

VAE-based Image Colorization

Khanh Duong, Tien-Huy Nguyen và Nhu-Tai Do

Ngày 2 tháng 4 năm 2024

Phần I: Giới thiệu

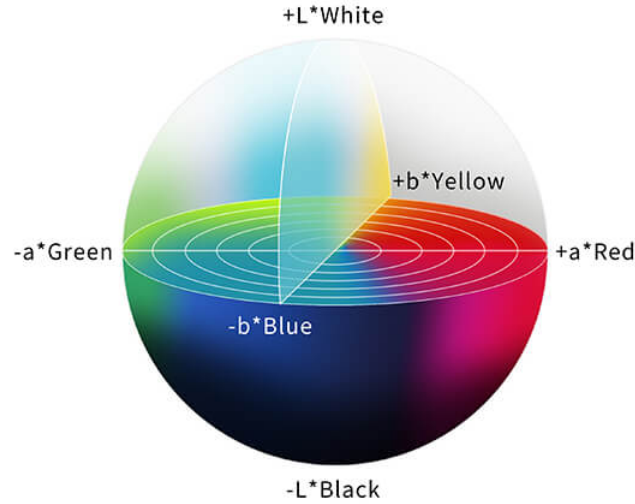
Image Colorization là quá trình dự đoán màu cho các ảnh đen trắng, giúp tái tạo lại hình ảnh thực tế từ dữ liệu đơn sắc, mang lại trải nghiệm hình ảnh phong phú và sống động. Với đầu vào là một ảnh xám, biểu thị cường độ sáng của ảnh, mô hình sẽ học cách ước tính các kênh màu của ảnh, tạo ra một hình ảnh hợp lý và hài hòa về mặt thị giác.



Hình 1: Ví dụ minh họa cho bài toán Image Colorization.

Vấn đề tô màu cho ảnh mang lại một sự thuận lợi đáng kể về mặt dữ liệu, khi mà việc gán nhãn bị loại bỏ hoàn toàn. Bởi, mỗi bức ảnh đều có thể được phân chia thành hai thành phần thông tin chính: gray channel và color channel. Trong đó, gray channel được sử dụng làm đầu vào cho mô hình, đóng vai trò là cơ sở cho quá trình dự đoán. Mô hình sẽ tiến hành dự đoán color channel, tức là thông tin về màu sắc, dựa trên gray channel. Khi mà mô hình đã hoàn thành việc dự đoán, color channel sẽ được kết hợp với gray channel để tạo ra một ảnh hoàn chỉnh.

Trong các bài toán tô màu cho ảnh, không gian màu phổ biến thường được sử dụng là *Lab* thay vì *RGB* như trong các bài toán xử lý ảnh thông thường, trong đó:



Hình 2: Minh họa cho không gian màu *Lab*

- **L (Lightness)**: Đây là thành phần đại diện cho độ sáng của một pixel trong ảnh. Thành phần *L* có thể được xem như phiên bản grayscale của ảnh.
- **a (Green-Red)**: Thành phần *a* biểu diễn sự khác biệt màu giữa các màu xanh lam và đỏ. Khi giá trị của *a* tăng, màu sắc trở nên đỏ hơn. Khi giá trị giảm, màu sắc trở nên xanh hơn.
- **b (Blue-Yellow)**: Thành phần *b* biểu diễn sự khác biệt màu giữa các màu xanh lá cây và màu vàng. Khi giá trị của *b* tăng, màu sắc trở nên vàng hơn. Khi giá trị giảm, màu sắc trở nên xanh hơn.

Bằng cách sử dụng không gian màu *Lab*, mô hình của chúng ta sẽ nhận đầu vào là kênh *L*, đại diện cho độ sáng, và sử dụng kênh *ab* như ground truth của mô hình. Vì vậy, số lượng kênh màu cần dự đoán là 2, bao gồm thành phần màu *a* và *b*, giúp giảm bớt độ phức tạp so với việc dự đoán 3 kênh màu khi sử dụng không gian màu *RGB*.

Có nhiều phương pháp được sử dụng để thực hiện việc tô màu cho ảnh, bao gồm sử dụng mạng CNN truyền thống, hay các mạng tạo sinh đã và đang phát triển trong những năm gần đây như GAN, VAE (Variational Autoencoder) và cả Diffusion Models. Mỗi phương pháp mang lại những ưu điểm và hạn chế riêng, đều hướng tới mục tiêu cuối cùng là tạo ra các ảnh màu tự nhiên và chân thực.

Trong phần này, chúng ta sẽ tập trung vào việc xây dựng một mô hình dựa trên VAE để giải quyết vấn đề Image Colorization. Input và output của chương trình như sau:

- **Input:** Ảnh xám *G* (*L* channel).
- **Output:** Trường ảnh màu *C* (*ab* channels).

Phần II: Nội dung

Trong phần này, chúng ta sẽ triển khai mô hình VAE-base Image Colorization dựa trên bài báo [Learning Diverse Image Colorization](#) để học cách tạo ra tập ảnh màu đa dạng về mặt kết quả. Cụ thể, ta sẽ xây dựng chương trình dựa trên bộ dữ liệu LFW ([Labeled Faces in the Wild Home](#)), một trong những bộ dữ liệu quan trọng và phổ biến trong lĩnh vực nhận dạng khuôn mặt. Bộ dữ liệu này chứa các hình ảnh của khuôn mặt được thu thập từ các bức ảnh chụp thực tế, bao gồm nhiều điều kiện ánh sáng, góc chụp và nền khác nhau.



Hình 3: Ảnh minh họa cho LFW dataset

Theo đó, nội dung thực nghiệm sẽ trình bày với các thành phần như sau:

- a) Data Preparation: Chuẩn bị dữ liệu cho tập huấn luyện.
- b) Models: Xây dựng mô hình VAE và mô hình MDN (Mixture Density Network).
- c) Loss Functions: Xây dựng hàm mất mát cho mô hình VAE và mô hình MDN.
- d) Trainer: Xây dựng các hàm để huấn luyện cho từng mô hình.
- e) Inference: Minh họa kết quả đạt được sau khi huấn luyện mô hình.

1. Data Preparation

Đầu tiên, chúng ta cần chuẩn bị bộ dữ liệu LFW thông qua dòng lệnh dưới đây. Bộ dữ liệu bao gồm hơn 12,000 ảnh train và 1,000 ảnh test. Cùng với đó là một bộ đặc trưng tương ứng với từng ảnh, trích xuất từ một mạng VGG được huấn luyện mạnh mẽ trên bộ dữ liệu lớn ImageNet.

Tải dữ liệu

```
1 !gdown 187x5YSXYibG4QwC5m_Hx8cNzPGVTXv6G
2 !unzip -q data.zip
3 !rm data.zip
```

Khai báo các thư viện:

```
1 import os
2 import numpy as np
3 from tqdm import tqdm
4 import torch
5 import torch.nn as nn
6 import torch.nn.functional as F
7 import torch.optim as optim
```

```

8 from torch.utils.data import DataLoader
9 import cv2
10 import numpy as np
11 from torch.utils.data import Dataset

```

Khởi tạo class ColorDataset:

```

1 class ColorDataset(Dataset):
2     def __init__(self, out_directory, listdir=None,
3                 featslistdir=None, shape=(64, 64),
4                 outshape=(256, 256), split="train"):
5
6         # Save paths to a list
7         self.img_fns = []
8         self.feats_fns = []
9
10        with open("%s/list.%s.vae.txt" % (listdir, split), "r") as ftr:
11            for img_fn in ftr:
12                self.img_fns.append(img_fn.strip("\n"))
13
14        with open("%s/list.%s.txt" % (featslistdir, split), "r") as ftr:
15            for feats_fn in ftr:
16                self.feats_fns.append(feats_fn.strip("\n"))
17
18        self.img_num = min(len(self.img_fns), len(self.feats_fns))
19        self.shape = shape
20        self.outshape = outshape
21        self.out_directory = out_directory
22
23        # Create a dictionary to save weight of 313 ab bins
24        self.lossweights = None
25        countbins = 1.0 / np.load("data/zhang_weights/prior_probs.npy")
26        binedges = np.load("data/zhang_weights/ab_quantize.npy").reshape(2, 313)
27        lossweights = {}
28        for i in range(313):
29            if binedges[0, i] not in lossweights:
30                lossweights[binedges[0, i]] = {}
31                lossweights[binedges[0, i]][binedges[1, i]] = countbins[i]
32        self.binedges = binedges
33        self.lossweights = lossweights
34
35    def __len__(self):
36        return self.img_num
37
38    def __getitem__(self, idx):
39        # Declare empty arrays to get values
40        color_ab = np.zeros((2, self.shape[0], self.shape[1]),
41                            dtype="f")
42        weights = np.ones((2, self.shape[0], self.shape[1]),
43                           dtype="f")
44        recon_const = np.zeros((1, self.shape[0], self.shape[1]),
45                                dtype="f")
46        recon_const_outres = np.zeros((1, self.outshape[0], self.outshape[1]),
47                                       dtype="f")
48        greyfeats = np.zeros((512, 28, 28), dtype="f")
49
50        # Read and reshape
51        img_large = cv2.imread(self.img_fns[idx])
52        if self.shape is not None:
53            img = cv2.resize(img_large, (self.shape[0], self.shape[1]))
54            img_outres = cv2.resize(img_large,

```

```

55         (self.outshape[0], self.outshape[1]))
56
57     # Convert BGR to LAB
58     img_lab = cv2.cvtColor(img, cv2.COLOR_BGR2LAB)
59     img_lab_outres = cv2.cvtColor(img_outres, cv2.COLOR_BGR2LAB)
60
61     # Normalize to [-1..1]
62     img_lab = ((img_lab * 2.0) / 255.0) - 1.0
63     img_lab_outres = ((img_lab_outres * 2.0) / 255.0) - 1.0
64
65     recon_const[0, :, :] = img_lab[..., 0]
66     recon_const_outres[0, :, :] = img_lab_outres[..., 0]
67
68     color_ab[0, :, :] = img_lab[..., 1].reshape(1, self.shape[0],
69                                                  self.shape[1])
70     color_ab[1, :, :] = img_lab[..., 2].reshape(1, self.shape[0],
71                                                  self.shape[1])
72
73     if self.lossweights is not None:
74         weights = self.__getweights__(color_ab)
75
76     # Load feature maps
77     featobj = np.load(self.feats_fns[idx])
78     greyfeats[:, :, :] = featobj["arr_0"]
79
80     return color_ab, recon_const, weights, recon_const_outres, greyfeats
81
82 def __getweights__(self, img):
83     """
84     Calculate weight values for each pixel of an image.
85     """
86     img_vec = img.reshape(-1)
87     img_vec = img_vec * 128.0
88     img_lossweights = np.zeros(img.shape, dtype="f")
89     img_vec_a = img_vec[: np.prod(self.shape)]
90     binedges_a = self.binedges[0, ...].reshape(-1)
91     binid_a = [binedges_a.flat[np.abs(binedges_a - v).argmin()]
92               for v in img_vec_a]
93     img_vec_b = img_vec[np.prod(self.shape) :]
94     binedges_b = self.binedges[1, ...].reshape(-1)
95     binid_b = [binedges_b.flat[np.abs(binedges_b - v).argmin()]
96               for v in img_vec_b]
97     binweights = np.array([self.lossweights[v1][v2] for v1, v2 in zip(binid_a
98 , binid_b)])
99     img_lossweights[0, :, :] = binweights.reshape(self.shape[0],
100 self.shape[1])
101     img_lossweights[1, :, :] = binweights.reshape(self.shape[0], self.shape
102 [1])
103     return img_lossweights
104
105 def saveoutput_gt(self, net_op, gt, prefix, batch_size,
106                  num_cols=8, net_recon_const=None):
107     """
108     Save images
109     """
110     net_out_img = self.__tiledoutput__(net_op, batch_size, num_cols=num_cols,
111                                       net_recon_const=net_recon_const)
112     gt_out_img = self.__tiledoutput__(gt, batch_size, num_cols=num_cols,
113                                       net_recon_const=net_recon_const)

```

```

113         num_rows = np.int_(np.ceil((batch_size * 1.0) / num_cols))
114         border_img = 255 * np.ones((num_rows * self.outshape[0], 128, 3),
115                                     dtype="uint8")
116         out_fn_pred = "%s/%s.png" % (self.out_directory, prefix)
117         cv2.imwrite(out_fn_pred,
118                     np.concatenate((net_out_img, border_img, gt_out_img), axis=1)
119     )
120
121     def __tilayoutput__(self, net_op, batch_size,
122                        num_cols=8, net_recon_const=None):
123         """
124         Generate a combined image from these inputs by stitching the images into
125         a large image.
126         """
127         num_rows = np.int_(np.ceil((batch_size * 1.0) / num_cols))
128         out_img = np.zeros((num_rows*self.outshape[0], num_cols*self.outshape[1],
129                             3),
130                             dtype="uint8")
131         img_lab = np.zeros((self.outshape[0], self.outshape[1], 3),
132                             dtype="uint8")
133         c = 0
134         r = 0
135
136         for i in range(batch_size):
137             if i % num_cols == 0 and i > 0:
138                 r = r + 1
139                 c = 0
140             img_lab[... , 0] = self.__decodeimg__(net_recon_const[i, 0, :, :].
141             reshape(self.outshape[0], self.outshape[1]))
142             img_lab[... , 1] = self.__decodeimg__(net_op[i, 0, :, :].reshape(self.
143             shape[0], self.shape[1]))
144             img_lab[... , 2] = self.__decodeimg__(net_op[i, 1, :, :].reshape(self.
145             shape[0], self.shape[1]))
146             img_rgb = cv2.cvtColor(img_lab, cv2.COLOR_LAB2BGR)
147             out_img[
148                 r * self.outshape[0] : (r + 1) * self.outshape[0],
149                 c * self.outshape[1] : (c + 1) * self.outshape[1],
150                 ...,
151             ] = img_rgb
152             c = c + 1
153
154         return out_img
155
156     def __decodeimg__(self, img_enc):
157         """
158         Denormalize from [-1..1] to [0..255]
159         """
160         img_dec = (((img_enc + 1.0) * 1.0) / 2.0) * 255.0
161         img_dec[img_dec < 0.0] = 0.0
162         img_dec[img_dec > 255.0] = 255.0
163         return cv2.resize(np.uint8(img_dec), (self.outshape[0], self.outshape[1])
164     )

```

Khởi tạo các siêu tham số toàn cục cho chương trình.

```

1 # Declare hyperparameters
2 args = {
3     "gpu": 1,
4     "epochs": 2,
5     "batchsize": 32,
6     "hiddensize": 64,

```

```

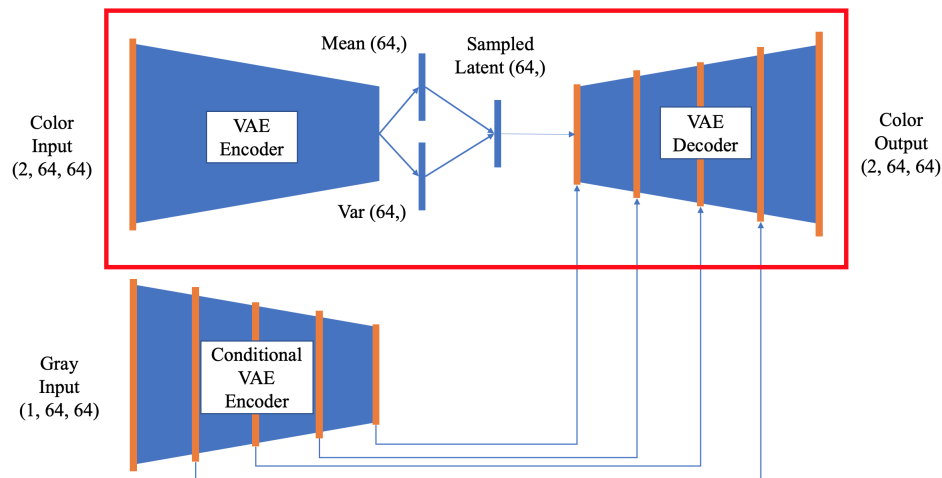
7     "nthreads": 2,
8     "epochs_mdn": 2,
9     "nmix": 8,
10    "logstep": 100,
11    "dataset_key": "lfw"
12 }
13
14 def get_dirpaths(args):
15     if args["dataset_key"] == "lfw":
16         out_dir = "data/output/lfw"
17         listdir = "data/imglist/lfw"
18         featslistdir = "data/featslist/lfw"
19     else:
20         raise NameError("[ERROR] Incorrect key: %s" % (args.dataset_key))
21     return out_dir, listdir, featslistdir

```

2. Models

Chúng ta sẽ tiến hành xây dựng mô hình VAE và mô hình MDN.

- (a) **VAE model:** Trong bài toán này, chúng ta sẽ sử dụng một biến thể của mô hình VAE, được gọi là Conditional VAE (CVAE). Mô hình này bao gồm ba phần: khối Encoder chính và khối Decoder chính (hai khối này tạo thành một mạng VAE cơ bản, được bao quanh bởi hình chữ nhật màu đỏ), cùng với một khối Conditional Encoder (khối này giúp mô hình tận dụng tối đa những trường thông tin có sẵn). Đầu vào của mạng VAE cơ bản là trường màu C có kích thước $(2 \times h \times w)$, và đầu ra là một feature map có kích thước tương tự $(2 \times h \times w)$. Đồng thời, ảnh xám G $(1 \times h \times w)$ cũng được sử dụng làm điểm khởi đầu cho khối Conditional Encoder để trích xuất các feature maps chứa thông tin cục bộ, và sau đó được sử dụng làm điều kiện để làm tăng khả năng cho khối Decoder.



Hình 4: Ảnh minh họa cho mô hình Conditional VAE.

```

1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4
5 class VAE(nn.Module):
6     def __init__(self):
7         super(VAE, self).__init__()

```

```

8         self.hidden_size = 64
9
10        # Encoder block
11        self.enc_conv1 = nn.Conv2d(2, 128, 5, stride=2, padding=2)
12        self.enc_bn1 = nn.BatchNorm2d(128)
13        self.enc_conv2 = nn.Conv2d(128, 256, 5, stride=2, padding=2)
14        self.enc_bn2 = nn.BatchNorm2d(256)
15        self.enc_conv3 = nn.Conv2d(256, 512, 5, stride=2, padding=2)
16        self.enc_bn3 = nn.BatchNorm2d(512)
17        self.enc_conv4 = nn.Conv2d(512, 1024, 3, stride=2, padding=1)
18        self.enc_bn4 = nn.BatchNorm2d(1024)
19        self.enc_fc1 = nn.Linear(4*4*1024, self.hidden_size*2)
20        self.enc_dropout1 = nn.Dropout(p=0.7)
21
22        # Conditional encoder block
23        self.cond_enc_conv1 = nn.Conv2d(1, 128, 5, stride=2, padding=2)
24        self.cond_enc_bn1 = nn.BatchNorm2d(128)
25        self.cond_enc_conv2 = nn.Conv2d(128, 256, 5, stride=2, padding=2)
26        self.cond_enc_bn2 = nn.BatchNorm2d(256)
27        self.cond_enc_conv3 = nn.Conv2d(256, 512, 5, stride=2, padding=2)
28        self.cond_enc_bn3 = nn.BatchNorm2d(512)
29        self.cond_enc_conv4 = nn.Conv2d(512, 1024, 3, stride=2, padding=1)
30        self.cond_enc_bn4 = nn.BatchNorm2d(1024)
31
32        # Decoder block
33        self.dec_upsamp1 = nn.Upsample(scale_factor=4, mode='bilinear')
34        self.dec_conv1 = nn.Conv2d(1024+self.hidden_size, 512, 3, stride=1,
padding=1)
35        self.dec_bn1 = nn.BatchNorm2d(512)
36        self.dec_upsamp2 = nn.Upsample(scale_factor=2, mode='bilinear')
37        self.dec_conv2 = nn.Conv2d(512*2, 256, 5, stride=1, padding=2)
38        self.dec_bn2 = nn.BatchNorm2d(256)
39        self.dec_upsamp3 = nn.Upsample(scale_factor=2, mode='bilinear')
40        self.dec_conv3 = nn.Conv2d(256*2, 128, 5, stride=1, padding=2)
41        self.dec_bn3 = nn.BatchNorm2d(128)
42        self.dec_upsamp4 = nn.Upsample(scale_factor=2, mode='bilinear')
43        self.dec_conv4 = nn.Conv2d(128*2, 64, 5, stride=1, padding=2)
44        self.dec_bn4 = nn.BatchNorm2d(64)
45        self.dec_upsamp5 = nn.Upsample(scale_factor=2, mode='bilinear')
46        self.dec_conv5 = nn.Conv2d(64, 2, 5, stride=1, padding=2)
47
48        def encoder(self, x):
49            x = F.relu(self.enc_conv1(x))
50            x = self.enc_bn1(x)
51            x = F.relu(self.enc_conv2(x))
52            x = self.enc_bn2(x)
53            x = F.relu(self.enc_conv3(x))
54            x = self.enc_bn3(x)
55            x = F.relu(self.enc_conv4(x))
56            x = self.enc_bn4(x)
57            x = x.view(-1, 4*4*1024)
58            x = self.enc_dropout1(x)
59            x = self.enc_fc1(x)
60            mu = x[..., :self.hidden_size]
61            logvar = x[..., self.hidden_size:]
62            return mu, logvar
63
64        def cond_encoder(self, x):
65            x = F.relu(self.cond_enc_conv1(x))
66            sc_feat32 = self.cond_enc_bn1(x)

```



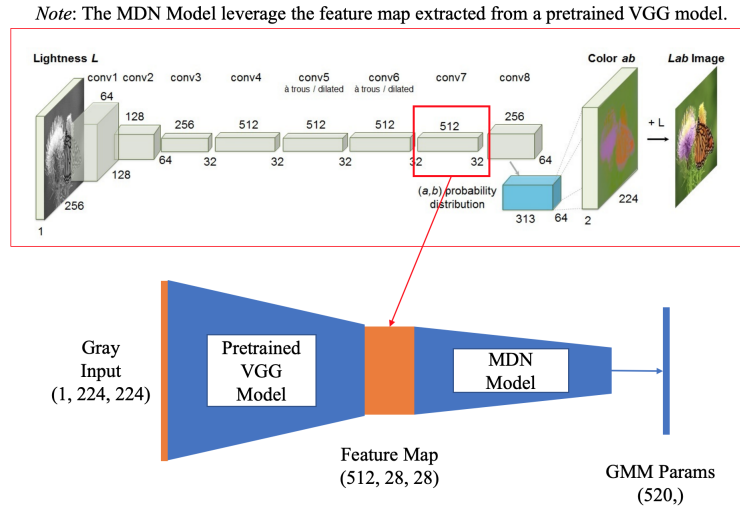
```

67         x = F.relu(self.cond_enc_conv2(sc_feat32))
68         sc_feat16 = self.cond_enc_bn2(x)
69         x = F.relu(self.cond_enc_conv3(sc_feat16))
70         sc_feat8 = self.cond_enc_bn3(x)
71         x = F.relu(self.cond_enc_conv4(sc_feat8))
72         sc_feat4 = self.cond_enc_bn4(x)
73         return sc_feat32, sc_feat16, sc_feat8, sc_feat4
74
75     def decoder(self, z, sc_feat32, sc_feat16, sc_feat8, sc_feat4):
76         x = z.view(-1, self.hidden_size, 1, 1)
77         x = self.dec_upsamp1(x)
78         x = torch.cat([x, sc_feat4], 1)
79         x = F.relu(self.dec_conv1(x))
80         x = self.dec_bn1(x)
81         x = self.dec_upsamp2(x)
82         x = torch.cat([x, sc_feat8], 1)
83         x = F.relu(self.dec_conv2(x))
84         x = self.dec_bn2(x)
85         x = self.dec_upsamp3(x)
86         x = torch.cat([x, sc_feat16], 1)
87         x = F.relu(self.dec_conv3(x))
88         x = self.dec_bn3(x)
89         x = self.dec_upsamp4(x)
90         x = torch.cat([x, sc_feat32], 1)
91         x = F.relu(self.dec_conv4(x))
92         x = self.dec_bn4(x)
93         x = self.dec_upsamp5(x)
94         x = torch.tanh(self.dec_conv5(x))
95         return x
96
97     def forward(self, color, greylevel, z_in=None):
98         sc_feat32, sc_feat16, sc_feat8, sc_feat4 = self.cond_encoder(
99             greylevel)
100         mu, logvar = self.encoder(color)
101         if self.training:
102             stddev = torch.sqrt(torch.exp(logvar))
103             eps = torch.randn_like(stddev)
104             z = mu + eps * stddev
105             z = z.to(greylevel.device)
106         else:
107             z = z_in
108             z = z.to(greylevel.device)
109         color_out = self.decoder(z, sc_feat32, sc_feat16, sc_feat8, sc_feat4)
110         return mu, logvar, color_out

```

(b) MDN model

Đầu vào của một mạng Conditional Variational Autoencoder (CVAE) yêu cầu thông tin về cả trường màu C và ảnh xám G. Trong quá trình huấn luyện, khối Encoder chính ánh xạ thông tin của trường màu C thành phân phối hậu nghiệm P, sau đó lấy mẫu từ phân phối P để làm điểm khởi đầu cho khối Decoder. Tuy nhiên, trong quá trình dự đoán, không có thông tin về trường màu C được cung cấp. Chính vì thế, một mạng MDN (Mixture Density Network) được đã được thiết kế. MDN nhận đầu vào là vector đặc trưng, được tạo ra bằng cách cho ảnh xám G đi qua mạng VGG đã được huấn luyện trước trong bài báo [Colorful Image Colorization](#). Kết quả đầu ra của mô hình MDN sau đó được sử dụng để tạo ra các tham số phân phối cho mô hình Gaussian Mixture Model, một mô hình thực hiện việc xấp xỉ phân phối P được tạo ra từ khối Encoder vừa được huấn luyện trước đó.



Hình 5: Ảnh minh họa cho mô hình MDN.

```

1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4
5 class MDN(nn.Module):
6     def __init__(self):
7         super(MDN, self).__init__()
8
9         self.feats_nch = 512
10        self.hidden_size = 64
11        self.nmix = 8
12        self.nout = (self.hidden_size + 1) * self.nmix
13
14        self.model = nn.Sequential(
15            nn.Conv2d(self.feats_nch, 384, 5, stride=1, padding=2),
16            nn.BatchNorm2d(384),
17            nn.ReLU(),
18            nn.Conv2d(384, 320, 5, stride=1, padding=2),
19            nn.BatchNorm2d(320),
20            nn.ReLU(),
21            nn.Conv2d(320, 288, 5, stride=1, padding=2),
22            nn.BatchNorm2d(288),
23            nn.ReLU(),
24            nn.Conv2d(288, 256, 5, stride=2, padding=2),
25            nn.BatchNorm2d(256),
26            nn.ReLU(),
27            nn.Conv2d(256, 128, 5, stride=1, padding=2),
28            nn.BatchNorm2d(128),
29            nn.ReLU(),
30            nn.Conv2d(128, 96, 5, stride=2, padding=2),
31            nn.BatchNorm2d(96),
32            nn.ReLU(),
33            nn.Conv2d(96, 64, 5, stride=2, padding=2),
34            nn.BatchNorm2d(64),
35            nn.ReLU(),
36            nn.Dropout(p=0.7)
37        )
38
39        self.fc = nn.Linear(4 * 4 * 64, self.nout)

```

```

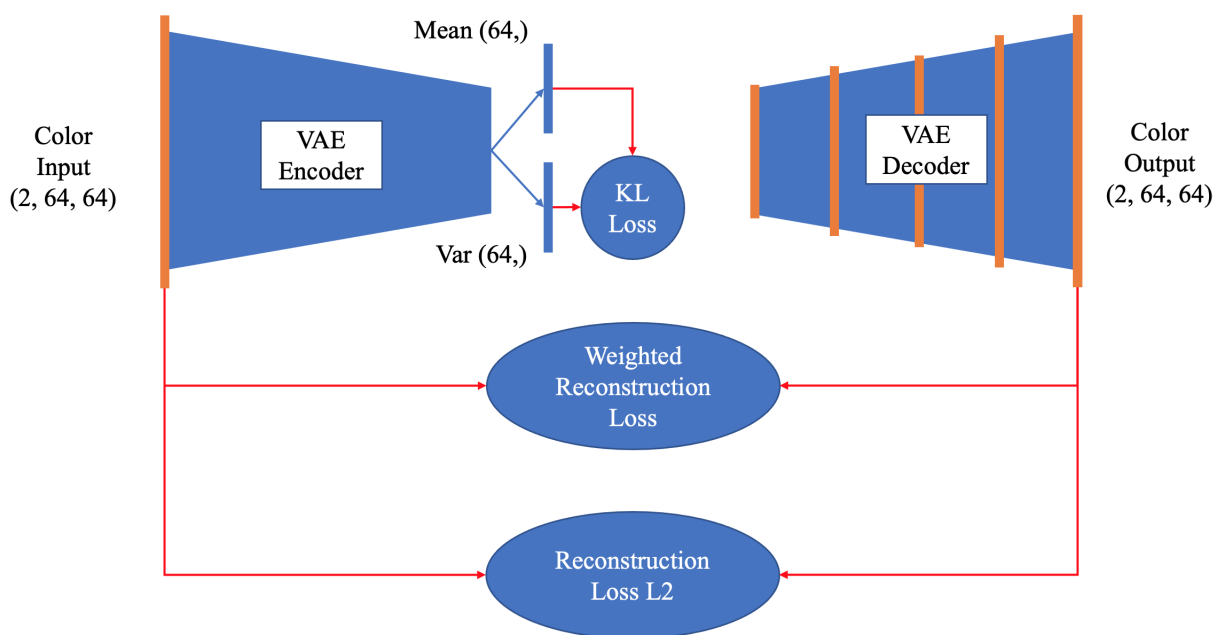
40
41     def forward(self, feats):
42         x = self.model(feats)
43         x = x.view(-1, 4 * 4 * 64)
44         x = F.relu(x)
45         x = F.dropout(x, p=0.7, training=self.training)
46         x = self.fc(x)
47         return x

```

3. Loss Functions

Trong phần này chúng ta xây dựng hàm mất mát cho các mô hình VAE và MDN.

VAE Loss



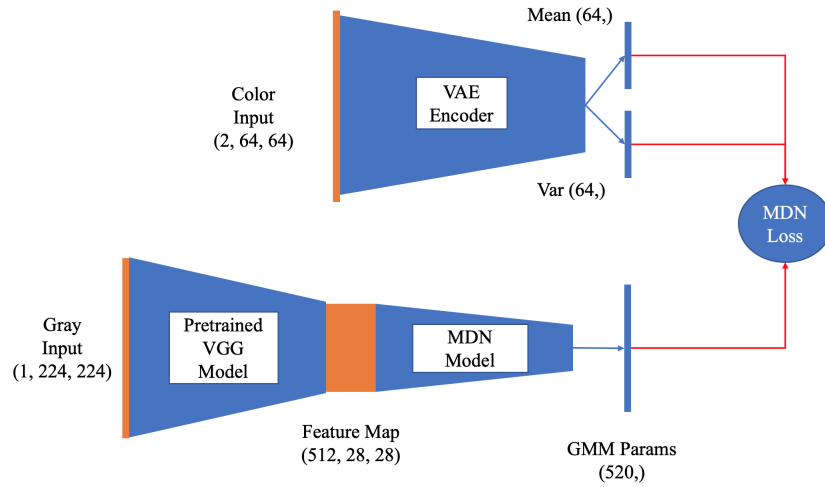
Hình 6: Ảnh minh họa cho VAE Loss.

```

1 def vae_loss(mu, logvar, pred, gt, lossweights, batchsize):
2     """
3     Return the loss values of the VAE model.
4     """
5     kl_element = torch.add(torch.add(torch.add(mu.pow(2), logvar.exp()), -1),
6                               logvar.mul(-1))
7     kl_loss = torch.sum(kl_element).mul(0.5)
8     gt = gt.reshape(-1, 64 * 64 * 2)
9     pred = pred.reshape(-1, 64 * 64 * 2)
10    recon_element = torch.sqrt(torch.sum(torch.mul(torch.add(gt, pred.mul(-1)).
11                                                       pow(2), lossweights), 1))
12    recon_loss = torch.sum(recon_element).mul(1.0 / (batchsize))
13
14    recon_element_l2 = torch.sqrt(torch.sum(torch.add(gt, pred.mul(-1)).pow(2),
15                                                       1))
16    recon_loss_l2 = torch.sum(recon_element_l2).mul(1.0 / (batchsize))
17
18    return kl_loss, recon_loss, recon_loss_l2

```

MDN Loss



Hình 7: Ảnh minh họa cho MDN Loss.

```

1 def get_gmm_coeffs(gmm_params):
2     """
3     Return the distribution coefficients of the GMM.
4     """
5     gmm_mu = gmm_params[..., : args["hiddensize"] * args["nmix"]]
6     gmm_mu.contiguous()
7     gmm_pi_activ = gmm_params[..., args["hiddensize"] * args["nmix"] :]
8     gmm_pi_activ.contiguous()
9     gmm_pi = F.softmax(gmm_pi_activ, dim=1)
10    return gmm_mu, gmm_pi
11
12 def mdn_loss(gmm_params, mu, stddev, batchsize):
13     """
14     Calculates the loss by comparing two distribution
15     - the predicted distribution of the MDN (given by gmm_mu and gmm_pi) with
16     - the target distribution created by the Encoder block (given by mu and
17     stddev).
18     """
19     gmm_mu, gmm_pi = get_gmm_coeffs(gmm_params)
20     eps = torch.randn(stddev.size()).normal_().cuda()
21     z = torch.add(mu, torch.mul(eps, stddev))
22     z_flat = z.repeat(1, args["nmix"])
23     z_flat = z_flat.reshape(batchsize * args["nmix"], args["hiddensize"])
24     gmm_mu_flat = gmm_mu.reshape(batchsize * args["nmix"], args["hiddensize"])
25     dist_all = torch.sqrt(torch.sum(torch.add(z_flat, gmm_mu_flat.mul(-1)).pow(2)
26     .mul(50), 1))
27     dist_all = dist_all.reshape(batchsize, args["nmix"])
28     dist_min, selectids = torch.min(dist_all, 1)
29     gmm_pi_min = torch.gather(gmm_pi, 1, selectids.reshape(-1, 1))
30     gmm_loss = torch.mean(torch.add(-1 * torch.log(gmm_pi_min + 1e-30), dist_min)
31     )
32     gmm_loss_l2 = torch.mean(dist_min)
33     return gmm_loss, gmm_loss_l2

```

4. Trainer

Trong phần này chúng ta xây dựng hàm huấn luyện cho từng mô hình.

Train VAE model

```

1 def test_vae(model):
2     model.eval()
3
4     # Load hyperparameters
5     out_dir, listdir, featslistdir = get_dirpaths(args)
6     batchsize = args["batchsize"]
7     hiddensize = args["hiddensize"]
8     nmix = args["nmix"]
9
10    # Create DataLoader
11    data = ColorDataset(
12        os.path.join(out_dir, "images"),
13        listdir=listdir,
14        featslistdir=featslistdir,
15        split="test",
16    )
17    nbatches = np.int_(np.floor(data.img_num / batchsize))
18    data_loader = DataLoader(
19        dataset=data,
20        num_workers=args["nthreads"],
21        batch_size=batchsize,
22        shuffle=False,
23        drop_last=True,
24    )
25
26    # Eval
27    test_loss = 0.0
28    for batch_idx, (
29        batch,
30        batch_recon_const,
31        batch_weights,
32        batch_recon_const_outres,
33        -,
34    ) in tqdm(enumerate(data_loader), total=nbatches):
35
36        input_color = batch.cuda()
37        lossweights = batch_weights.cuda()
38        lossweights = lossweights.reshape(batchsize, -1)
39        input_greylevel = batch_recon_const.cuda()
40        z = torch.randn(batchsize, hiddensize)
41
42        mu, logvar, color_out = model(input_color, input_greylevel, z)
43        _, _, recon_loss_l2 = vae_loss(
44            mu, logvar, color_out, input_color, lossweights, batchsize
45        )
46        test_loss = test_loss + recon_loss_l2.item()
47
48    test_loss = (test_loss * 1.0) / nbatches
49    model.train()
50
51    return test_loss
52
53
54 def train_vae():
55     # Load hyperparameters
56     out_dir, listdir, featslistdir = get_dirpaths(args)
57     batchsize = args["batchsize"]
58     hiddensize = args["hiddensize"]
59     nmix = args["nmix"]

```

```

60     nepochs = args["epochs"]
61
62     # Create DataLoader
63     data = ColorDataset(
64         os.path.join(out_dir, "images"),
65         listdir=listdir,
66         featslistdir=featslistdir,
67         split="train",
68     )
69     nbatches = np.int_(np.floor(data.img_num / batchsize))
70     data_loader = DataLoader(
71         dataset=data,
72         num_workers=args["nthreads"],
73         batch_size=batchsize,
74         shuffle=True,
75         drop_last=True,
76     )
77
78     # Initialize VAE model
79     model = VAE()
80     model.cuda()
81     model.train()
82
83     optimizer = optim.Adam(model.parameters(), lr=5e-5)
84
85     # Train
86     itr_idx = 0
87     for epochs in range(nepochs):
88         train_loss = 0.0
89
90         for batch_idx, (
91             batch,
92             batch_recon_const,
93             batch_weights,
94             batch_recon_const_outres,
95             -,
96         ) in tqdm(enumerate(data_loader), total=nbatches):
97
98             input_color = batch.cuda()
99             lossweights = batch_weights.cuda()
100             lossweights = lossweights.reshape(batchsize, -1)
101             input_greylevel = batch_recon_const.cuda()
102             z = torch.randn(batchsize, hiddensize)
103
104             optimizer.zero_grad()
105             mu, logvar, color_out = model(input_color, input_greylevel, z)
106             kl_loss, recon_loss, recon_loss_l2 = vae_loss(
107                 mu, logvar, color_out, input_color, lossweights, batchsize
108             )
109             loss = kl_loss.mul(1e-2) + recon_loss
110             recon_loss_l2.detach()
111             loss.backward()
112             optimizer.step()
113
114             train_loss = train_loss + recon_loss_l2.item()
115
116             if batch_idx % args["logstep"] == 0:
117                 data.saveoutput_gt(
118                     color_out.cpu().data.numpy(),
119                     batch.numpy(),

```

```

120         "train_%05d_%05d" % (epochs, batch_idx),
121         batchsize,
122         net_recon_const=batch_recon_const_outres.numpy(),
123     )
124
125     train_loss = (train_loss * 1.0) / (nbatches)
126     print("VAE Train Loss, epoch %d has loss %f" % (epochs, train_loss))
127     test_loss = test_vae(model)
128     print("VAE Test Loss, epoch %d has loss %f" % (epochs, test_loss))
129
130     # Save VAE model
131     torch.save(model.state_dict(), "%s/models/model_vae.pth" % (out_dir))
132
133     print("Complete VAE training")
134
135 train_vae()

```

Train MDN model

```

1 def test_mdn(model_vae, model_mdn):
2     # Load hyperparameters
3     out_dir, listdir, featslistdir = get_dirpaths(args)
4     batchsize = args["batchsize"]
5     hiddensize = args["hiddensize"]
6     nmix = args["nmix"]
7
8     # Create DataLoader
9     data = ColorDataset(
10         os.path.join(out_dir, "images"), listdir, featslistdir, split="test"
11     )
12     nbatches = np.int_(np.floor(data.img_num / batchsize))
13     data_loader = DataLoader(
14         dataset=data,
15         num_workers=args["nthreads"],
16         batch_size=batchsize,
17         shuffle=True,
18         drop_last=True,
19     )
20
21     optimizer = optim.Adam(model_mdn.parameters(), lr=1e-3)
22
23     # Eval
24     model_vae.eval()
25     model_mdn.eval()
26     itr_idx = 0
27     test_loss = 0.0
28
29     for batch_idx, (batch, batch_recon_const, batch_weights, _, batch_feats) in
30         tqdm(enumerate(data_loader), total=nbatches):
31
32         input_color = batch.cuda()
33         input_greylevel = batch_recon_const.cuda()
34         input_feats = batch_feats.cuda()
35         z = torch.randn(batchsize, hiddensize)
36         optimizer.zero_grad()
37
38         # Get the parameters of the posterior distribution
39         mu, logvar, _ = model_vae(input_color, input_greylevel, z)
40
41         # Get the GMM vector

```

```

42     mdn_gmm_params = model_mdn(input_feats)
43
44     # Compare 2 distributions
45     loss, _ = mdn_loss(mdn_gmm_params, mu, torch.sqrt(torch.exp(logvar)),
46                       batchsize)
47
48     test_loss = test_loss + loss.item()
49
50     test_loss = (test_loss * 1.0) / (nbatches)
51     model_vae.train()
52     return test_loss
53
54 def train_mdn():
55     # Load hyperparameters
56     out_dir, listdir, featslistdir = get_dirpaths(args)
57     batchsize = args["batchsize"]
58     hiddensize = args["hiddensize"]
59     nmix = args["nmix"]
60     nepochs = args["epochs_mdn"]
61
62     # Create DataLoader
63     data = ColorDataset(
64         os.path.join(out_dir, "images"), listdir, featslistdir, split="train"
65     )
66     nbatches = np.int_(np.floor(data.img_num / batchsize))
67     data_loader = DataLoader(
68         dataset=data,
69         num_workers=args["nthreads"],
70         batch_size=batchsize,
71         shuffle=True,
72         drop_last=True,
73     )
74
75     # Initialize VAE model
76     model_vae = VAE()
77     model_vae.cuda()
78     model_vae.load_state_dict(torch.load("%s/models/model_vae.pth" % (out_dir)))
79     model_vae.eval()
80
81     # Initialize MDN model
82     model_mdn = MDN()
83     model_mdn.cuda()
84     model_mdn.train()
85
86     optimizer = optim.Adam(model_mdn.parameters(), lr=1e-3)
87
88     # Train
89     itr_idx = 0
90     for epochs_mdn in range(nepochs):
91         train_loss = 0.0
92
93         for batch_idx, (
94             batch,
95             batch_recon_const,
96             batch_weights,
97             -,
98             batch_feats,
99         ) in tqdm(enumerate(data_loader), total=nbatches):
100             input_color = batch.cuda()

```



```

101         input_greylevel = batch_recon_const.cuda()
102         input_feats = batch_feats.cuda()
103         z = torch.randn(batchsize, hiddensize)
104         optimizer.zero_grad()
105
106         # Get the parameters of the posterior distribution
107         mu, logvar, _ = model_vae(input_color, input_greylevel, z)
108
109         # Get the GMM vector
110         mdn_gmm_params = model_mdn(input_feats)
111
112         # Compare 2 distributions
113         loss, loss_l2 = mdn_loss(
114             mdn_gmm_params, mu, torch.sqrt(torch.exp(logvar)), batchsize
115         )
116
117         loss.backward()
118         optimizer.step()
119         train_loss = train_loss + loss.item()
120
121         train_loss = (train_loss * 1.0) / (nbatches)
122         test_loss = test_mdn(model_vae, model_mdn)
123         print(
124             f"End of epoch {epochs_mdn:3d} | Train Loss {train_loss:8.3f} | Test
125             Loss {test_loss:8.3f}"
126         )
127
128         # Save MDN model
129         torch.save(model_mdn.state_dict(), "%s/models_mdn/model_mdn.pth" % (
130             out_dir))
131
132         print("Complete MDN training")
133
134 train_mdn()

```

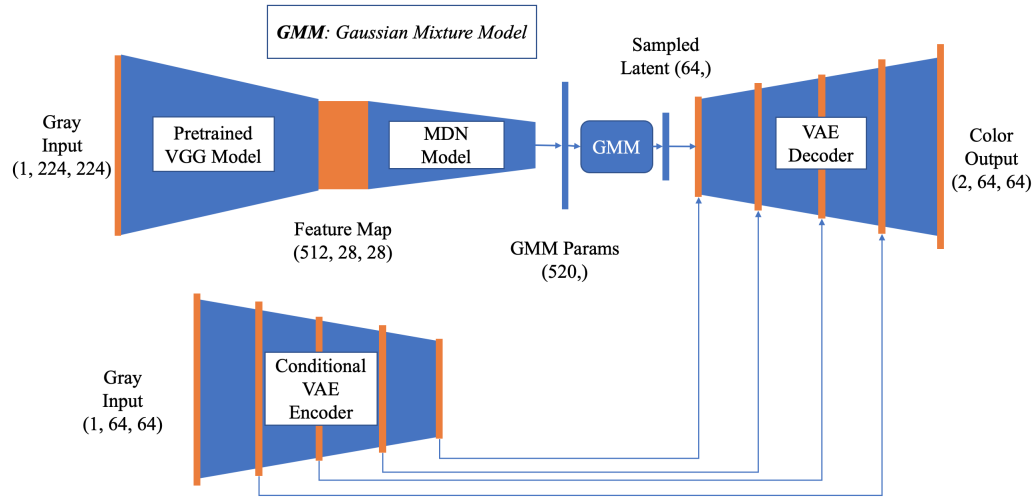
5. Inference

Bạn có thể sử dụng checkpoint sẵn có để tiến hành quá trình suy luận thử nghiệm.

```

1 # Download VAE checkpoint
2 !gdown 1wdyK198lXwwZ04NIB7DzJmA5arwUVWUDU
3 # Download MDN checkpoint
4 !gdown 1AhiLMrR_C04v7_sysuf5ffEVsQ1lo2W6

```



Hình 8: Ảnh minh họa cho Big Model ở giai đoạn inference

```

1 def inference():
2     # Load hyperparameters
3     out_dir, listdir, featslistdir = get_dirpaths(args)
4     batchsize = args["batchsize"]
5     hiddensize = args["hiddensize"]
6     nmix = args["nmix"]
7
8     # Create DataLoader
9     data = ColorDataset(
10         os.path.join(out_dir, "images"),
11         listdir=listdir,
12         featslistdir=featslistdir,
13         split="test",
14     )
15
16     nbatches = np.int_(np.floor(data.img_num / batchsize))
17
18     data_loader = DataLoader(
19         dataset=data,
20         num_workers=args["nthreads"],
21         batch_size=batchsize,
22         shuffle=False,
23         drop_last=True,
24     )
25
26     # Load VAE model
27     model_vae = VAE()
28     model_vae.cuda()
29     model_vae.load_state_dict(torch.load("%s/models/model_vae.pth" % (out_dir)))
30     model_vae.eval()
31
32     # Load MDN model
33     model_mdn = MDN()
34     model_mdn.cuda()
35     model_mdn.load_state_dict(torch.load("%s/models/model_mdn.pth" % (out_dir)))
36     model_mdn.eval()
37
38     # Infer
39     for batch_idx, (

```

```

40     batch,
41     batch_recon_const,
42     batch_weights,
43     batch_recon_const_outres,
44     batch_feats,
45 ) in tqdm(enumerate(data_loader), total=nbatches):
46
47     input_feats = batch_feats.cuda()
48
49     # Get GMM parameters
50     mdn_gmm_params = model_mdn(input_feats)
51     gmm_mu, gmm_pi = get_gmm_coeffs(mdn_gmm_params)
52     gmm_pi = gmm_pi.reshape(-1, 1)
53     gmm_mu = gmm_mu.reshape(-1, hiddensize)
54
55     for j in range(batchsize):
56         batch_j = np.tile(batch[j, ...].numpy(), (batchsize, 1, 1, 1))
57         batch_recon_const_j = np.tile(
58             batch_recon_const[j, ...].numpy(), (batchsize, 1, 1, 1)
59         )
60         batch_recon_const_outres_j = np.tile(
61             batch_recon_const_outres[j, ...].numpy(), (batchsize, 1, 1, 1)
62         )
63
64         input_color = torch.from_numpy(batch_j).cuda()
65         input_greylevel = torch.from_numpy(batch_recon_const_j).cuda()
66
67         # Get mean from GMM
68         curr_mu = gmm_mu[j * nmix : (j + 1) * nmix, :]
69         orderid = np.argsort(
70             gmm_pi[j * nmix : (j + 1) * nmix, 0].cpu().data.numpy().reshape
71             (-1)
72         )
73
74         # Sample from GMM
75         z = curr_mu.repeat(int((batchsize * 1.0) / nmix), 1)
76
77         # Predict color
78         _, _, color_out = model_vae(input_color, input_greylevel, z)
79
80         # Save images
81         data.saveoutput_gt(
82             color_out.cpu().data.numpy()[orderid, ...],
83             batch_j[orderid, ...],
84             "divcolor_%05d_%05d" % (batch_idx, j),
85             nmix,
86             net_recon_const=batch_recon_const_outres_j[orderid, ...],
87         )
88
89     print("Complete inference")
90
91 vae_ckpt = "model_vae.pth"
92 mdn_ckpt = "model_mdn.pth"
93 inference(vae_ckpt, mdn_ckpt)

```

Kết quả thực nghiệm mô hình sau khi huấn luyện



Hình 9: Kết quả thực nghiệm mô hình sau khi huấn luyện.

Phần III: Câu hỏi trắc nghiệm

1. VAE có thể được sử dụng trong các ứng dụng nào?
 - (a) Tô màu cho ảnh xám
 - (b) Nén ảnh
 - (c) Sinh ảnh mới
 - (d) Tất cả các phương án trên
2. Mô hình VAE cơ bản có bao nhiêu block chính?
 - (a) 1
 - (b) 2
 - (c) 3
 - (d) 4
3. Trong VAE, đối tượng cần được mã hóa được biểu diễn như thế nào?
 - (a) Dưới dạng một giá trị số thực.
 - (b) Dưới dạng một phân phối xác suất
 - (c) Dưới dạng một véc-tơ nhị phân
 - (d) Dưới dạng một ma trận đặc trưng
4. Trong VAE, khi huấn luyện mô hình, ta muốn KL Divergence Loss đạt giá trị bằng bao nhiêu?
 - (a) 0
 - (b) 1
 - (c) Không có giá trị nhất định
 - (d) Càng lớn càng tốt
5. Trong Image Colorization, vì sao không gian màu Lab thường được sử dụng hơn không gian màu RGB (chọn phương án đúng nhất)?
 - (a) **Phân biệt rõ ràng giữa độ sáng và màu sắc:** Không gian màu Lab phân chia màu sắc và độ sáng thành hai kênh riêng biệt (L, a, và b), giúp mô hình tập trung vào việc tái tạo màu sắc một cách chính xác hơn. Trong khi đó, ảnh RGB có thể làm mất thông tin về độ sáng khi thêm màu vào, gây ra hiệu ứng không mong muốn.
 - (b) **Khả năng lưu trữ thông tin màu sắc chi tiết:** Các giá trị trong không gian màu Lab là liên tục, ngược lại với các giá trị rời rạc trong không gian màu RGB. Điều này dẫn đến việc không gian màu Lab có khả năng lưu trữ lớn hơn và chính xác hơn trong việc biểu diễn hình ảnh
 - (c) **Độ phức tạp thấp hơn khi dự đoán màu sắc:** Với không gian màu Lab, chúng ta chỉ cần dự đoán hai kênh màu a và b, thay vì cả ba kênh màu như trong không gian màu RGB. Điều này giúp giảm độ phức tạp của bài toán và tăng hiệu suất của mô hình.
 - (d) Tất cả đáp án trên.
6. Trong quá trình inference, mô hình CVAE có sự tham gia của những thành phần nào?
 - (a) Encoder, Conditional Encoder, Decoder
 - (b) Encoder, Decoder

(c) Conditional Encoder, Decoder

(d) Decoder

- *Hết* -