luminance, i.e., $\mu_{\mathbf{A}} \approx \mu_{\mathbf{B}}$, and SSIM values in the $[0.2, 0.8]$ range, this function is approximately linear, which indicates that for any other distortion than a luminance shift, $\mathrm{SSIM}^*$ is qualitatively equivalent to PSNR.

These findings question the validity of claims that SSIM is a perception-based index, since MSE is not a perception-based metric [9, 18]. The small discrepancies in correlation between MSE and SSIM were shown to stem partly from the fact that SSIM is derived for a spatial subregion of the whole images, and an effect of the $C_i$ constants [7, 8].
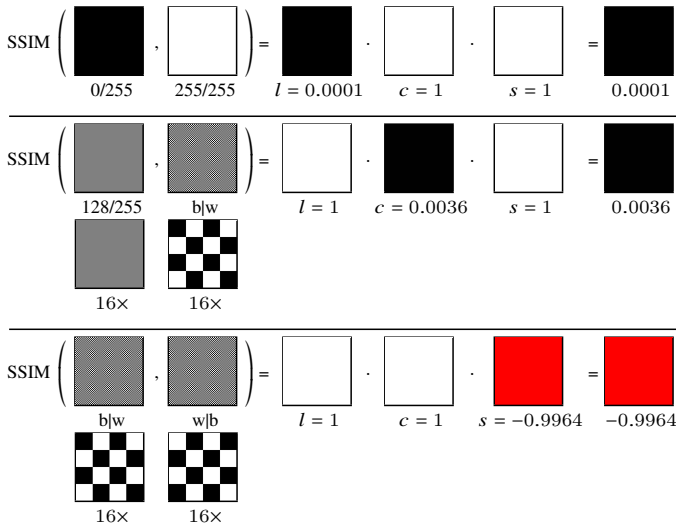
Figure 2: Minimum values of the SSIM components: $l_{\min} = 0.0001$ (top), $c_{\min} = 0.0036$ (middle), and $s_{\min} = -0.9964$ (bottom). The images at the bottom in the middle and bottom examples have been zoomed by a factor 16×.

## 3 Mathematical Properties

This section will analyze the components of SSIM from a mathematical standpoint. The behavior of the quality index itself will be scrutinized in the next section.

We use the notation $x/y$ to denote a pixel of a particular grayscale value. For example, $128/255$ corresponds to a gray value of 50.2%. The shorthand form $x|y$ designates a pixel-sized checkerboard pattern with colors $x/255$ and $y/255$. Letters $b$ and $w$ are shorthand for black ($0/255$) and white ($255/255$), respectively.

### 3.1 Minimum Values of the SSIM Factors

To understand the workings of the components of SSIM, as we will see later in this section, and since it, to our knowledge, has not been done before, we explain how to minimize $l$, $c$, and $s$, one at a time. For $l(x, y)$, shown in Equation 2, we can differentiate and solve for zero, with the assumptions that $\mu_A, \mu_B \in [0, L]$. This gives us a minimum when $\mu_A = 0$ and $\mu_B = L$ (or vice versa). Hence, we have $l_{\min} = C_1/(L^2 + C_1) = K_1^2/(K_1^2 + 1) = 0.0001$, for $K_1 = 0.01$.

The $c$ and $s$ components are slightly more involved and first we note that maximum variance and covariance for variables in $[0, L]$ is $(L/2)^2$. For the $c$ component (Equation 3), we can differentiate in the same manner as for $l$, and find that the minimum occurs when $\sigma_A = 0$ and $\sigma_B = (L/2)^2$ (or vice versa), which gives $c_{\min} = C_2/(L^2/4 + C_2) = K_2^2/(K_2^2 + 0.25) = 0.0036$, for $K_2 = 0.03$. Brunet [2] has ascertained that $s \in [-1, 1]$, but one can do slightly better than that, since $s$ (Equation 4) is minimized for minimum covariance $\sigma_{AB}$ and maximum variances ($\sigma_A^2, \sigma_B^2$). This gives the minimum at $s_{\min} = (-L^2/4 + C_3)/(L^2/4 + C_3) = (2K_2^2 - 1)/(2K_2^2 + 1) = -0.9964$.

To give an example that the above minima can occur, we show examples of when SSIM becomes $l_{\min}$, $c_{\min}$, and $s_{\min}$ in Figure 2. The first row uses a black and a white image, which results in $l = l_{\min}$ and $\sigma_A = \sigma_B = \sigma_{AB} = 0$, implying that $c = s = 1$, which gives SSIM=$l_{\min}$. The second row uses a $128/255$ image and a $b|w$ checkerboard, which thus has $\mu_A = \mu_B$, while at the same time making $\sigma_A = 0$ and $\sigma_B = (L/2)^2$. As a result,

SSIM=$c_{\min}$. Looking at the second row of Figure 2 with a low dot pitch, it is hard for the human visual system (HVS) to discern differences between the images, while SSIM values are close to zero, which indicates low quality contrary to the actual experience. The last row in Figure 2 achieves SSIM=$s_{\min}$ by using two inverted checkerboard images. Note that all three minima are independent of $L$, as expected, and it is clear that the range of SSIM is $(-1, 1]$.

While most people use the simplified version (Equation 6) of SSIM, which is also what the reference implementation [19] uses, the constants $\alpha$, $\beta$, and $\gamma$ have been used to find a more optimized SSIM expression [15] for antialiasing detection in games, but also in MS-SSIM [21], and for optimizing SSIM parameters using machine learning [4]. Therefore, it is important to take a look at the components in the full SSIM expression (Equation 1) as well and see what their ranges are. The $s$ component (Equation 4) deserves additional attention. If $C_3$ is zero, $s(x, y)$ is equivalent to the sample Pearson correlation coefficient, usually denoted $r_{xy}$. We have found no scientific support that this coefficient correlates with human perception of "structure." As mentioned in the introduction, SSIM has selected $C_3 = C_2/2 = (K_2 L)^2/2 = (0.03 \cdot 1)^2/2 = 0.00045$, for pixel values in the range $[0, 1]$. Since by definition we have $\sigma_A \geq 0$ and $\sigma_B \geq 0$, the denominator will always be positive. The covariance term $\sigma_{AB}$ can take on negative values, which means that $s(x, y)$ can be negative. Note that raising a negative number to a positive, non-integer number, $\gamma$, results in a complex number (with a real and an imaginary part), which in practice, e.g., in programming languages, makes the result *undefined*. The `std::pow()` function gives NaN (not a number) when a negative number is raised to a number (even to one). MS-SSIM [21] uses non-integer $\gamma$-values, as do the work of Čadík et al. [4] and Piórkowski & Mantiuk [15], so it seems that this is not well-known.

Undefined results (or complex numbers), unless properly defined, should not be an outcome of an image quality index. Even if the range of SSIM is allowed to be complex, there are no descriptions on how to interpret such values. To our knowledge, this problem has never been identified before and means that SSIM implementations can generate undefined results for some inputs and parameter settings.

### 3.2 Perceptual Properties

The purpose of deriving these minima is not only out of curiosity, but also hints at a deeper problem with the index itself and the claims of it being based on perception. Referring to the two bottom examples in Figure 2, it could be argued that detecting the difference between the images leading to a minimum value for $c$ and $s$ is hard (please look at the images in the supplemental material). Depending on the monitor dot pitch and viewing distance, any difference between the images $128/255$ and $b|w$ can be indistinguishable to a human viewer. The same is true for the bottom row, images $b|w$ and $w|b$. This property will be further investigated in Section 4.

As far as we can see, the only perception related claims made in the SSIM paper [19] is that the $l$ component is qualitatively consistent with Weber's law and that the $c$ component is consistent with the contrast-masking feature of the HVS. This is somewhat contradictory, since in the original UQI paper, the authors explicitly state that "the new index is mathematically
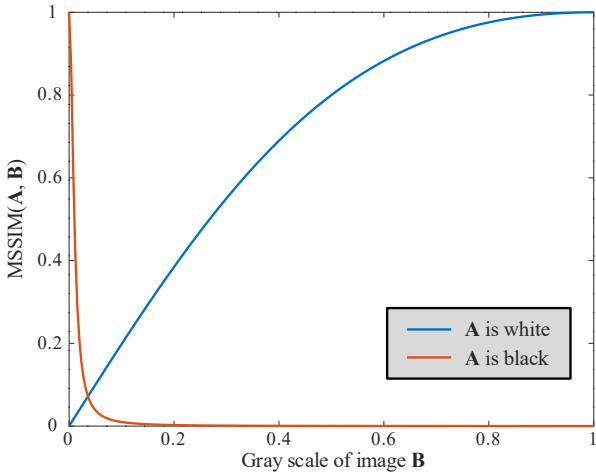
Figure 3: MSSIM as a function of an increasingly brighter, constant-valued image, **B**, against a black image (orange curve) and a white image (blue curve).
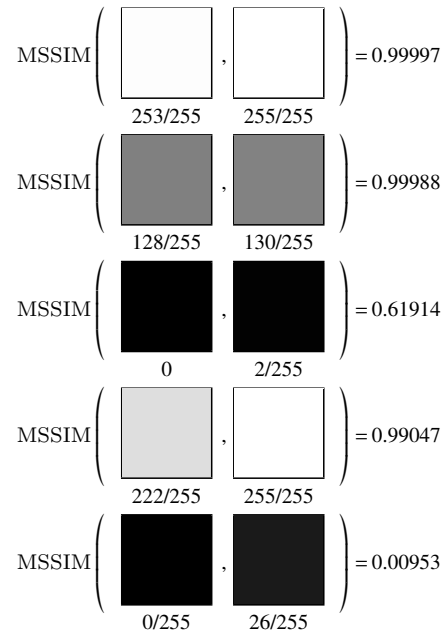


Figure 4: MSSIM behavior for constant grayscale images. The grayscale value is reported below each image. Note that in the third row, it is difficult to see any differences between the images, and still, MSSIM values indicate low quality. In the fourth row, it is clearly easy to see a difference between the images, while MSSIM indicate that they are similar. At the bottom, we see that when comparing a black image to an image with 26/255, MSSIM is almost zero.

defined and no human visual system model is explicitly employed." Weber's law [11] states that the ratio between the difference in stimulus against a *background* signal, to achieve the same psychophysical sensation, is approximately constant ($\frac{\Delta S}{S} = k$). Contrast-masking is the destructive interference between (transient) stimuli closely coupled in space and time [13]. Both these psychophysical phenomena are well known. However, these effects are both defined only *within the same frame of reference*, i.e., when introducing stimuli onto some sort of background. Consequently, in the context of image comparisons, these effects hold true for a local pixel value against its background *within* the same image, and not true for variations *between* images.

The claim that the $l$ factor is perceptually motivated is problematic. Evidence for this can easily be found and first, we point to Figure 4 to get a feel for this. To explain these results, we refer to Figure 3, which shows a plot of SSIM as a function of a constant colored image from black to white, against both a black and a white image. Note that, as before, $c = s = 1$. As can be seen, when the comparison is against a black image, **A**, and the image **B** is close to black, a minuscule change in **B** triggers a huge difference in $l$, and thus in SSIM. Furthermore, when **A** is black and **B** > 0.2, $l$ is always close to zero. The other case, when **A** is white, is not as radical, but for nearly white images **B**, relatively large changes in **B** do not change the SSIM value much. Section 4.1 reveals several interesting, nonintuitive results based on this diagram, and as a consequence, shows that the $l$ component is not perceptually based.

## 4 Evaluation

All results showing MSSIM values and the images of SSIM index maps were computed using the Matlab script (`www.cns.nyu.edu/~lcv/ssim/`) of Wang et al. All images are made available as supplemental material, as well as the scripts to generate our results. Since images ideally should be viewed on a display at 100% scale, instead of in a PDF viewer or on paper, we urge the reader to look at the images in our supplemental material. SSIM is visualized using a heatmap where white is SSIM = 1 (identical images), black is SSIM = 0, and SSIM values in $(-1, 0]$ map to (red, green].

## 4.1 Luminance

SSIM was designed so that $\mathrm{MSSIM}(\mathbf{A}, \mathbf{A}) = 1$, that is, if the images are the same, MSSIM will be one, and in general, a value $\geq 0.99$ indicates that the images are indistinguishable. Here, we will explore how MSSIM behaves for images that only contain a single grayscale value. For such images, $c = 1$ and $s = 1$, and therefore, it is only the $l$ component that affects the values in this experiment. In Figure 4 we reveal some results that have previously been unknown (to the best of our knowledge). The first row compares a white (255/255) image against a nearly-white (253/255) image, and MSSIM is almost one, which makes sense, since it is hard to see any difference between these two images. The difference in grayscale values is $2/255$. On the second row, we do the same but for mid-gray images with difference $2/255$ and the result is similar. However, the third row compares a black image to a nearly-black (2/255) image, again with a difference of $2/255$, and in this case, MSSIM values are low, indicating that the images are *not* similar, when, in fact, it is difficult to see any difference between the two.

The fourth row is, perhaps, even more surprising, since MSSIM values indicate that the images are similar, while they visually are not. In the fifth row, the MSSIM values indicate that these two images are dissimilar, and when increasing the value from 26 up to 255, the MSSIM value decreases to 0.0. This means that for 90% of the range, MSSIM is close to zero, which unreasonably compresses the resolution of the index. Mathematically, this stems from the quadratic forms in the normalization denominator ($\mu_{\mathbf{A}}^2 + \mu_{\mathbf{B}}^2 + C_1$) of Equation 2, which exaggerates differences near black. Regardless, SSIM does not seem to be aligned with the HVS's ability to detect luminance differences. The orange and blue curves in Figure 3 predict the results shown in Figure 4, confirming that $l$ component of SSIM can be misleading.
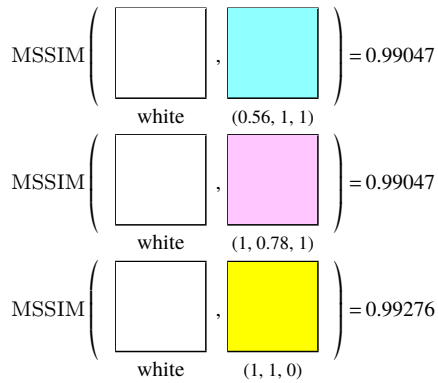
Figure 5: Color images and SSIM do not work well together, as shown here. We have reduced the red (top), green (middle), and blue (bottom) color components until the error became ≈ 0.99. Numbers in parentheses represent RGB. Clearly, the differences are visible even though SSIM values indicate that the image pairs are extremely similar.
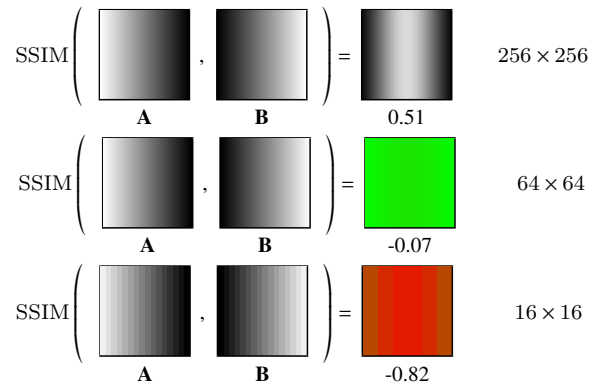


Figure 6: SSIM images of gradients versus mirrored gradients at different resolutions, shown to the right. Grayscale indicates positive SSIM, while green and red indicate negative (the same heatmap is used here as in Figure 1). The numbers under the images are the MSSIM values. Note that the original SSIM implementation, uses only "valid" values for the filter, which means that it removes a border of 5 pixels around the SSIM image. Hence, the bottom SSIM image is only $16 - 5 - 5 = 6$ pixels wide, for example.

## 4.2 Color

The SSIM authors present results using only JPEG and JPEG2000 color images [19], but add that using other color components does not significantly change the performance of the model. The index has still, nevertheless, been used on color images by first converting to grayscale values, using, for example, the color encoding standard Rec. 601 (which is used in Matlab's `rgb2gray`, and is the procedure recommended on the SSIM website):

$$Y = 0.2989r + 0.5810g + 0.1140b. \qquad (8)$$

Another approach is to convert from $RGB$ to $YC_rC_b$, and apply SSIM to $Y$, $C_r$, and $C_b$ individually, and then use $0.8\text{SSIM}_Y + 0.1\text{SSIM}_{C_r} + 0.1\text{SSIM}_{C_b}$ [20] as the quality index.

It is well-known that there are many colors of equal luminance and furthermore that any mapping between $RGB$ and grayscale value is many-to-one. As a consequence, SSIM can generate a high value, indicating similarity, even though the colors are visibly dissimilar. This is visualized for three color pairs in Figure 5, where the original `rgb2gray` function in Matlab has been used. It is evident that simply converting from RGB to grayscale values can give erroneous SSIM results, as would any metric relying on a reduction function from color to grayscale. Even the original SSIM paper [19] does this, and based on our findings here, we advise not to use SSIM with color images. A better solution would be to use a metric that inherently handles color.

## 4.3 Gradients

Next, we compare an image containing a gradient against a horizontally mirrored version of the same image at different resolutions. The results are summarized in Figure 6. In all examples, $c = 1$ and the image of $l$ (not shown) contains vertical lines of constant values, starting with a low value at the left, peaking in the middle, and going down to a low value at the right. The $s$ component, which is not shown in the figure, starts at $0.86$ for the pair with $256 \times 256$ pixels, but goes down to $-0.10$ for the middle row ($64 \times 64$), and even further down to $-0.90$ for $16 \times 16$ pixels, which means that it is the $s$ component that makes MSSIM values in Figure 6 negative toward the bottom. This was surprising since the **A** images are similar at all resolutions,

as are the **B** images. Assuming the $256 \times 256$ and the $16 \times 16$ gradients are part of a large image, the perceived error in the $16 \times 16$ region will surely be less glaring than in the $256 \times 256$ region. This is the opposite of the results, generated by SSIM, shown in Figure 6.

Recall that a negative $s$ to the power of a non-integer value generates a complex number or an undefined result (see the last part of Section 3.1), and this problem will occur for the two bottom rows in Figure 6. These examples might seem overly contrived, but serve to demonstrate that it is indeed possible, and not highly unlikely, for SSIM to generate negative values for simple distortions (see also Figure 1c, 1e, and 1f, which are discussed below).

## 4.4 Complex Images

The average MSSIM value has been correlated with the subjective quality (five levels, from "Bad" to "Excellent") of JPEG-compressed images [19]. The resulting SSIM image has however, to our knowledge, never been evaluated with subjective observers. It is certainly true that SSIM in many cases finds image differences that are also detected by a human observer. Had it not, its use would have been less widespread. With well-behaved image pairs, the corresponding SSIM map tells a convincing story. Synthetic examples (seen above, for example) serve to clarify situations where SSIM behaves contradictory to human perception, and it can be challenged whether these situations arise for more general images with reasonable distortions. In the following, we will demonstrate several such cases. All observations have been carried out on a Dell UP3216Q monitor, calibrated for sRGB (IEC 61966-2-1:1999 standard), with a D65 standard illuminant.

We start by discussing the images in Figure 1 in more detail. In image pair $a$, the surface of the moon has substantially different luminance, but that error is detected as lower than the text "Houston, we have a problem!", which is difficult to see for most people. As we have seen earlier, this can be explained using Figure 3. The test image in $b$ has its chrominance shifted compared the reference, and it is clear that SSIM does not react much to this, and the explanation to this is given in Section 4.2. Image pair $c$ is from the LocVis [23] database and shows a situ-
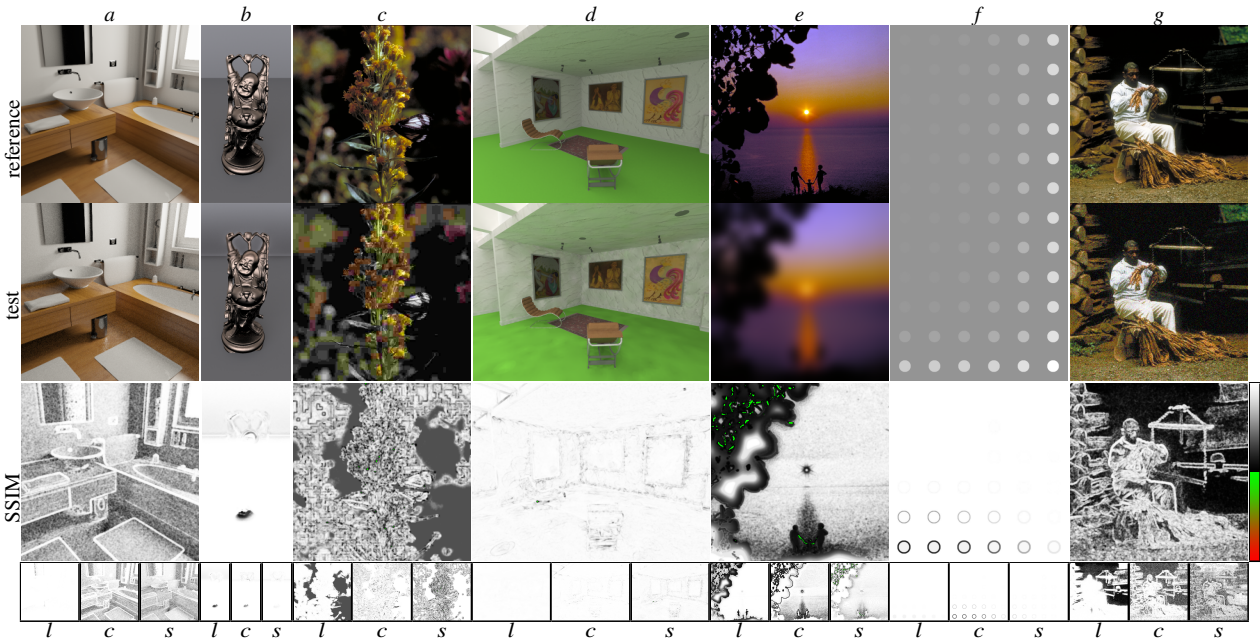
Figure 7: Images should ideally be viewed on a display at 100% scale, so we urge the reader to look at the images in our supplemental material. The same heatmap is used here as in Figure 1. All image pairs are discussed in Section 4.4. Note that all reference and test images, but not SSIM images, are reduced in size for space considerations.



Figure 8: Zooming in on the trash can in Figure 7a, we see that SSIM values tends to have very high values, i.e., low error, in a relatively large vicinity of contrast edges, which is unexpected since the Monte Carlo noise is of the same magnitude on both sides.

ation where the SSIM values are high on the red door and on the red window shutters, while it is clear that there are differences. The explanation is that after grayscale conversion, the grayscale values on the door and window shutters are the same in the reference and test images.

Image pair *d*, from the CS-IQ database [12], shows that SSIM sometimes can generate high SSIM values on edges in images, even though the error is evenly spread out over the entire image. As can be seen below the SSIM images, the *l*, *c*, and *s* components are all white on edges, so none of the terms detect the induced error. In image pair *e*, we have introduced a dithering that in a checkerboard pattern adds and subtracts 6 from the pixel (8-bit) grayscale value, clamped to 0 and 255, respectively. The test image uses the same, but inverse, dithering. To the human observer, we claim that these distortions have little, or no, visual impact. SSIM, however, finds these two images very dissimilar with many pixels generating negative (green) SSIM values. As seen in Section 3.1, a negative value to the power of a floating-point exponent is undefined in most programming languages or otherwise generates a complex number. Again, this is difficult to interpret in an image quality measure. The images in Figure 1*f*, with lowered contrast compared to reference, also from the CS-IQ database, show an

SSIM image that is particularly counterintuitive—in the bright regions, where visible differences can be argued to be largest, the SSIM image is bright as well, while in the dark regions of the reference and test image, SSIM values are much lower. As can be seen below the SSIM image, it is mostly the *l* component that makes this so, which has been explained in Section 4.1.

Turning to Figure 7, image pair *a* is a path traced rendering with different sample counts per pixel, for which the SSIM values at first seems to correspond well with the actual image differences. Closer inspection reveals two important exceptions, namely, high similarity along all high contrast edges, and an unreasonable dissimilarity for darker regions, mainly attributed to the *c* component. In Figure 8, we have zoomed in on the trash can in these images, to further illustrate how large the region around an edge is, with very high SSIM values. Image pair *b*, from LocVis, with different lighting conditions for reference and test, is noteworthy, since it demonstrates both a false positive (spot under the Buddha, nearly invisible) and a false negative (failure to detect luminance shift in the upper background) in the same image. The JPEG-compressed test image in *c*, from CS-IQ, is highly distorted in the detailed regions. A slight difference in gray level in the dark regions between the reference and test images, however, dominates the final output. Additionally, a few negative values show up for the *s* component.

Another example of what we consider is a false negative, is the blotchy distortion of image *d* (from LocVis), which most notable on the green floor. SSIM fails to detect this disturbing variation in intensity. Turning to the heavily blurred test image in *e*, from CS-IQ, the different components of SSIM again highlights dark regions where the images only differ by a small amount. The large drop in intensity of the sun is nearly ignored and the indication of high similarity in its middle is exaggerated. For the synthetic image pair in *f*, from LocVis, the *c* component shows dissimilarity around the edges of the circles, while the *l* component, containing information about the intensity levels within the circles, shows high similarity. We consider the lower

left circle in the test image to be significantly different from the reference image, but it does not show up much inside the circle. The final image pair *g*, from CS-IQ, is reported as strikingly dissimilar by SSIM across all dark regions, which is indicated by both the *l* and *c* components. The additive noise distortion is spread evenly across the whole test image, but in our opinion, these errors are harder to detect in darker areas, which seems to be the inverse behavior of SSIM.

## 5 Conclusions

We have demonstrated the mathematical properties of SSIM and shown that it is not adhering to properties of the human visual system. This is not surprising, as it was not a goal stated in the original work [17] upon which SSIM is based. However, over time, the similarity index has grown in number of uses, and popular belief of the index's capabilities has significantly widened with respect to this original scope. The purpose of this paper is to moderate this belief, since it can guide research in the wrong direction. Even though SSIM generates useful results in some cases, it can generate counterintuitive results in many others, as we have seen. This opens the door to research improved metrics to describe how humans detect differences between images, be they synthetic, rendered, or natural.

SSIM is at its core a statistical measure [17], a product of three local dissimilarity factors, namely, luminance, variance, and correlation. We have derived these factors' minima and shown how their ranges and normalization are creating non-intuitive results. This occurs, for example, for low luminance values or when the local distribution of pixel values visually differ very little, though regularly. We have also shown that the original SSIM formulation with certain parameters can output undefined results. For one of the major areas of use for SSIM, namely rendering, these results constitute the core of the index's weaknesses—both as a qualitative indication, using the pooled MSSIM value, and as a quantitative value, when using the SSIM map to understand the visual performance of rendering algorithms.

Current graphics and rendering research has a major focus on Monte Carlo ray tracing and denoising and reconstruction algorithms using neural networks. Such networks often introduce small variations during training and could potentially suffer disproportionately from the shortcomings of SSIM. We thus encourage further graphics research to employ the index with care and caution, or preferably replace it, since its use may distort or bias image quality assessment. The *difference evaluator for alternating images* (ꟻLIP) [1], is a step toward such a replacement.

## Acknowledgements

## References

[1] Andersson, P., Nilsson, J., Akenine-Möller, T., Oskarsson, M., Åström, K., and Fairchild, M. D. ꟻLIP: A Difference Evaluator for Alternating Images. In *High Performance Graphics* (2020).

[2] Brunet, D. *A Study of the Structural Similarity Image Quality Measure with Applications to Image Processing*. PhD thesis, 2012.

[3] Brunet, D., Vrscay, E. R., and Wang, Z. On the Mathematical Properties of the Structural Similarity Index. *IEEE Transactions on Image Processing 21*, 4 (April 2012), 1488–1499.

[4] Čadík, M., Herzog, R., Mantiuk, R., Mantiuk, R., Myszkowski, K., and Seidel, H.-P. Learning to Predict Localized Distortions in Rendered Images. *Computer Graphics Forum 32*, 7 (2013), 401–410.

[5] Čadík, M., Herzog, R., Mantiuk, R., Myszkowski, K., and Seidel, H.-P. New Measurements Reveal Weaknesses of Image Quality Metrics in Evaluating Graphics Artifacts. *ACM Transactions of Graphics 31*, 6 (2012), 147:1–147:10.

[6] Chen, S., Zhang, Y., Li, Y., Chen, Z., and Wang, Z. Spherical Structural Similarity Index for Objective Omnidirectional Video Quality Assessment. In *IEEE International Conference on Multimedia & Expo* (07 2018), pp. 1–6.

[7] Dosselmann, R., and Yang, X. D. A Formal Assessment of the Structural Similarity Index. Technical report tr-cs 2008-2, Department of Computer Science, University of Regina, Canada, 2008.

[8] Dosselmann, R., and Yang, X. D. A Comprehensive Assessment of the Structural Similarity Index. *Signal, Image, and Video Processing 5*, 1 (2011), 81–91.

[9] Girod, B. Digital images and human vision. MIT Press, Cambridge, MA, USA, 1993, ch. What's Wrong with Mean-squared Error?, pp. 207–220.

[10] Horé, A., and Ziou, D. Image Quality Metrics: PSNR vs. SSIM. In *International Conference on Pattern Recognition* (2010), pp. 2366–2369.

[11] Kandel, E. R. *Principles of Neural Science*. McGraw-Hill, 2013.

[12] Larson, E. C., and Chandler, D. M. Most Apparent Distortion: Full-Reference Image Quality Assessment and the Role of Strategy. *Journal of Electronic Imaging 19*, 1 (2010).

[13] Legge, G., and Foley, J. Contrast Masking in Human Vision. *Journal of the Optical Society of America 70*, 12 (1981), 1458–1471.

[14] Liu, H., and Wang, Z. Perceptual Quality Assessment of Medical Images. In *Encyclopedia of Biomedical Engineering*, R. Narayan, Ed., vol. 2. 2018, pp. 588–596.

[15] Piórkowski, R., and Mantiuk, R. Calibration of Structural Similarity Index Metric to Detect Artefacts in Game Engines. In *Computer Vision and Graphics* (2016), pp. 86–94.

[16] Sampat, M. P., Wang, Z., Gupta, S., Bovik, A. C., and Markey, M. K. Complex Wavelet Structural Similarity: A New Image Similarity Index. *IEEE Transactions on Image Processing 18*, 11 (2009), 2385–2401.

[17] Wang, Z., and Bovik, A. C. A Universal Image Quality Index. *IEEE Signal Processing Letters 9*, 3 (March 2002), 81–84.

[18] Wang, Z., and Bovik, A. C. Mean squared error: Love it or leave it? A new look at Signal Fidelity Measures. *IEEE Signal Processing Magazine 26*, 1 (January 2009), 98–117.

[19] WANG, Z., BOVIK, A. C., SHEIKH, H. R., AND SIMONCELLI, E. P. Image Quality Assessment: from Error Visibility to Structural Similarity. *IEEE Transactions on Image Processing 13*, 4 (2004), 600–612.

[20] WANG, Z., LU, L., AND BOVIK, A. C. Video Quality Assessment Based on Structural Distortion Measurement. *Signal Processing: Image Communication 19*, 1 (2004), 121–132.

[21] WANG, Z., SIMONCELLI, E. P., AND BOVIK, A. C. Multiscale Structural Similarity for Image Quality Assessment. In *The Thirty-Seventh Asilomar Conference on Signals, Systems Computers, 2003* (2003), vol. 2, pp. 1398–1402.

[22] WHITTLE, J., JONES, M. W., AND MANTIUK, R. Analysis of Reported Error in Monte Carlo Rendered Images. *The Visual Computer 33*, 6 (2017), 705–713.

[23] WOLSKI, K., GIUNCHI, D., YE, N., DIDYK, P., MYSZKOWSKI, K., MANTIUK, R., SEIDEL, H.-P., STEED, A., AND MANTIUK, R. K. Dataset and Metrics for Predicting Local Visible Differences. *ACM Transactions on Graphics 37*, 5 (2018), 1–14.

[24] ZENG, K., AND WANG, Z. 3D-SSIM for Video Quality Assessment. In *IEEE International Conference on Image Processing* (Sept 2012), pp. 621–624.

[25] ZHAO, H., GALLO, O., FROSIO, I., AND KAUTZ, J. Loss Functions for Image Restoration with Neural Networks. *IEEE Transactions on Computational Imaging 3*, 1 (March 2017), 47–57.