
Image Colorization using Convolutional Neural Networks

Alex Aranburu*

Arghyadeep Giri†

Rene Gutierrez‡

Steven Reeves§

Abstract

1 For the culmination of the course CMPS 242, Machine Learning, the authors
2 present a method for image colorization using convolutional neural networks.
3 Colorization, taking a black and white image and turning into a color (RGB) image,
4 is inherently an underdetermined problem. Because of this we aim to generate
5 plausible colorizations using the technology of convolutional neural networks.

6 1 Introduction

7 Image colorization is a difficult problem to solve, since one must predict a higher dimensional object
8 from a lower dimensional one. Our approach to this underdetermined problem is to use a 7 layer
9 convolutional neural network to generate the color channels.

10 1.1 Color spaces

11 The Red-Green-Blue(RGB) space is only one way to represent a color image. Going from a traditional
12 RGB representation of an image to grayscale is a trivial matter, in which the grayscale image is just a
13 linear combination of the red, green, and blue channels. However, going the opposite way has many
14 challenges, that is we must predict a three dimensional object from a one dimensional input.

15 Before we construct our neural network, we change the problem. Since going from one channel to
16 three is very difficult, we change our bases into something more manageable - the LAB color space.
17 Here L is the lightness, and is analogous to the grayscale image. A and B, is a two dimensional map-
18 ping of RGB. We can transition between RGB and LAB using the following method: First transition
19 to the XYZ specific illuminants and observers and then construct the "Hunter-Lab" colorspace:

$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = [M] \begin{bmatrix} r \\ g \\ b \end{bmatrix}$$

20 where $[M]$ is calculated from the RGB reference white primaries (X_r, Y_r, Z_r) of the computer system.
21 After obtaining the XYZ space from the original RGB image, the LAB space is calculated via the
22 relationship

$$\begin{aligned} L &= 116f_y - 16 \\ A &= 500(f_x - f_y) \\ B &= 200(f_y - f_z) \end{aligned}$$

*aaranbur@ucsc.edu

†argiri@ucsc.edu

‡rgutie17@ucsc.edu

§sireeves@ucsc.edu

23 where

$$f_d = \begin{cases} \sqrt[3]{d_r} & \text{if } d_r > \epsilon \\ \frac{\kappa d_r + 16}{116} & \text{else} \end{cases} \quad d_r = D/D_r$$

$$d = x, y, z; \quad D = X, Y, Z$$

$$\epsilon = \frac{216}{24389}; \quad \kappa = \frac{24389}{27}$$

24 as per the CIE standard [1].

25 We do this transformation to reduce the dimension of our problem. In the original situation, the
 26 grayscale image is used to predict the RGB channels. After the transformation, we use the L channel,
 27 analogous to grayscale, to predict only the A and B channels. This is a common strategy used in the
 28 literature[2,3,4]. The opencv library offers routines to translate RGB to LAB using the formalism
 29 above, we used this library in our preprocessing and post processing stages in our algorithm.

30 2 Neural network architecture

31 To create a model for color images given only a black and white one, we use a 7 layer convolutional
 32 neural network to extract the A and B channels in the LAB color space. Each layer contains a
 33 convolution, a rectified linear activation and a batch renormalization stage. We include the batch
 34 renormalization to reduce internal covariate shift while training, this helps accelerate the training
 35 stage of the model [5]. After each layer we downsample using a stride of two. At the output layer, we
 36 apply a deconvolution to upsample the data back to a useable dimension. Figure 1 is a schematic
 37 of our network architecture, adapted from Zhang et al's technique [2]. To implement this neural
 38 network, we use tensorflow and the opencv library.

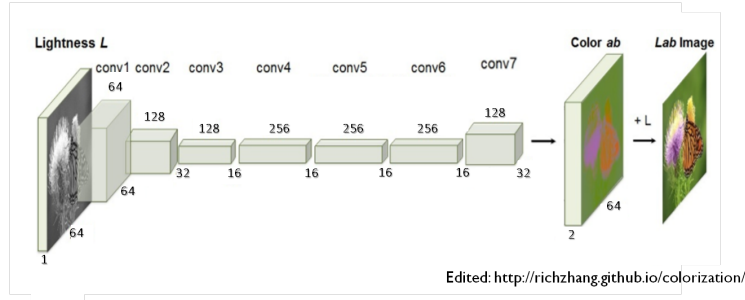


Figure 1: Schematic of the convolutional neural network.

39 2.1 Loss and tuning parameters

40 To train our network defined in figure 1, we use the classic L_2 loss:

$$Loss = ||X_{pred} - X_{labels}||^2$$

41 where X_{pred} are the predicted AB channels from the L input and X_{label} are the true AB channels.

42 2.2 Tune-in

43 There are several hyperparameters that are tunable in our tensorflow implementation. The batch size
 44 of training data, dropout rate of weights, kernel size, the dilation rate or stride and the final activation
 45 function are all things to tune within the model. In table 1, we evaluated different configurations
 46 of this hyperparameters individually, choosing the best performers. Clearly, this approach is very
 47 limited, as ideally each hyperparameter configuration should be considered by itself, and is very
 48 unlikely that the combination of the best performance on each parameter will result on the best
 49 performance overall. Furthermore, the number of configurations tried is very small, however we were
 50 very limited by time and resources, and the tune-in should be improved in the future. For the results
 51 presented we use the following values: a batch size of 10, dropout rate of 0.4, kernel are 3 by 3, and
 52 the dilation rate was 1. For the final activation we found that the hyperbolic tangent to be best suited
 53 for our loss function.

Table 1: Performance Statistics for Different Parameter Configurations.

<i>Batch Size</i>		<i>Dropout Rate</i>		<i>Kernel Size</i>		<i>Stride</i>	
<i>10</i>	<i>20</i>	<i>0.4</i>	<i>0.6</i>	<i>3x3</i>	<i>5x5</i>	<i>1</i>	<i>2</i>
195.55	217.08	215.69	221.48	200.09	209.24	206.2	214.87

3 Results

This section illustrates the results of our model using the above hyperparameters and loss function.

3.1 Data

We trained our model on 32 by 32 and 64 by 64 pixel images downloaded from Imagenet. The data set for each resolution contained 256232 images. We used 90% of the data for training, and 10% to test our model.

3.2 Results for the 32×32 pixel model

We find that the model performs well on certain images and not so well on others. Figure 2 contains images that were the model was more successful.

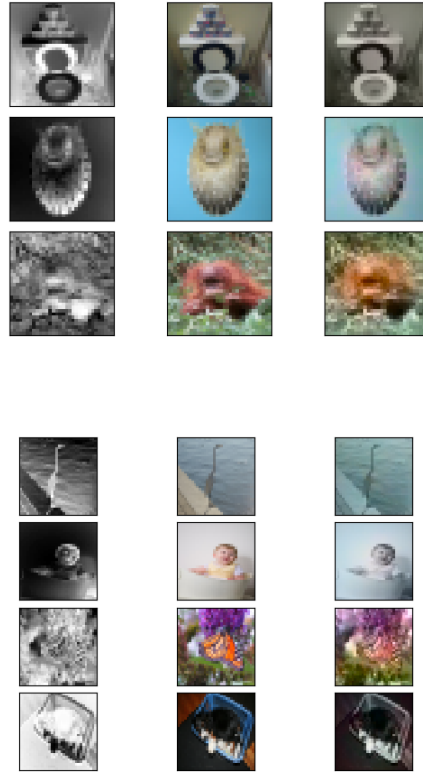


Figure 2: Relative successes on 32×32 pixel images. Left: Grayscale, middle: ground truth, left: model prediction.

In contrast, there were images where the model did not predict the color very well. Figure 3 shows two images where the model was not successful.



Figure 3: Failures of the model on 32×32 pixel images. Left: Grayscale, middle: ground truth, left: model prediction.

65 The first image in figure 3 of the African Grey Parrot was not well resolved. They grey in the
 66 prediction does not match the type of grey on the bird and the red hues of the branch and tail are not
 67 resolved. In the second image the model blurs all colors into different shades of green, losing the
 68 brown color of the turkey.

69 3.3 Test on 64×64 images

70 Perhaps some of the colors are not predicted because of the small resolution of each image in the
 71 previous section. The next figure illustrate the capabilities of the model trained on the 64 by 64 pixel
 72 images.

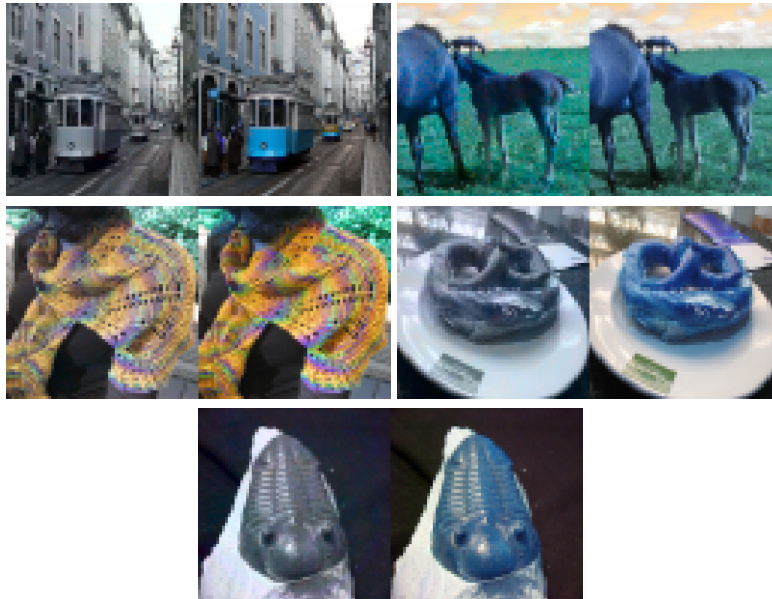


Figure 4: Accuracy of the model on 64×64 pixel images. Left: model prediction, right: ground truth.

73 The images in figure 4 gives more insight into the model than the lower resolution images. Here we
 74 see the effects of the Euclidean loss function. In images with sparse or very localized vibrancy, the
 75 model tends to blur the color, for example images 1,4 and 5 in figure 4. However, in images with

76 consistent color the model performs well, see images 2 and 3. Notice that in image 3, the vibrancy is
77 still diminished, giving the predicted photo a sepia like tone.

78 4 Discussion

79 The above results posit some insight into our methodology. The minimization of the Euclidean norm
80 results in averaging out the vibrancy and color the images, resulting in a sepia-like tone. Furthermore,
81 this type of model seems to perform better on larger resolution images.

82 To move forward, we plan to move to higher resolution images, 256×256 , and change to a more
83 exotic loss function. In [2], they created a loss function that allows for any "believable" configuration
84 of colors in their colorization, essentially allowing for a multimodal loss. Another strategy would
85 be to subtract out the mean luminosity from each training image use that to try and fix the vibrancy
86 issue. As we increase the resolution, we aim to make the CNN into a deep neural network increasing
87 the color feature extraction for better prediction.

88 Additionally, we investigated the use of the cross entropy loss. This would use the softmax function
89 as the ending activation and could be weighted easily to favor vibrancy. However, we found it to be
90 too computationally cost prohibitive for our study.

91 Other techniques could be used to solve this problem. Support vector machines would use 256
92 possible AB vectors using a K-means discretization, where K is chosen via validation. This would
93 extract the SURF features from each pixel, along with a local min and variance. Another solution to
94 remove the effects of a researcher chosen loss function would be to use an adversarial network.

95 References

- 96 [1] http://www.bruceindbloom.com/index.html?Eqn_RGB_XYZ_Matrix.html
- 97 [2] Richard Zhang, Phillip Isola, Alexei A. Efros. Colorful Image Colorization. University of California,
98 Berkeley. <https://arxiv.org/abs/1603.08511>
- 99 [3] Gustav Larsson, Michael Maire, Gregory Shakhnarovich. Learning Representations for Automatic Coloriza-
100 tion. University of Chicago, Toyota Technological Institute at Chicago. <https://arxiv.org/abs/1603.06668>
- 101 [4] Jeff Hwang, You Zhou. Image Colorization with Deep Convolutional Neural Networks. Stanford University.
102 cs231n.stanford.edu/reports/2016/pdfs/219_Report.pdf
- 103 [5] Sergey Ioffe, Christian Szegedy. Batch Normalization: Accelerating Deep Network Training By Reducing
104 Internal Covariate Shift. <https://arxiv.org/abs/1502.03167>