

SIGK2025Z

Projekt 2

Mapowanie tonalne

Oskar Nowak i Paulina Staszewska

17 listopada 2025

1 Cel projektu

Celem projektu było stworzenie mapowania tonalnego, które wykorzystuje sieci neuronowe oraz przeprowadzenie nad nim eksperymentów.

2 Zbiór danych

Zbiór danych to SI-HDR, w którym wykorzystano pliki z folderu *reference* oraz odpowiadające im pliki z folderu *clip_97* według nazwy. Zbiór został podzielony na następujące podzbiory:

- treningowy (70%, czyli 127 plików z każdego folderu),
- walidacyjny (10%, czyli 18 plików z każdego folderu),
- testowy (20%, czyli 36 plików z każdego folderu).

W obu folderach było w sumie po 181 plików.

3 Model

Architektura sieci składa się z 3 koderów, fusion blocku i dekodera (listing 1). Każdy z tych elementów posiada:

- kodery:
 - wejściem do każdego z kodera to obraz 3 kanałowy
 - 3 warstwy konwolucji z *kernel* i *paddingem* równymi kolejno 3 i 1,
 - po każdej tych warstw użyta jest funkcja ReLU,
 - w ostatniej warstwie otrzymywana jest mapa cech z ilością kanałów równą 64
- fusion block:
 - łączy wyjścia koderów w jedno,
 - zmniejsza liczbę kanałów do 64 (krok dodany dodatkowo ze względu na zbyt duże wykorzystanie pamięci RAM w systemie)
 - 2 warstwy konwolucji
- dekoder:
 - 3 warstwy konwolucji z *paddingem* równymi 1,
 - po każdej tych warstw użyta jest funkcja ReLU, prócz ostatniej, po której występuje funkcja sigmoid

- ilość kanałów odpowiada ilości kanałów z kodera w odwróconej kolejności
- wyjściem dekodera jest obraz o 3 kanałach

```

1  # -----
2  class Encoder(nn.Module):
3      def __init__(self):
4          super().__init__()
5          self.conv1 = nn.Conv2d(3,16,3,padding=1)
6          self.conv2 = nn.Conv2d(16,32,3,padding=1)
7          self.conv3 = nn.Conv2d(32,64,3,padding=1)
8      def forward(self,x):
9          x = F.relu(self.conv1(x))
10         x = F.relu(self.conv2(x))
11         x = F.relu(self.conv3(x))
12         return x
13
14 class Decoder(nn.Module):
15     def __init__(self):
16         super().__init__()
17         self.conv1 = nn.Conv2d(64,32,3,padding=1)
18         self.conv2 = nn.Conv2d(32,16,3,padding=1)
19         self.conv3 = nn.Conv2d(16,3,3,padding=1)
20     def forward(self,x, img1,img2,img3):
21         x = F.relu(self.conv1(x))
22         x = F.relu(self.conv2(x))
23         x = self.conv3(x)
24         x = torch.sigmoid(x + img1 + img2 + img3)
25         return x
26
27 class TMONet(nn.Module):
28     def __init__(self):
29         super().__init__()
30         self.encoder = Encoder()
31         self.decoder = Decoder()
32         # fusion convs
33         self.fusion_conv1 = nn.Conv2d(64*3, 64, 3, padding=1) # zmniejszone ze
34             # wzgl du na zbyt duze uzycie RAM w porowaniu do oryginalnej architektury
35         self.fusion_conv2 = nn.Conv2d(64, 64, 1)
36
37     def forward(self, i0, i1, i2):
38         # encode each input
39         o0 = self.encoder(i0)
40         o1 = self.encoder(i1)
41         o2 = self.encoder(i2)
42
43         # concat features
44         o0 = torch.cat([o0, o1, o2], dim=1) # nadpisywanie dla optymalizacji
45             # wykorzystania RAM
46         o1 = o2 = None
47
48         # fusion block
49         o0 = F.relu(self.fusion_conv1(o0))
50
51         o0 = self.fusion_conv2(o0)
52
53         # decode
54         o0 = self.decoder(o0, i0, i1, i2)
55         return o0

```

Listing 1: Architektura sieci

3.1 Funkcja straty

Sieć korzysta z funkcji straty o nazwie *fcm_loss*, która korzysta z *feature contrast masking loss*, *VGG*, lokalnej średniej i odchyłonej standardowej oraz maskowania (listing 2). Strata ta bierze pod uwagę percepcję wzrokowej człowieka.

```

1  # -----
2  # VGG feature extractor

```

```

3 vgg_model = models.vgg19(pretrained=True).features.to(device).eval()
4 for p in vgg_model.parameters():
5     p.requires_grad = False
6
7 def vgg_features(x, layers=['0','5','10']): # odpowiada VGG11,21,31 w TF
8     features = []
9     h = x
10    for idx, layer in enumerate(vgg_model):
11        h = layer(h)
12        if str(idx) in layers:
13            features.append(h)
14    return features # lista tensorow
15
16 # -----
17 # Gaussian kernel
18 _gaussian_cache = {} # optymalizacja aby raz wczytac
19 def gaussian_kernel(size=13, sigma=2.0, channels=3, device='cpu'):
20     key = (size, sigma, channels)
21     if key in _gaussian_cache:
22         return _gaussian_cache[key].to(device)
23     coords = np.arange(size) - size//2
24     x, y = np.meshgrid(coords, coords)
25     kernel = np.exp(-(x**2 + y**2)/(2*sigma**2))
26     kernel = kernel / kernel.sum()
27     kernel = torch.tensor(kernel, dtype=torch.float32, device=device)
28     kernel = kernel.view(1,1,size,size).repeat(channels,1,1,1)
29     _gaussian_cache[key] = kernel
30     return kernel
31
32 def local_mean_std(x, kernel_size=13, sigma=2.0):
33     C = x.shape[1]
34     w = gaussian_kernel(kernel_size, sigma, C).to(x.device)
35     x_pad = F.pad(x, (kernel_size//2,)*4, mode='reflect')
36     mean_local = F.conv2d(x_pad, w, groups=C)
37     mean_sq = F.conv2d(x_pad**2, w, groups=C)
38     std_local = torch.sqrt(torch.clamp(mean_sq - mean_local**2, min=1e-8))
39     return mean_local, std_local
40
41 # -----
42 # Feature contrast masking
43 def sign_num_den(x, gamma=0.5, beta=0.5, sigma=2.0, kernel_size=13):
44     mean_local, std_local = local_mean_std(x, kernel_size, sigma)
45     norm_num = torch.sign(x - mean_local) * torch.abs((x - mean_local)/(mean_local.abs()
46     ()+1e-8))**gamma
47     norm_den = 1.0 + (std_local / (mean_local.abs()+1e-8))**beta
48     return norm_num, norm_den
49
50 def feature_contrast_masking(x, gamma=0.5, beta=0.5, sigma=2.0, kernel_size=13):
51     num, den = sign_num_den(x, gamma, beta, sigma, kernel_size)
52     return num / den
53
54 def masking_loss(pred, target, gamma=0.5, beta=0.5, sigma=2.0, kernel_size=13):
55     f_pred = feature_contrast_masking(pred, gamma=1.0, beta=beta, sigma=sigma,
56                                         kernel_size=kernel_size)
57     f_target = feature_contrast_masking(target, gamma=gamma, beta=beta, sigma=sigma,
58                                         kernel_size=kernel_size)
59     return F.l1_loss(f_pred, f_target)
60
61 # -----
62 # FCM loss
63 def fcm_loss(pred, target, gamma=0.5, beta=0.5, sigma=2.0, kernel_size=13):
64     feats_pred = vgg_features(pred, layers=['0','5','10']) # odpowiada VGG11,21,31
65     with torch.no_grad():
66         feats_target = vgg_features(target, layers=['0','5','10'])
67         loss_total = 0.0
68         for f_pred, f_target in zip(feats_pred, feats_target):
69             loss_total += masking_loss(f_pred, f_target, gamma=gamma, beta=beta, sigma=
70                                         sigma, kernel_size=kernel_size)
71         loss_total /= len(feats_pred)

```

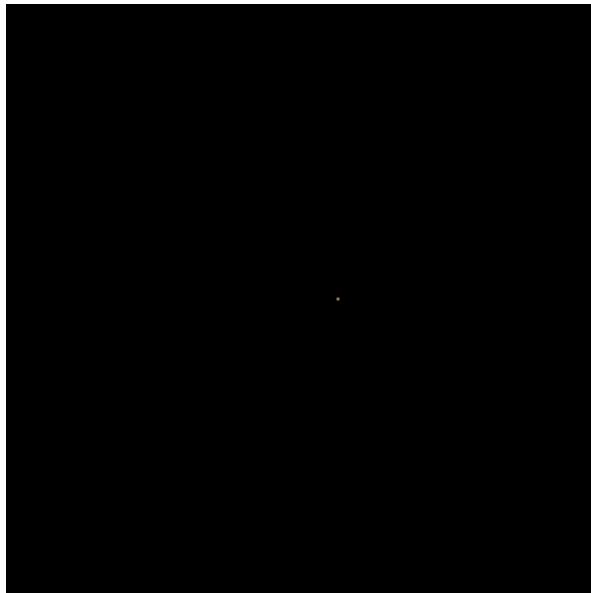
```
68     return loss_total
```

Listing 2: Funkcja straty

Model został wytrenowany według następujących parametrów:

- rozmiar batcha: 1,
- ilość epok: 49,
- funkcja straty: *fcm_loss*,
- optymalizator: Adam
- learning rate: 2^{-5}
- checkpoint: *model_11_17_16_1_14_48.pth*.

4 Wizualizacja wyników



Rysunek 1: Obraz 156_Mantiuk.png otrzymany metodą Mantiuk

Na podstawie zdjęcia z folderu *clip_97* (rysunek 4) można zauważyc, że obrazy otrzymane metodą Mantiuk (rysunek 1) jest wręcz czarny. Natomiast obrazy otrzymane przez sieć neurnową (rysunek 3) i metodę Reinhard (rysunek 2) są prawie podobne do siebie.

5 Wyniki

Ze zbioru wybrano 20% plików i na podstawie ich obliczono wartości metryk, które zostały zaprezentowane w tabeli 1. Z niej warto podkreślić, że wynik idealny równy 1 dla SSIM dla metody *input_97* to

Metoda	SSIM	BRISQUE
<i>input_97</i>	1,000	18,757
<i>Reinhard_Operator</i>	0,839	37,427
<i>Mantiuk_Operator</i>	0,779	82,894
<i>Learned_Operator</i>	0,661	19,747

Tabela 1: Metryki



Rysunek 2: Obraz 156_Reinhard.png otrzymany metodą Reinhard

tak naprawdę punkt odniesienia na podstawie, którego wykonano porównanie otrzymanych obrazów SDR, które zostały zmapowane z HDR. Według współczynnika SSIM najlepszą metodą jest metoda Reinhard, która uzyskała wartość równą 0,839. Kolejno najlepszą metodą była wykorzystująca sposobu Mantiuk. Model z siecią neuronową uzyskał najgorszy wynik z pośród wszystkich sposobów, a patrząc jedynie w obrębie współczynnika SSIM jest to wynik średni. Natomiast według metryki BRISQUE najlepszy wynik osiągnęła metoda wykorzystująca sieć neuronową. Średnią metodą jest metoda Reinharda. Najgorszy wynik otrzymała metoda Mantiuk. Biorąc pod uwagę metrykę BRISQUE to model z siecią neuronową otrzymuje wartość zbliżoną do baseline.

Wszystkie otrzymane zdjęcia z modelu w fazie treningu zapisano do foldera o nazwie *results*, a w folderze o nazwie *results/TEST* są pliki stworzone w czasie testowania. W folderze o nazwie *compare* są zdjęcia otrzymane z pozostałych metod oraz pliki HDR. Foldery te są umieszczone w folderze dla danego modelu, którego one dotyczą.

6 Eksperymenty

W projekcie początkowo wersja modelu podczas treningu miała następujące parametry:

- rozmiar batcha: 1,
- ilość epok: 50,
- funkcja straty: *fcm_loss*,
- optymalizator: Adam
- learning rate: 2^{-4}
- checkpoint: *model_11_16_4_25_47_49.pth*.

Finalna wersja modelu posiada learning rate równe 2^{-5} ze względu na to, że większa wartość powodowała zbyt duże wahania w otrzymywaniu wag sieci, a cały trening był niestabilny.

Pierwotna wersja modelu nie wykonywała ona zmniejszania obrazów do 255 na 255 z interpolacją ustawioną na *INTER_AREA*. Osiągnięte wartości metryk tego modelu zostały umieszczone w tabeli 2.

Na podstawie 2 wersji modelu można powiedzieć, że ostateczna wersja modelu uzyskuje lepsze wyniki pod względem współczynnika BRISQUE. Natomiast, biorąc pod uwagę parametr SSIM to dla metody Mantiuk wypada jedynie lepiej nowsza wersja modelu niż starsza. Jednak oryginalny model uzyskał wartość BRISQUE znacznie gorszą od swojej nowszej wersji i od baseline.



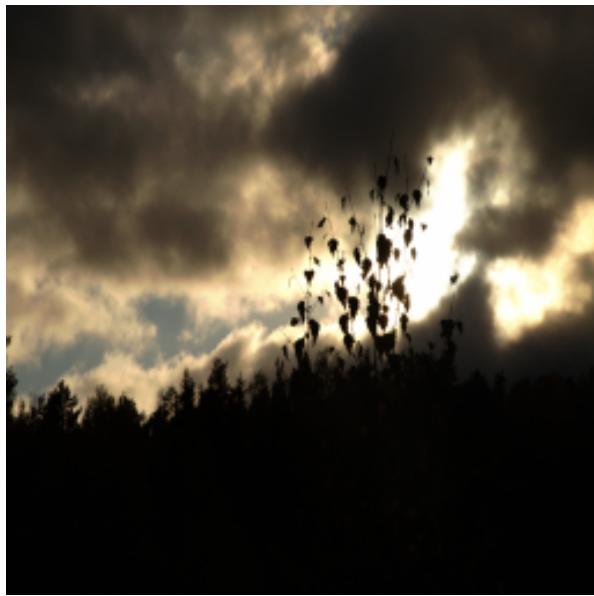
Rysunek 3: Obraz 156.png otrzymany wyuczonym modelem

Metoda	SSIM	BRISQUE
<i>input_97</i>	1,000	19,657
<i>Reinhard_Operator</i>	0,716	55,721
<i>Mantiuk_Operator</i>	0,837	103,632
<i>Learned_Operator</i>	0,633	57,798

Tabela 2: Metryki

7 Wnioski

Zaimplementowana sieć neuronowa pod względem metryki BRISQUE daje najlepszy rezultat, prawie wręcz porównywalny z orginałem pod względem precepcji widzenia człowieka. Jednak, jeśli chodzi o SSIM to model zwraca obrazy, które są średniej jakości m. in. pod względem struktury i szumu. Finalna wersja modelu jest lepsza niż jej pierwówzór w kwesti mapowania pliku HDR na SDR oraz w procesie trenowania.



Rysunek 4: Obraz 156_SDR_input97.png - baseline