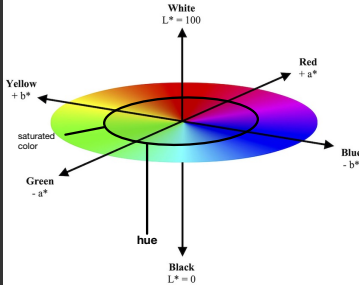


CIE $L^*a^*b^*$ color space :

L = luminance (0 to 100)

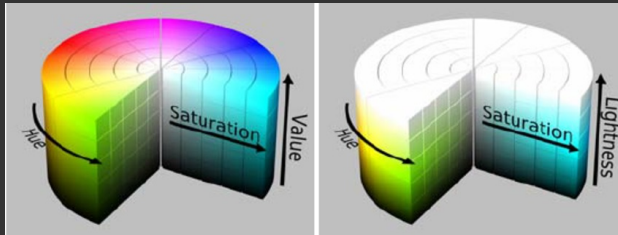
a^* (green-red) (-128 to 127)
 b^* (blue-yellow) (-128 to 127)



Cylindrical view:

Think of chroma (here a^* , b^*)
defining a planar disc at
each luminance level (L)

follow up : HSL (hue saturation lightness)
HSV (hue saturation value)



1 dimensional signal \longrightarrow 3 dimensional signals

Grayscale Image : L channel
 $x \in \mathbb{R}^{H \times W \times 1}$

input

f (high level abstraction / semantics)

Color Information : ab channels

$$\hat{Y} \in \mathbb{R}^{H \times W \times 2}$$

output

image
synthesis
problem

free supervisory signal

Concatenate (L, ab)

$$(x, \hat{Y})$$

refers to the idea that
the data distribution
has a single peak
(mode)

Better Loss Function

Colors in ab space (continuous)

Regression with L2 loss inadequate:

$$L_2(\hat{Y}, Y) = \frac{1}{2} \sum_{h,w} \left\| Y_{h,w} - \hat{Y}_{h,w} \right\|_2^2$$

The L2 loss assumes a unimodal distribution of errors, meaning it is robust when the residuals are centered around zero with a single mode, following a Gaussian (normal) distribution.

If the underlying data or error distribution is multi-modal or has heavy tails, other robust loss functions might be more appropriate.

Colors in ab space (discrete)

Use multinomial classification

discrete
bins of size $\boxed{10}$
(here)

$$L(\hat{Z}, Z) = -\frac{1}{H \cdot W} \sum_{h,w} \sum_q Z_{h,w,q} \log(\hat{Z}_{h,w,q})$$

Class rebalancing to encourage learning of rare colors :

$$L(\hat{Z}, Z) = -\frac{1}{H \cdot W} \sum_{h,w} \nu(Z_{h,w}) \sum_q Z_{h,w,q} \log(\hat{Z}_{h,w,q})$$

Previous works :

(i) non-parametric framework :

reference image (colors) $\xrightarrow{+}$ grayscale image

fail to generalize

(ii) parametric framework :

- L2 regression (i) hand engineered features
(may give sepia, or desaturated result) (ii) deep networks
- Classification

Network Architecture :

x = lightness (L channel)

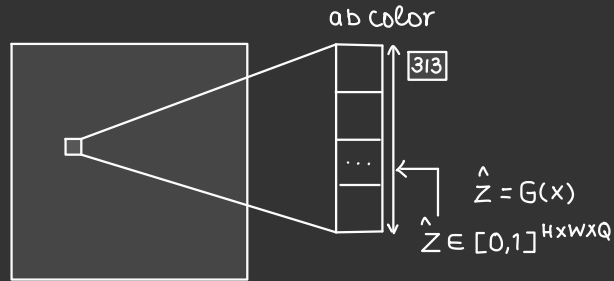
\hat{y} = ab color (ab channel) $\xrightarrow{+L}$ L^*a^*b color image

$$\hat{y} = F(x)$$

probability distribution — pixel level

VGG network :

Layers	dimensions
input x	
conv1	$224 \times 224 \times 1$
conv2	$224 \times 224 \times 64$
conv3	$112 \times 112 \times 128$
conv4	$56 \times 56 \times 256$
conv5	$28 \times 28 \times 512$
fc1	$14 \times 14 \times 512$
fc2	$1 \times 1 \times 4096$
	$1 \times 1 \times 4096$



VGG modified architecture:

Layers	dimensions
input X	
conv1	224 x 224 x 1
conv2	224 x 224 x 64
conv3	112 x 112 x 128
conv4	56 x 56 x 256
conv5	28 x 28 x 512
conv6	28 x 28 x 512
conv7	28 x 28 x 512
conv8	28 x 28 x 512
	256 x 256 x 56

Variation of standard convolutions in deep learning, where the convolutional kernel is expanded by inserting holes.
(larger receptive field without increasing the number of numbers or computation cost)



→ atrous/dilated
Convolutions
(spatial resolution
addition)

single point estimate :

$$\hat{y} = H(\hat{z})$$

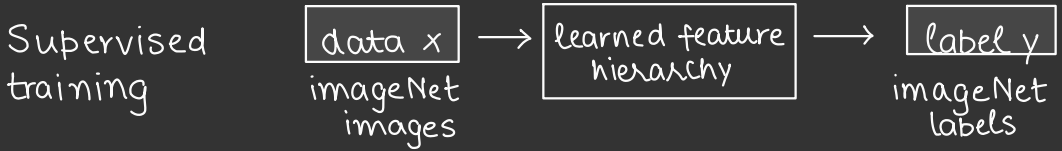
interpolation between the mean and the mode, allows us to keep the vibrancy of the output colors while maintaining some spatial consistency.

Evaluation :

joint interaction b/w pixels
↑
& overall perceptual quality

	visual quality	representation learning
quantitative	<ul style="list-style-type: none"> per pixel accuracy perceptual realism semantic interpretability 	<ul style="list-style-type: none"> task generalization task & dataset generalization
qualitative	<ul style="list-style-type: none"> low-level stimuli legacy grayscale photos 	<ul style="list-style-type: none"> hidden activation units

Predicting Labels from Data :



Predicting Data from Data :

