

---

# CUD-NET: Color Universal Design Neural Filter for the Color Weakness

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Information on images should be visually understood to anyone, including the color  
2 weakness. However, it is not recognizable if color that seems distorted to the color  
3 weakness meets an adjacent object. We suggest CUD-NET<sup>1</sup> based on convolutional  
4 deep neural network to generate color universal design (CUD) images that satisfy  
5 both color preservation and distinguishment of color for input images. CUD-NET  
6 regresses the node point of the piecewise linear function based on information of  
7 input images and comprises a specific filter per image. We present the following  
8 methods to generate CUD images for the color weakness. First, we refine the CUD  
9 dataset on specific criteria by color experts. Second, the input image information  
10 is expanded through the pre-processing specialized on the color weakness vision.  
11 Third, we suggest a multi-modal feature fusion architecture that combines features  
12 to process expanded images. Finally, we suggest a deformable loss function by the  
13 composition of the predicted image through the model to avoid the one-to-many  
14 problems of the dataset.

## 15 1 Introduction

### 16 1.1 Motivation

17 The green and red color blindness are made up of 8% of males and 0.5% of females in Northern  
18 European descent[Won11], which is almost up to rate of one person in 20 people. Green and red  
19 blindness is the most common pattern, followed by blue, yellow, and total color blindness. In this  
20 paper, we generate Color Universal Design (CUD) images, which are color weakness friendly design  
21 forms, through deep learning around the aspect of the red color weakness (protanopia) and green  
22 color weakness (deutanopia) vision. Protanopia is insensitive to red color and deutanopia is  
23 insensitive to green color, although it varies depending on individual color weakness extent.

24 There are studies that help color discrimination to the color weakness, including wearable devices  
25 and surgeries[VZCR20]. However, since these research require time and cost, we simply generate  
26 CUD images with an image enhancement method based on deep learning to make the corresponding  
27 color visible for the color weakness. For an example of the left-above image  $I$  in Figure 1, the people  
28 who are not color weakness can distinguish the letter ‘5’ in the image. But as a deutanopia vision in  
29 left-below image  $I^d$ , the surrounding color and the letter ‘5’ are very analogous, making it ambiguous  
30 to distinguish the bound of adjacent object. The right-bottom target image  $T^d$ , refined image by color  
31 expert designers, shows that the letter ‘5’ appeared well at the deutanopia vision. Here, we define  
32 the non-CUD objects as the letter ‘5’ and surroundings invisible to deutanopia vision in the image  
33  $I$ , and define the CUD objects as the letter ‘5’ and surroundings visible to deutanopia vision in the  
34 image  $T$ . In other words, CUD object means that adjacent objects are distinguishable on both the

---

<sup>1</sup>Code available at <https://github.com/Anonymous68864576/CUD-NET-anonymous>

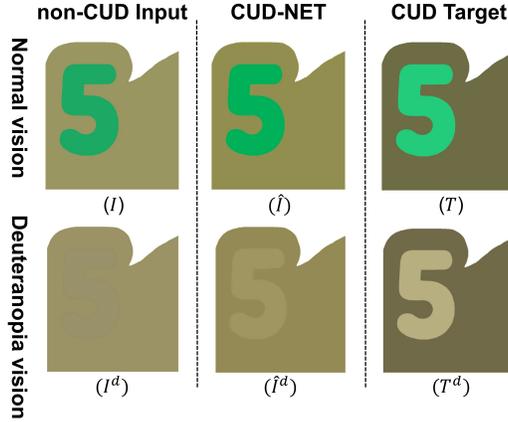


Figure 1: Comparisons of the non-CUD image, our CUD-NET’s predicted image, and CUD image. The above row is represented in normal vision, and below row is represented in deuteranopia vision.

35 normal vision and the color weakness vision. The non-CUD object means that adjacent objects are  
 36 distinguishable on normal vision but not the color weakness vision. Consequently, we generate  $\hat{I}$  that  
 37 satisfies CUD with a specific filter to the image  $I$ .

38 We want to apply as weak filter as possible to CUD objects to preserve color, which requires a  
 39 certain level of object comprehension mechanism to do so. There are various studies from classic  
 40 PCA[WEG87] to machine learning-based object segmentation methods[TSC20, ZGL\*20] to define  
 41 specific objects or areas in image. The research on semantic segmentation, which even provides labels  
 42 between objects, seems that deep learning still does not have a complete comprehension of all objects  
 43 in the real-world. The visual question answering to arbitrary questions about object’s interactions,  
 44 the most general issue on comprehension of object, does not have high transmission power to be  
 45 practical uses[AHB\*18, KZG\*17, LYL\*20]. Therefore, we expand feature of the input image around  
 46 the information of color weakness vision and define the robust neural filter. In summary, we suggest  
 47 a CUD-NET that generates an image suitable for CUD, while complying with the color preservation  
 48 for the source image.

49 In this paper, we suggest the Color Universal Design Network (CUD-NET) to satisfy both color  
 50 preservation and contrast of non-CUD objects (CUD suitability). We introduce 4 core contributions  
 51 of CUD-NET.

- 52 • **Dataset refinement criteria for CUD image** We refine training data into two groups, the  
 53 one with a simple color tone image based on H and V in the HSV color space, the other with  
 54 two or more non-CUD objects that must be distinguished in publications.
- 55 • **Image pre-processing for CUD-NET** We carry out pre-processing to expand the infor-  
 56 mation of the input image. Input image  $I$  is reconstructed with three expanded feature  
 57 information with noise removed.
- 58 • **Multi-modal feature fusion architecture** We define a feature layer, the fusion layer, and  
 59 a regression layer to handle pre-processed images. The three features from the feature  
 60 extracting layer are combined into the one fusion feature, and finally a filter is constructed  
 61 by regressing the node point of the piecewise linear function, or indicator of filter.
- 62 • **Variational loss function** We suggest a deformable loss function by the composition of the  
 63 predicted image through the model. Our data have a problem of one-to-many, where the  
 64 specific color in input image  $I$  is mapped into multiple colors in target image  $T$ .

## 65 1.2 Related Works

66 **Image-to-Image translation based on GAN** GAN is used in various image translation areas,  
 67 including image generation, style transfer, and colorization[KWK21, IZZE17]. In a preliminary  
 68 experiment, Cycle-GAN[PEZZ20] has reached the best performance in maximizing the contrast of

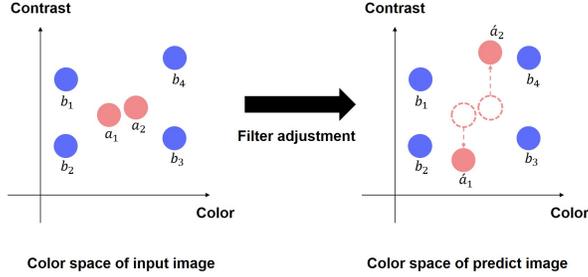


Figure 2: The ideal color conversion of predicted image between contrast and color preservation. Non-CUD object  $a$  should increase the gap compared to the input image and preserve its original color, while the CUD object  $b$  maintain both contrast and color.

69 non-CUD objects. However, our goal is to keep the color preservation of the input image as well,  
 70 so in the case of black color, which has lost all its color of the input image, it is considered the  
 71 worst case for color preservation. Enlighten-GAN[JGL\*21] complements those instability, enabling  
 72 them to generate more stable results on color preservation. But since most of the GAN-based image  
 73 translation fixes the size of the predicted image, reshaping a high-resolution image causes information  
 74 loss of source image. Also, it is difficult to reconstruct the complete geometry for the source image  
 75 as it generates images through the dilated convolution layer.

76 **Image enhancement based on neural filter estimation** Unlike GAN, there are researches that  
 77 scale the pixel values of images based on neural filter estimation[WZF\*19, DLT18, BCPS19]. Zero-  
 78 DCE[GLG\*20] is a low-light image enhancement research that provides a brighter visual display  
 79 of input image. It estimates pixel-wise and high-order filter for dynamic range adjustment of input  
 80 images with lightweight deep network, DCE-Net. DeepLPF[MMM\*20] tried to solve the problem by  
 81 using a graduated filter, elliptical filter, and polynomial filter. The authors not only tried to visually  
 82 enhance the contrast of images but also to comprise stable filters that are easy to understand for  
 83 the spectators while keeping the color preservation. In our problem, however, the contrast factor is  
 84 almost same results as the input image in both visions, while complying the high color preservation,  
 85 resulting over-stable filter. It is assumed that the inability in comprehension of object’s interaction  
 86 leads to over-stable filter.

## 87 2 Methodology

88 We define the ideal predicted image as an increase in the contrast between non-CUD objects and the  
 89 color preservation for the input image. The non-CUD object  $a$  should be mapped into  $\hat{a}$  and CUD  
 90 object  $b$  should preserve its color and contrast like an ideal example of Figure 2. However, as our  
 91 neural filter affects the whole pixels throughout the image, we have the constraint of applying the  
 92 same filter to objects  $a$  and  $b$ . It is very hard to make the contrast and color of object  $b$  exactly the  
 93 same as before the filter adjustment while maximizing the contrast of object  $a$ . Therefore, we propose  
 94 a deep learning-based regression to comprise the specific filter per image that maximizes the contrast  
 95 of object  $a$  while minimizing the adjustment of features on object  $b$ .

96 First, we propose a solution to maximize the contrast of the L channel values in CIELab color  
 97 space[RG19]. We empirically confirmed that protanopia and deuteranopia, which account for the  
 98 most proportion of color weakness, can distinguish the difference by L channel values in common  
 99 when the non-CUD objects are adjacent to each other. To illustrate Figure 1 again, the L channel  
 100 value of letter ‘5’ in image  $I$  is 61 and the surrounding color is 61. The distinguishment between  
 101 the two objects is easy to normal vision, however the image  $I^d$ , the deuteranopia vision, is very  
 102 ambiguous. On the contrary, the CUD target image  $T$  and  $T^d$  have a difference of L channel value 75  
 103 for the letter ‘5’ and 45 for the surroundings, making it easy to distinguish between the normal and  
 104 the deuteranopia vision. Due to the characteristics of these data, we refine a data pair by defining a  
 105 criterion that separates two invisible non-CUD objects by L channel values.

106 Secondly, we propose a variational loss function and multi-modal feature fusion network for color  
 107 preservation. It can be said that the increase in the contrast of L channel values between non-CUD

108 objects is quantitatively superior, but not in the case of increasing the differences in color preservation  
 109 of input images. When non-CUD objects exist, as a simple example, the most likely way to maximize  
 110 contrast is to polarize the color of the object black and white. But it is the result of complete ignorance  
 111 for color preservation, so just enabling to distinguish between non-CUD objects is not always a good  
 112 answer. A strong filter must be applied to distinguish non-CUD object, but its impact should not be  
 113 too extensive to leading the loss of information in CUD objects. In this paper, we solve this problem  
 114 by taking an appropriate trade-off of color preservation and contrast of L channel value.

## 115 2.1 Dataset refinement criteria for CUD image

116 The training data is refined by two groups. The one is vectorized image with two colors divided  
 117 by value V and hue degree H in HSV color space[HMKO19], and the other is image with two  
 118 or more objects that must be distinguished while preserving the color of non-CUD objects. The  
 119 training data is grouped about 1,600 color combinations into the same V and then simulates them  
 120 with the deuteranopia vision, converting to the adjacent color family to comply color preservation.  
 121 All conversions are scaled within only S and V in HSV color space to increase at least 15 difference  
 122 in the L channel value of selected non-CUD objects. The colors are combined with the 10 essential H  
 123 and tones, and the similar color simulated with the deuteranopia vision was converted. Consequently,  
 124 the key part of refining training data is preservation of color, allowing the models to comply with the  
 125 same approach on learning.

## 126 2.2 Image pre-processing

127 Our model regresses node points of piecewise linear function, which will be described in the model  
 128 architecture section, and the final filter is a multiplication operation for the input image. Therefore,  
 129 the multiplication operations of less than the number 1 tend to fade the color saturation. The image  
 130 without color inversion converges the white color value to 1, so if the multiplying value is in the [0,  
 131 1] range, the white color is shifted to the black. By inverting the color of input image, it ignores the  
 132 multiplication operations for white value with 0.

133 We generate the map image  $I^m$  based on original RGB input images calculating the difference value  
 134 between the image with an aspect of normal vision and the image with an aspect of deuteranopia vision.  
 135 Recent studies have been conducted to augment the information or expanded the models' perspective  
 136 through transformer models[JSZK16, RFB15]. In our experiment, however, the transformer model  
 137 tends to generate the predicted image ignoring the source color, which result in the polarized color to  
 138 black and white like Cycle-GAN's.

$$I^m = |invert(I^n) - invert(I^d)| \quad (1)$$

$$I = \delta(cat^{channel}(I^n, I^d, I^m)) \quad (2)$$

139 After applying color inversion from the original RGB input image  $I^n$ , we generate the image  $I^d$  with  
 140 an aspect of deuteranopia vision in equation 1. From these two generated images, we can get the  
 141 absolute difference value to compose the map image  $I^m$ . In equation 2, the final input  $I$  concatenated  
 142 with  $9 \times H \times W$  dimensions passes through the model. The  $\delta(\cdot)$  clips output to a range of [0, 1].

## 143 2.3 Model Architecture

144 CUD-NET regresses the node points of piecewise linear filter function from the input. The value  
 145 of each node points computes the multiplication operation and generates the predicted image. The  
 146 input  $I$  is compressed into 3 feature blocks matching each input through convolution layer, pooling  
 147 layer, and global pooling layer. The input with 9 channels is separated into  $3 \times 3$  channels before  
 148 passing the model. The first 3 channels are literally used as the main inputs, where the multiplication  
 149 operation takes place, while the remaining 6 channels are used as features.

150 First of all, we use multi-modal fusion architecture for three separate inputs to extract expanded  
 151 features. The three inputs converted to the HSV color space pass through a weights-sharing convolu-  
 152 tion layer to extract a feature block corresponding to the inputs. Each convolution layer consists of  
 153 kernel size=3, stride=1, and padding=1, reducing dimension through average pooling. We empirically

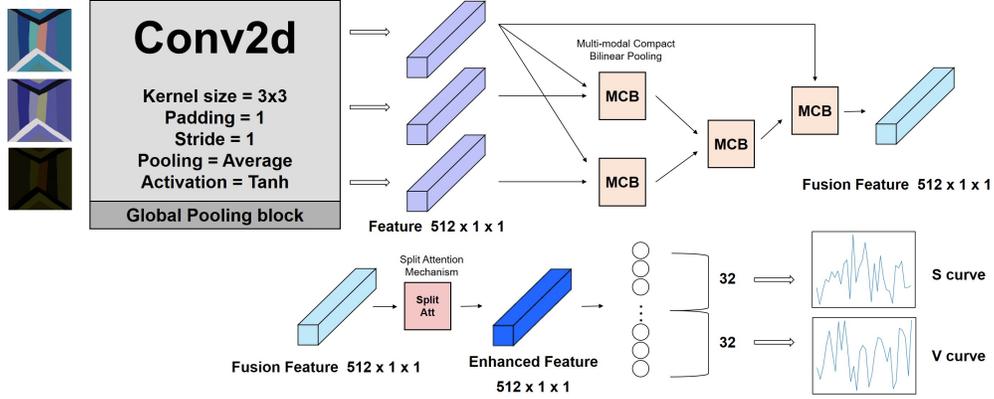


Figure 3: Overview structure of CUD image generation

154 noticed that the most values of output feature have distribution within the range of  $[-1, 1]$  with valid  
 155 values for constructing the node points, so we use hyperbolic tan for activation function. Since we  
 156 use inputs with unstructured image size, the last global pooling block holds the size of the feature  
 157 instead of the average pooling block[LCY14].

158 The three feature blocks are combined through the multi-modal compact bilinear pooling gate  
 159 (MCB)[FPY\*16], following the fusion process shown in Figure 3. The MCB gate allows both  
 160 features to interact in a multiplicative way with low memory consumption and computation  
 161 times. The fusion features are complemented to enhanced feature through the split attention  
 162 mechanism[ZWZ\*20]. At the beginning of the experiment, we have applied the convolutional  
 163 block attention mechanism[WPLK18] of each MCB gate, but we found that it does not make sense  
 164 of understanding the feature itself, so we apply only one attention block to the last fusion feature.

165 The enhanced feature pass through the fully-connected regression layer. We picked the 64 points  
 166 to be regressed to compose the piecewise linear function, which is empirically confirmed to the  
 167 optimized number of points in this research. The first half of the values construct the node points of  
 168 the S channel and the other half comprise the V channel in HSV color space. Finally, node points  
 169 become the scaling factors to generate predicted image in equation 3[MMS19].

$$S(I_i^{s,v}) = k_0 + \sum_{m=0}^{M-1} (k_{m+1} - k_m) \delta(MI_i^{s,v} - m) \quad (3)$$

170 The total number of node point  $M$ , each pixel values of S, V channel in input image  $I_i^{s,v}$  are  
 171 multiplied with the slope of actual regressed value  $k_m$ , the  $m$ -th generated node point. The  
 172 specific node points  $M$  is scaled through a multiplication operation to pixel value of the input image  
 173 according to each node point.

## 174 2.4 Loss function

175 Our dataset has one-to-many problems between input and target data. In dataset pair  
 176  $(I_1, T_1), (I_2, T_2), \dots, (I_n, T_n)$ , for example, the red color in  $I_1$  can be targeted to purple color in  
 177  $T_1$ , and the red color in  $I_2$  can be targeted to orange color in  $T_2$ . With these one-to-many dataset  
 178 structures, we design the loss function  $\mathcal{L}$  that expresses the potential and the diversity of predicted  
 179 image in equation 4.

$$\mathcal{L} = \sum_{i=1}^N Lab_{loss} \left( V \left( \Phi \left( \hat{I}_i \right) \right) \right) + H_{loss} \left( \Phi \left( \hat{I}_i \right) \right) \quad (4)$$

180 **Stencil Masking** As explained in the dataset refining criteria, we do not proceed with color  
 181 conversion for all areas in the target images, but only for areas with color combinations that are

182 invisible to the deuteranopia (non-CUD object). For this reason, the input image has color regions of  
 183 converting color and unconverted color, which also can be referred to as non-CUD object and CUD  
 184 object. To imply the color bound to model, the stencil masking method is introduced.

$$\Phi(\hat{I}_i) = \hat{I}_{ij} \parallel (I_{ij} \cdot T_{ij}) \quad (5)$$

185 We consist a stencil maps through the logical and operations ‘.’ of each pixel value  $I_{ij}, T_{ij}$ . Stencil  
 186 map can specify the non-CUD area and be computed with predicted image  $\hat{I}_{ij}$  of logical or operation  
 187 ‘||’ in equation 5. Consequently, CUD object of the image adjusted with a stencil mask does not carry  
 188 out the neural filter computation, such as the same way we refine the target image. This refined image  
 189 is calculated on the loss function.

190 **CIELab Loss** We use the CIELab channel loss function to maximize the contrast of color on  
 191 deuteranopia vision. To stabilize the contrast and brightness of the predicted image, we calculate the  
 192 MS-SSIM(multi-scale structural similarity[WSB03]) of L channel.

$$Lab_{loss} = \left\| Lab(\hat{I}_i^{rgb}) - Lab(T_i^{rgb}) \right\|_1 + MS\text{-}SSIM\left(Lab(\hat{I}_i^L), Lab(T_i^L)\right) \quad (6)$$

193 The  $Lab(\cdot)$  expression in equation 6 returns the CIELab channel corresponding to the RGB channel,  
 194 and all calculations are made only on the L channel.

195 **Histogram Loss** We use the histogram loss function to comply with the color preservation of the  
 196 image. The RGB channel is used to preserve its color, contrary to using only the L channel in other  
 197 loss functions. Handling the RGB channel as a loss function rather than using Lab’s ab channels has  
 198 shown better results on color preservation.

$$H_{loss} = -\omega_{hist} \int N(\hat{I}_i^{rgb}; \sigma) - N(T_i^{rgb}; \sigma) \quad (7)$$

199 When simply designing a loss function with the L1 distance of the RGB channel pixel values, it was  
 200 very sensitive to certain values and the gradients are diverged, resulting in an untrainable experiment.  
 201 Therefore, we used a gaussian expansion method[SAC\*17] denoted by  $N(\cdot)$  to infer a differentiable  
 202 histogram loss function in equation 7. We compute the difference of the RGB channel of the  
 203 differentiable histogram function, which can be altered to mean squared error or cosine similarity.  
 204 The scaler  $\omega_{hist}$  is determined in inverse proportion to the size of the input image. By maintaining  
 205 the RGB similarity between the predicted image and the target image, we can comply with the color  
 206 preservation.

207 **Variational Prediction** There are various ways to maximize difference of the L channel in the  
 208 image. And the target image is converted at least two colors compared to the input image. However,  
 209 the predicted image of the model is generated by the neural filter, so it is unpredictable which area  
 210 of color is modified. Therefore, if the color in predicted image is over-shifted or in the color value  
 211 of opposite shifts to the target, the loss will rather increase. In addition to one-to-many problem  
 212 that the data pair itself does not matches one-to-one in a particular color, it is necessary to generate  
 213 alternative predicted image with the same aspect of the data pair. We calculate the loss function with  
 214 a variational prediction based on the predicted image for potential color shifts.

215 The first potential is the case of excessive shifts. Assume that  $I_i^L = \{74, 41, 79\}$ ,  $\hat{I}_i^L = \{97, 10, 70\}$ ,  
 216  $T_i^L = \{50, 41, 80\}$  in L channel value. The first and third components of each image is non-CUD  
 217 objects, and second component is CUD object. Therefore, we refined data paired with a difference of  
 218 15 on L channel. Here, we clip the excessive L channel value in  $\hat{I}_i^L$  by equation 8. Up to this point,  
 219 no calculation is made as no value is exceeded in this example. The second potential is the case of  
 220 opposite shifts. It can be said that a complete neural filter has been proceeded for value 97, 10, 70  
 221 where L channel difference is 27. However, if we actually calculate the mean square error between  
 222  $\hat{I}_i^L$  and  $T_i^L$ , it will be an large value over 1k. Here we can generate alternative predicted image from  
 223 equation 9 and 10.

$$\text{clip}(\hat{I}_{ij}) = \begin{cases} \max(\hat{I}_{ij}, T_{ij}), & I_{ij} > T_{ij} \\ \min(\hat{I}_{ij}, T_{ij}), & I_{ij} \leq T_{ij} \end{cases} \quad (8)$$

$$R_1 = 2I_{ij} - \hat{I}_{ij}, R_2 = \hat{I}_{ij} \quad (9)$$

$$V(\hat{I}_{ij}) = \text{argmin}(\|\text{clip}(R_{1,2}) - T_{ij}\|_2) \quad (10)$$

224 As mentioned above, we define thresholds by the maximum and minimum value of each corresponding  
 225 pixel position of  $\hat{I}_i^L$  and  $T_i^L$ . By computing a difference of residual map and the input image, we  
 226 induce the alternative two images  $R_1, R_2$ . As a result,  $\hat{I}_i^L$  with a smaller L2 distance is selected to  
 227 alternative predicted image in equation 10, and it is finally computed with loss function compared to  
 228 the  $T_i^L$ . The above equation establishes  $V(\Phi(\hat{I}_i)) = \{54, 41, 80\}$  and the mean square error to the  
 229 target image is approximately 5, which is agreeable loss value respect to  $\hat{I}_i^L$  itself.

230 **Identity Loss[ZPIE17, TPW16]** We use  $\mathcal{L}_{identity}(T_i)$  to apprehend the CUD object to the model.  
 231 In the case of target image that already satisfy the CUD, the filter should be relatively weakly applied  
 232 than input image. The input of identity loss is target image  $T_{ij}$  instead of input image  $I_{ij}$ , and the  
 233 reference of the loss function is also target image  $T_{ij}$  to maintain the value itself. In computing  
 234 identity loss, we do not require variational prediction as we cannot judge the potential region by  
 235 equation 10.

### 236 3 Experiments

237 The experiment was performed with Tesla V100 SXM2 and Intel Xeon Gold 5120 and the computation  
 238 speed was about 40 images per minutes. We refined a dataset with Adobe Photoshop to maximize  
 239 contrast in the L channel by adjusting saturation and brightness for areas that require color conversion  
 240 based on deuteranopia vision simulation. Color experts has refined about 1,500 vectorized image for  
 241 the training data and 300 publication images for the validation data. All the comparative experimental  
 242 models used the same train, test, validation data in this paper. We used the inference data in  
 243 publications, which is almost composed of vectorized images, as colors often appear distorted in a  
 244 gradation-rich image. The Figure 4 is arranged in descending order of the number of combinations in  
 245 colors from the top image.

246 Both structure similarity (SSIM)[ZBSS04] and peak signal to noise ratio (PSNR) in Table 1 can  
 247 indicate whether the image is suitable for CUD or not. As a notable aspect, the result has shown that  
 248 comparative models with lower metrics are sensitive to high-gradation input images, which generated  
 249 color-heterogeneous image. SSIM and PSNR itself can determine the increase in contrast compared  
 250 to the target image but do not determine whether the color preservation complied. Therefore, we  
 251 evaluated SSIM and PSNR with three references, inputs images  $I$ , predicted images  $\hat{I}$ , and target  
 252 images  $T$ . The higher the estimation of the  $\hat{I}$  and  $I$ , the more color preservation factor worked.  
 253 The higher the estimation of the  $\hat{I}$  and  $T$ , the more increase in contrast can be considered. We also  
 254 define SSIM mean absolute error, PSNR mean absolute error to measure the extent of the conversion  
 255 between the  $F : I \rightarrow T$  and  $F : I \rightarrow \hat{I}$  in equation 11 and 12, respectively. The  $N$  is total number  
 256 of inference data.

$$SSIM \cdot MAE = \frac{1}{N} \sum_{i=1}^N \left| SSIM(\hat{I}_i, T_i) - SSIM(I_i, T_i) \right| \quad (11)$$

$$PSNR \cdot MAE = \frac{1}{N} \sum_{i=1}^N \left| PSNR(\hat{I}_i, T_i) - PSNR(I_i, T_i) \right| \quad (12)$$

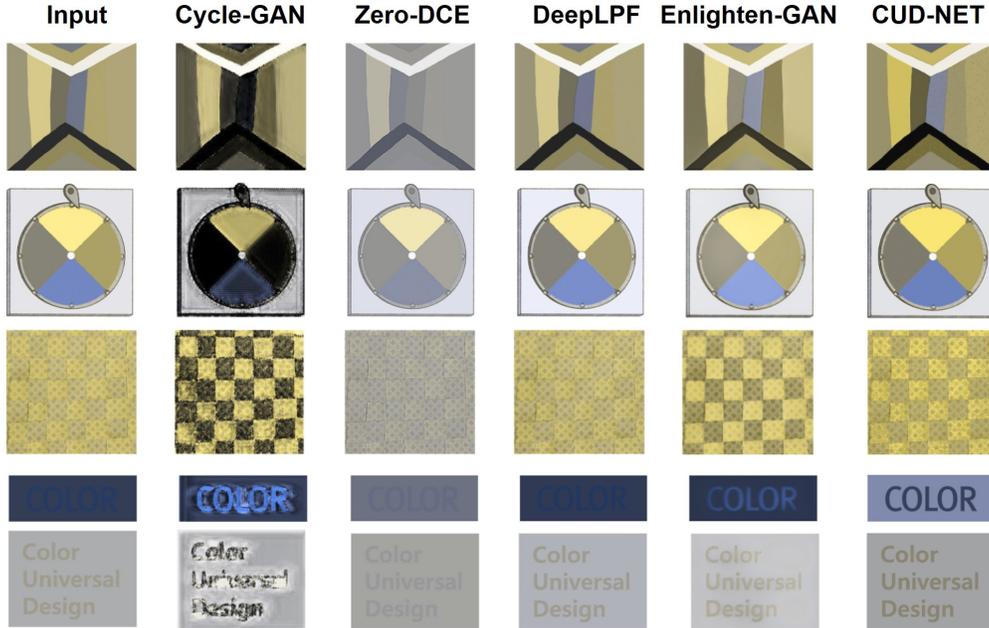


Figure 4: Comparisons of predicted images in deuteranopia vision. The color experts selected the validation data that do not satisfy the CUD in publications.

Architecture	$SSIM(\hat{I}, I)$	$SSIM(\hat{I}, T)$	$PSNR(\hat{I}, I)$	$PSNR(\hat{I}, T)$	SSIM-MAE	PSNR-MAE
Cycle-GAN	0.630	0.634	13.28	13.83	0.3191	8.4430
Zero-DCE	0.924	0.888	21.95	18.67	0.0661	3.7300
DeepLPF	0.850	0.831	26.31	20.34	0.1220	2.0566
Enlighten-GAN	0.820	0.808	21.85	19.58	0.1470	3.9983
Enlighten-GAN(scaled)	<b>0.966</b>	0.921	24.86	<b>21.36</b>	0.0392	3.4937
CUD-NET(low bottle-neck feature)	0.897	0.866	27.77	21.01	0.0901	2.0826
CUD-NET	0.962	<b>0.924</b>	<b>29.54</b>	21.19	<b>0.0312</b>	<b>1.4760</b>

Table 1: Evaluation table of comparison experiment. The CUD-NET with a low bottle-neck feature achieves better results in the experiment of the deuteranopia and the protanopia subjects(Figure 5), although the evaluation metrics are lower than that of CUD-NET.

257 Cycle-GAN and Zero-DCE showed worse result than others. Cycle-GAN model had difficulty  
 258 reconstructing a geometry of a particular object, and overall color had low saturation and brightness,  
 259 resulting in color conversion into an almost grey scale image. Zero-DCE is faded in color, and the  
 260 contrast was not much different from the input image. The overall image lost its color preservation,  
 261 which we focused to solve in this paper.

262 DeepLPF meets both the color preservation and contrast that we deal with for. However, DeepLPF  
 263 tends to color be over-stably filtered for the images with fewer color combinations. Although the  
 264 color preservation has complied better than other experiments, there were many failed results from  
 265 the perspective of contrast, which the over-stable filter leads to by DeepLPF.

266 Remarkably, predicted images of Enlighten-GAN showed reasonable results. However, simple color  
 267 combinations or the images with already satisfying the CUD often showed results degenerated with  
 268 low CUD suitability. Enlighten-GAN was able to generate the results we targeted, but its deviation of  
 269 filter is so high that it sometimes failed to satisfy the contrast even on simple images or decreased  
 270 the contrast. As the problem of GAN-based method including Enlighten-GAN, moreover, model  
 271 fixes the width and height of the predicted image. If width and height of  $T$  and  $I$  down-scaled

272 to size of Enlighten-GAN  $\hat{I}$  (approximately 25K pixels in this experiment), the  $SSIM(\hat{I}, I)$  and  
 273  $PSNR(\hat{I}, T)$  showed higher estimation in some metrics than CUD-NET. In the opposite case of  $\hat{I}$   
 274 up-scaled to size of  $T$  and  $I$ , the significantly low estimation was recorded due to the information  
 275 loss of up-scaling problem.

276 CUD-NET showed stable and robust predicted images in both color preservation and increase in  
 277 the contrast compared to other experiments. In comparing the values in the same region of  $I$  and  
 278  $\hat{I}$ , the model scaled two L channel values with opposite side in the most of case, the one goes up  
 279 and the other goes down. When we reduced the number of bottle-neck feature of model, it tends to  
 280 record relatively high deviation of filter scales according to the number of combinations of colors. In  
 281 summary, the CUD-NET showed the highest estimations for 4 evaluation metrics. Moreover, as our  
 282 model adopted a neural filter unlike generation models, there is no loss of information regarding the  
 283 scaling of predicted images.

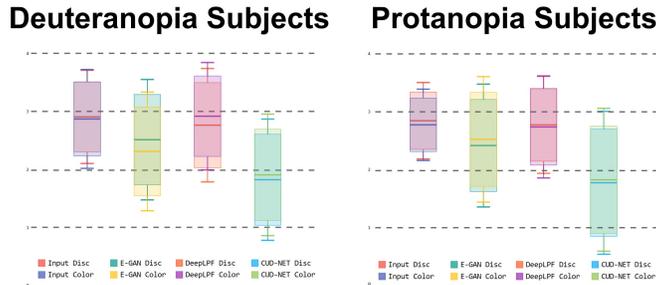


Figure 5: The box bar is ordered to the left side, input image  $I$ , Enlighten-GAN, DeepLPP, CUD-NET. The y position of box bar represents a mean and length of the box bar represents a deviation of each experiment. The lower the graph is, the higher the rank is.

284 The figure 5 shows the evaluation of the deuteranopia and the protanopia. The evaluation metrics  
 285 consist of object distinguishability and color harmony in order of input image  $I$ , predicted image  
 286 of Enlighten GAN, DeepLPP, and CUD-NET. User study has tested upon the total of 6 subjects, 4  
 287 deuteranomaly and 2 protanomaly. The subjects were asked to list the ranks of object distinguishability  
 288 and color harmony of 4-paired-image for each model-blinded item. As the experimental results,  
 289 the deuteranopia subject ranked the 1-st in the object distinguishability of CUD-NET at an average  
 290 rank of 1.821, followed by Enlighten-GAN at an average rank of 2.512. Similarly, the protanopia  
 291 subject also ranked the 1-st in CUD, followed by Enlighten-GAN, DeepLPP, and input images. The  
 292 evaluation of color harmony showed that the subjects tend to assume that the image with a good  
 293 object distinguishability has good color harmony preferentially. For a total of six subjects, the five  
 294 subjects chose the CUD-NET, with the exception of one who ranked Enlighten-GAN by a subtle gap

## 295 4 Conclusion

296 In this paper, we proposed deep network to generate CUD images from non-CUD input images. The  
 297 pre-processing and multi-modal fusion layer could comprehend the information for color weakness,  
 298 and the variational loss function makes the model further adapt to CUD dataset. Compared to other  
 299 research, we are able to maintain high-resolution images and both stable color preservation and  
 300 contrast with neural filter per images.

301 Our current research shows a robust filter for a single color, such as vectorized images, but it is  
 302 difficult to expect stable results in the case of a real-world image with high gradation in hues. We  
 303 consider the same limitation of our work when the certain pixel values react sensitively, making noise  
 304 appear more prominent in the predicted image. In the future, we plan to create additional datasets  
 305 with gradation on the vectorized image and focus on the fusion layer to improve performance of the  
 306 model.

## 307 References

- 308 [AHB\*18] ANDERSON P., HE X., BUEHLER C., TENEY D., JOHNSON M., GOULD S., ZHANG L.: Bottom-  
309 up and top-down attention for image captioning and visual question answering. In *Proceedings of*  
310 *the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (June 2018).
- 311 [BCPS19] BIANCO S., CUSANO C., PICCOLI F., SCHETTINI R.: Content-preserving tone adjustment for  
312 image enhancement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
313 *Recognition (CVPR) Workshops* (June 2019).
- 314 [DLT18] DENG Y., LOY C. C., TANG X.: Aesthetic-driven image enhancement by adversarial learning.  
315 In *Proceedings of the 26th ACM International Conference on Multimedia* (New York, NY, USA,  
316 2018), MM '18, Association for Computing Machinery, p. 870–878. URL: [https://doi.org/](https://doi.org/10.1145/3240508.3240531)  
317 [10.1145/3240508.3240531](https://doi.org/10.1145/3240508.3240531), doi:10.1145/3240508.3240531.
- 318 [FPY\*16] FUKUI A., PARK D. H., YANG D., ROHRBACH A., DARRELL T., ROHRBACH M.: Multimodal  
319 compact bilinear pooling for visual question answering and visual grounding, 2016. arXiv:  
320 1606.01847.
- 321 [GLG\*20] GUO C., LI C., GUO J., LOY C. C., HOU J., KWONG S., CONG R.: Zero-reference deep curve  
322 estimation for low-light image enhancement. In *Proceedings of the IEEE/CVF Conference on*  
323 *Computer Vision and Pattern Recognition (CVPR)* (June 2020).
- 324 [HMKO19] HIRA S., MATSUMOTO A., KIHARA K., OHTSUKA S.: Hue rotation (hr) and hue blending (hb):  
325 Real-time image enhancement methods for digital component video signals to support red-green  
326 color-defective observers. *Journal of the Society for Information Display* 27, 7 (2019), 409–426.  
327 URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/jsid.758>, arXiv:<https://onlinelibrary.wiley.com/doi/pdf/10.1002/jsid.758>, doi:[https://doi.org/10.](https://doi.org/10.1002/jsid.758)  
328 [1002/jsid.758](https://doi.org/10.1002/jsid.758).
- 330 [IZZE17] ISOLA P., ZHU J.-Y., ZHOU T., EFROS A. A.: Image-to-image translation with conditional  
331 adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern*  
332 *Recognition (CVPR)* (July 2017).
- 333 [JGL\*21] JIANG Y., GONG X., LIU D., CHENG Y., FANG C., SHEN X., YANG J., ZHOU P., WANG Z.:  
334 Enlightengan: Deep light enhancement without paired supervision. *IEEE Transactions on Image*  
335 *Processing* 30 (2021), 2340–2349. doi:10.1109/TIP.2021.3051462.
- 336 [JSZK16] JADERBERG M., SIMONYAN K., ZISSERMAN A., KAVUKCUOGLU K.: Spatial transformer  
337 networks, 2016. arXiv:1506.02025.
- 338 [KWK21] KUMAR M., WEISSENBORN D., KALCHBRENNER N.: Colorization transformer. In *International*  
339 *Conference on Learning Representations* (2021). URL: [https://openreview.net/forum?id=](https://openreview.net/forum?id=5NA1PinlGFu)  
340 [5NA1PinlGFu](https://openreview.net/forum?id=5NA1PinlGFu).
- 341 [KZG\*17] KRISHNA R., ZHU Y., GROTH O., JOHNSON J., HATA K., KRAVITZ J., CHEN S., KALANTIDIS  
342 Y., LI L.-J., SHAMMA D. A., BERNSTEIN M. S., FEI-FEI L.: Visual genome: Connecting  
343 language and vision using crowdsourced dense image annotations. *Int. J. Comput. Vision* 123, 1  
344 (May 2017), 32–73. URL: <https://doi.org/10.1007/s11263-016-0981-7>, doi:10.1007/  
345 [s11263-016-0981-7](https://doi.org/10.1007/s11263-016-0981-7).
- 346 [LCY14] LIN M., CHEN Q., YAN S.: Network in network, 2014. arXiv:1312.4400.
- 347 [LYL\*20] LI X., YIN X., LI C., ZHANG P., HU X., ZHANG L., WANG L., HU H., DONG L., WEI F.,  
348 CHOI Y., GAO J.: Oscar: Object-semantics aligned pre-training for vision-language tasks. In  
349 *Computer Vision – ECCV 2020* (Cham, 2020), Vedaldi A., Bischof H., Brox T., Frahm J.-M.,  
350 (Eds.), Springer International Publishing, pp. 121–137.
- 351 [MMM\*20] MORAN S., MARZA P., MCDONAGH S., PARISOT S., SLABAUGH G.: Deeplpf: Deep local  
352 parametric filters for image enhancement. In *Proceedings of the IEEE/CVF Conference on*  
353 *Computer Vision and Pattern Recognition (CVPR)* (June 2020).
- 354 [MMS19] MORAN S., MCDONAGH S., SLABAUGH G.: Curl: Neural curve layers for global image  
355 enhancement, 2019. arXiv:1911.13175.
- 356 [PEZZ20] PARK T., EFROS A. A., ZHANG R., ZHU J.-Y.: Contrastive learning for unpaired image-to-image  
357 translation. In *Computer Vision – ECCV 2020* (Cham, 2020), Vedaldi A., Bischof H., Brox T.,  
358 Frahm J.-M., (Eds.), Springer International Publishing, pp. 319–345.
- 359 [RFB15] RONNEBERGER O., FISCHER P., BROX T.: U-net: Convolutional networks for biomedical image  
360 segmentation. In *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*  
361 (Cham, 2015), Navab N., Hornegger J., Wells W. M., Frangi A. F., (Eds.), Springer International  
362 Publishing, pp. 234–241.

- 363 [RG19] RIBEIRO M., GOMES A. J. P.: Recoloring algorithms for colorblind people: A survey. *ACM*  
364 *Comput. Surv.* 52, 4 (Aug. 2019). URL: <https://doi.org/10.1145/3329118>, doi:10.1145/  
365 3329118.
- 366 [SAC\*17] SCHÜTT K. T., ARBABZADAH F., CHMIELA S., MÜLLER K. R., TKATCHENKO A.: Quantum-  
367 chemical insights from deep tensor neural networks. *Nature Communications* 8, 1 (2017).
- 368 [TPW16] TAIGMAN Y., POLYAK A., WOLF L.: Unsupervised cross-domain image generation, 2016.  
369 arXiv:1611.02200.
- 370 [TSC20] TAO A., SAPRA K., CATANZARO B.: Hierarchical multi-scale attention for semantic segmentation,  
371 2020. arXiv:2005.10821.
- 372 [VZCR20] VARIKUTI V. N., ZHANG C., CLAIR B., REYNOLDS A. L.: Effect of enchroma glasses on color  
373 vision screening using ishihara and farnsworth d-15 color vision tests. *Journal of American Associ-*  
374 *ation for Pediatric Ophthalmology and Strabismus* 24, 3 (2020), 157.e1–157.e5. URL: <https://www.sciencedirect.com/science/article/pii/S1091853120301002>, doi:<https://doi.org/10.1016/j.jaapos.2020.03.006>.
- 377 [WEG87] WOLD S., ESBENSEN K., GELADI P.: Principal component analysis. *Chemo-*  
378 *metrics and Intelligent Laboratory Systems* 2, 1 (1987), 37–52. Proceedings of  
379 the Multivariate Statistical Workshop for Geologists and Geochemists. URL: <https://www.sciencedirect.com/science/article/pii/0169743987800849>, doi:[https://doi.org/10.1016/0169-7439\(87\)80084-9](https://doi.org/10.1016/0169-7439(87)80084-9).
- 382 [Won11] WONG B.: Points of view: Color blindness, 2011. doi:<https://doi.org/10.1038/nmeth.1618>.
- 384 [WPLK18] WOO S., PARK J., LEE J.-Y., KWEON I. S.: Cbam: Convolutional block attention module. In  
385 *Proceedings of the European Conference on Computer Vision (ECCV)* (September 2018).
- 386 [WSB03] WANG Z., SIMONCELLI E. P., BOVIK A. C.: Multiscale structural similarity for image quality  
387 assessment. In *The Thrity-Seventh Asilomar Conference on Signals, Systems Computers, 2003*  
388 (2003), vol. 2, pp. 1398–1402 Vol.2. doi:10.1109/ACSSC.2003.1292216.
- 389 [WZF\*19] WANG R., ZHANG Q., FU C.-W., SHEN X., ZHENG W.-S., JIA J.: Underexposed photo  
390 enhancement using deep illumination estimation. In *Proceedings of the IEEE/CVF Conference on*  
391 *Computer Vision and Pattern Recognition (CVPR)* (June 2019).
- 392 [ZBSS04] ZHOU WANG, BOVIK A. C., SHEIKH H. R., SIMONCELLI E. P.: Image quality assessment:  
393 from error visibility to structural similarity. *IEEE Transactions on Image Processing* 13, 4 (2004),  
394 600–612. doi:10.1109/TIP.2003.819861.
- 395 [ZGL\*20] ZOPH B., GHIASI G., LIN T.-Y., CUI Y., LIU H., CUBUK E. D., LE Q. V.: Rethinking  
396 pre-training and self-training, 2020. arXiv:2006.06882.
- 397 [ZPIE17] ZHU J.-Y., PARK T., ISOLA P., EFROS A. A.: Unpaired image-to-image translation using  
398 cycle-consistent adversarial networks. In *Proceedings of the IEEE International Conference on*  
399 *Computer Vision (ICCV)* (Oct 2017).
- 400 [ZWZ\*20] ZHANG H., WU C., ZHANG Z., ZHU Y., LIN H., ZHANG Z., SUN Y., HE T., MUELLER  
401 J., MANMATHA R., LI M., SMOLA A.: Resnest: Split-attention networks, 2020. arXiv:  
402 2004.08955.

## 403 Checklist

- 404 1. For all authors...
- 405 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
406 contributions and scope? [Yes]
- 407 (b) Did you describe the limitations of your work? [Yes] See on conclusion section
- 408 (c) Did you discuss any potential negative societal impacts of your work? [N/A] This work  
409 is for positive societal impacts
- 410 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
411 them? [Yes]
- 412 2. If you ran experiments...
- 413 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
414 mental results (either in the supplemental material or as a URL)? [Yes] See on github:  
415 <https://github.com/Anonymous68864576/CUD-NET-anonymous>

- 416 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
417 were chosen)? [Yes] on github repository
- 418 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
419 ments multiple times)? [Yes] See on section 2.3
- 420 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
421 of GPUs, internal cluster, or cloud provider)? [Yes] See on Experiment section
- 422 3. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 423 (a) If your work uses existing assets, did you cite the creators? [Yes] Cooperated with  
424 Co-author
- 425 (b) Did you mention the license of the assets? [Yes] on github repository
- 426 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
427 on github repository
- 428 (d) Did you discuss whether and how consent was obtained from people whose data you're  
429 using/curating? [Yes] on github repository
- 430 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
431 information or offensive content? [N/A]