
Edge Enhanced Three-Stage Deep Network for Image Super-Resolution

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053
054

Farhan Sadik¹

Abstract

Single Image Super-Resolution (SISR) poses an ill-posed problem, involving the generation of high-resolution images from low-resolution input. This implies a scenario with a higher number of variables to address compared to the equations available for solving. Historically, SISR challenges have been addressed through interpolation techniques. However, the advent of deep learning has significantly enhanced super-resolution performance. Data-driven approaches leverage multiple convolution layers to acquire feature representations from diverse images, enabling more accurate interpolation of low-resolution data. However, in conventional deep learning methods, low-resolution images contain high-frequency information that is uniformly handled across all channels. As a result, these techniques do not fully take advantage of the high-frequency information available. In this paper, I reimplemented a deep learning-based architecture that exploits high-frequency information in a separate stage and incorporated that information as a loss constraint. The entire structure consists of a dual-stage design, with one stage dedicated to an end-to-end network for SISR and the other dedicated to retaining high-frequency information. Moreover, I incorporated a third stage, where the edge information of the image is extracted through the Sobel kernel and incorporated as a loss to boost the overall end-to-end performance. The proposed method achieved a PSNR (Peak Signal-to-Noise Ratio) of 30.39 dB on the DIV2K dataset. Moreover, Incorporating the Sobel loss also increased the PSNR to 30.87 dB.

1. Introduction

Reducing an image to its low resolution can greatly enhance data transfer efficiency. However, at the user end, it needs to be reconstructed to high resolution. Such techniques find applications in diverse fields. One such area is MRI acquisition which involves reconstructing high-resolution data from low-resolution noisy images (Yang et al., 2018). Another example can be smartphone cameras where hardware might have limitations in resolution, but post-processing super-resolution techniques can significantly enhance image quality (Chen et al., 2019).

In this paper, I reimplemented the full work of (Han et al., 2021), where a dual-stage deep neural network is incorporated for SISR task. A two-stage network is employed where the first stage is an end-to-end network that converts a low-resolution (LR) input to an LR output with high-frequency information. Specifically, residual channel attention blocks are integrated to perform this task, i.e., retaining high-frequency information. During the second stage, the LR image, containing high-frequency information, is transformed into a high-resolution (HR) image while preserving the high-frequency details. A dual-stage loss is employed to collectively acquire knowledge about high-frequency information across both stages. The authors did not supply the code, prompting me to re-implement it entirely. Additionally, I expanded on the work by introducing an extra stage that enhanced edge information through the application of a loss constraint based on the Sobel operator.

2. Review of paper to implement or extend

2.1. Storyline

High-level motivation/problem Single Image Super Resolution is an intriguing field, as it involves an ill-posed problem where you have more variables than equations to solve for (Yang et al., 2019). In this particular challenge, the task is to reconstruct a high-quality image, which involves a larger number of pixels (variables) from low-resolution images, where there are fewer pixels (equations) to guide the process. Due to limitations in hardware capabilities, most of the information is transmitted after compressing the data. This expedited data transfer comes at the cost of losing some data. Fast and accurate Single Image Super Resolution

¹Purdue University, West Lafayette, Indiana, USA. Correspondence to: Farhan Sadik <fsadik@purdue.edu>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

(SISR) facilitates quicker transmission and ensures a precise portrayal of the image at the user's end.

Prior work on this problem Existing Single Image Super Resolution (SISR) methods can be classified into three main categories: interpolation-based approaches (Chang et al., 2004; Bevilacqua et al., 2012), reconstruction-based algorithms (Chakrabarti et al., 2007), and learning-based techniques (Glasner et al., 2009; Huang et al., 2015; Freedman & Fattal, 2011; Johnson et al., 2016). Both interpolation-based techniques and reconstruction-based algorithms experience reduced efficiency and increased sensitivity to scaling factors. Deep Learning (DL), as a subset of learning-based methods, grasps the nonlinear relationship between low-resolution input and high-resolution output using convolutional neural networks (CNNs). Don et al. (Dong et al., 2016) introduced the Super Resolution CNN (SRCNN), which utilizes basic CNN operations for extraction, non-linear mapping, and reconstruction of images in alignment with sparse coding (Yang et al., 2010) to learn high-resolution pixels. More intricate networks such as LapSRN (Lai et al., 2017), DRRN (Tai et al., 2017), SR-ResNet (Ledig et al., 2017), EDSR (Lim et al., 2017) and RCAN (Zhang et al., 2018b). were employed to reconstruct superior-quality images, involving the utilization of large models with increased parameters to grasp complex non-linear mappings. In recent times, attention-based networks have been employed in super-resolution tasks such as non-local attention learning (Zhang et al., 2019) which helps to retain low-level features.

Research gap Nonetheless, blindly enhancing the complexity and depth of the network marginally improves performance, but the associated computational costs may hinder many potential applications (Han et al., 2021). Furthermore, a majority of CNN-based approaches overlook the significance of high-frequency features (Han et al., 2021). Consequently, many of the reconstructed images exhibit significant blurriness (Han et al., 2021).

Contributions In this paper (Han et al., 2021), a Two-Stage Network (TSN) is proposed to reconstruct high-resolution images from low-resolution images by exploiting the high-frequency information present in the high-resolution images. The proposed TSN in this paper captures high-frequency details from high-resolution images and then learns to reconstruct high-resolution images from low-resolution inputs.

2.2. Proposed solution

To manage the network's complexity while also extracting high-frequency details, this paper proposes the use of Residual Channel Attention Groups (RCAGs) (Zhang

et al., 2018b) to obtain a high-frequency low-resolution image from a high-frequency high-resolution image. In the first stage of the network, the low-resolution image is transformed into a low-resolution image that retains high-frequency information. Subsequently, in the second stage, this low-resolution, high-frequency image is further transformed into a high-resolution image with enhanced high-frequency details. Additionally, a two-stage learning loss is introduced to guide the super-resolution process by simultaneously controlling the learning in both stages.

The first stage is divided into two simultaneous operations. Firstly, several RCAGs followed by a downscale module (collectively they named it H_{down}) are used to transform the high-resolution images (I_{HR}) to low-resolution images with high-frequency information ($I_{LR'}$) which will be later used for calculating the loss.

$$I_{LR'} = H_{down}(I_{HR}) \quad (1)$$

Secondly, the network is initiated by a convolution layer (H_{SF}) to extract shallow feature (F_o) from the input LR image (I_{LR}).

$$F_o = H_{SF}(I_{LR}) \quad (2)$$

This shallow feature, F_o is passed through four consecutive RCAGs (H_{RCAG4}) to obtain the intermediate low-resolution high-frequency image (F_{HF}).

$$F_{HF} = H_{RCAG4}(F_o) \quad (3)$$

In the second stage, the low-resolution high-frequency image (F_{HF}) is passed through four consecutive RCAGs to extract more deep features (F_{DF}). Finally, in the reconstruction layer, F_{DF} is upsampled to get the super-resolution image (I_{SR}).

$$F_{DF} = H_{RCAG8}(F_{HF}) \quad (4)$$

$$I_{SR} = H_{\uparrow}(F_{DF}) \quad (5)$$

$I_{LR'}$ from Equation (1) is used to calculate the joint loss given by

$$L_{\Theta} = \frac{1}{N} \sum_{n=1}^N \lambda \|I_{HR} - I_{SR}\|_1 + \|I_{HF} - I_{LR'}\|_1 \quad (6)$$

where, $\theta \in \Theta$ are the network parameters and I_{HR} is the high resolution ground truth image.

Table 1. Effects of different modules. We report the PSNR on Set5 datasets in the 200 epoch.

Method	R_a	R_b	R_c	R_d	R_e
Two-stage		✓	✓	✓	✓
RCAG			✓	✓	✓
$\lambda=1$		✓	✓		
$\lambda=5$				✓	
$\lambda=10$					✓
PSNR	32.23	32.45	33.16	33.23	33.21

Figure 1. Table 1 is taken from page 5, section 4.2 of the original paper (Han et al., 2021) as a screenshot.

2.3. Claims-Evidence

Claim 1 The introduction of RCAG to transform the high-resolution space into a low-resolution space led to a significant improvement in the network’s performance.

Evidence 1 Table 1 from the paper shows that after incorporating the loss where RCAG is used to calculate $I_{LR'}$, the performance increased by 1 dB PSNR (R_c in Figure 1).

Claim 2 The paper criticized existing methods because of higher complexity with marginal performance and proposed a model with low parameters but high performance.

Evidence 2 Table 3 (Figure 2) from the paper represents that TSN has significantly lower parameters but its performance is better than the existing state-of-the-art methods.

Table 3. Computational and parameter comparison (2X) Set5.

	EDSR	MemNet	MSRN	SeaNet	TSN
Para.	43M	677k	6M	8M	9.6M
PSNR	38.11	37.78	38.07	38.08	38.19

Figure 2. Table 3 is taken from page 5, section 4.4 of the original paper (Han et al., 2021) as a screenshot.

Claim 3 Methods described in existing literature tend to overlook high-frequency details, leading to blurry artifacts in the reconstructed images. The method presented integrates high-frequency information into both the network architecture and the loss function, resulting in higher-quality images.

Evidence 3 As illustrated in Figure 2 from the paper (Figure 3), the proposed TSN exhibits notably superior sharpness compared to other methods at a 4x scaling factor. The zoomed-in image of the selected section distinctly highlights the inability of most of the super-resolution methods to capture sharp details.

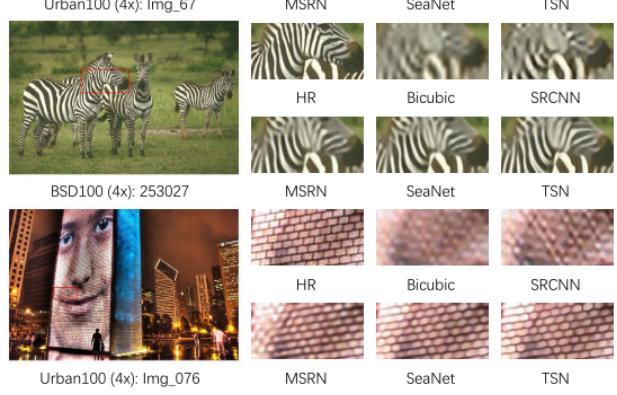


Figure 2. Visual comparison for 4x SR with BI model.

Figure 3. Figure 2 is taken from page 6, section 4.3 of the original paper (Han et al., 2021) as a screenshot.

2.4. Critique and Discussion

The paper is skillfully written, with a particularly impactful storytelling element that captivates its readers. The most enlightening aspect was their assertion regarding the effectiveness of employing attention blocks, substantiated through experimental evidence. Moreover, the claims made throughout the papers are justified through proper tables or figures.

However, the paper suffers from several issues. First of all, the implementation details of this paper are vague. The authors did not mention the number of epochs to train, and the training time, so that the reader could expect the training time while re-implementing their work. Nonetheless, the paper also failed to provide a proper justification for some elements. For example, the authors used 15 Residual Channel Attention Blocks without offering any rationale for their choice of this particular quantity.

Apart from that, L_1 loss is used for optimization instead of L_2 , which is more common in vision applications. They could provide experimental results showing why L_1 loss is better than L_2 loss in this case.

3. Review of 2nd paper

3.1. Storyline

High-level motivation/problem Single Image Super Resolution (SISR) refers to the process of generating a high-resolution image from a low-resolution input. The problem is ill-posed as one low-resolution image can have multiple high-resolution solutions. Nonetheless, the majority of Single Image Super Resolution (SISR) techniques demand significant memory resources, leading to subpar performance which hinders their potential for commercialization. This paper aims to work on putting emphasis on second-order

channel attention which is different from all the previous works as well as focuses on improving the visual quality of the images.

Prior work on this problem The strong capability of Convolutional Neural Networks (CNNs) in representing features has made them widely utilized in single-image super-resolution tasks. CNNs function as a method to map features from low-resolution input images to high-resolution output images. While CNNs (Zhang et al., 2018a), (Zhang et al., 2018b), (Zhang et al., 2018c) have attained cutting-edge performance in Single-Image Super-Resolution (SISR), they encounter issues such as generating redundant information or hallucinations. Nonetheless, the reconstruction performance of the CNNs is a great improvement over the previous interpolation methods such as non-local similarity prior (Zhang et al., 2012) and sparsity prior (Dong et al., 2011).

Research gap Most of the CNN-based SISR methods incorporate deeper networks for unrolling the images but that in turn avoids feature relation between intermediate layers. Therefore, these methods do not fully exploit the information that lies within the low-resolution images which leads to poor performance. Moreover, deepening the layers only tends to learn the high-level features hindering the representational power of CNNs. Nevertheless, transitioning away from CNNs isn't practical, as interpolation-based methods also require extensive time for optimization. Additionally, these techniques quickly become ineffective when the model processes input data with different statistics than the prior dataset.

Contributions In this work (Dai et al., 2019), the author introduces a second-order attention network (SAN) for more robust feature expression and feature correlation between intermediate layers. Specifically, second-order channel attention is proposed to rescale feature statistics and enhance discriminative learning power. To build a deeper network, a non-locally enhanced residual group (NLRG) is utilized that shares contextual information. Additionally, skip connections are utilized to transfer learned representations across numerous layers within the neural network. These connections enable the network to effectively utilize the bypassed skip connections, facilitating the flow of information from low-resolution input images to the subsequent layers and simplifying the training process.

3.2. Proposed solution

The proposed SAN network can be divided into four parts. 1) shallow feature extraction, 2) non-locally enhanced residual group (NLRG) based deep feature extraction, 3) up-scale module, and 4) reconstruction part. Let, I_{LR} and I_{SR} are

the input and output of the network. At first, the network extracts shallow feature F_0 from the LR input

$$F_0 = H_{SF}(I_{LR}) \quad (7)$$

where $H_{SF}(\cdot)$ corresponds to the convolution operation. Then the shallow features are fed into NLRG-based deep feature extraction and the deep features F_{DF} are found.

$$F_{DF} = H_{NLRG}(F_0) \quad (8)$$

Here, H_{NLRG} enlarges feature learning by sequentially passing the input through several non-local modules. The output of NLRG is then upscaled via upscale module F_{\uparrow} .

$$F_{\uparrow} = H_{\uparrow}(F_{DF}) \quad (9)$$

In the reconstruction part, the upscaled features are converted to the super-resolution output via one convolution block H_R .

$$I_{SR} = H_R(F_{\uparrow}) = H_{SAN}(I_{LR}) \quad (10)$$

Here, H_{SAN} is the overall SAN network.

Optimization is one important aspect of neural networks and this work proposes a n_1 -norm loss function between the output super-resolution image and the ground truth for this purpose.

$$L_{\Theta} = \frac{1}{N} \sum_{i=1}^N \lambda \|I_{HR}^i - H_{SAN}(I_{LR}^i)\|_1 \quad (11)$$

where, $\theta \in \Theta$ are the network parameters and I_{HR}^i is the high resolution ground truth image.

3.3. Claims-Evidence

Claim 1 According to the author, the proposed second-order channel attention (SOCA) adaptively learns features using second-order feature statistics instead of first-order ones which in turn puts emphasis on informative feature and improve discriminative learning ability rather than the other methods that incorporate first-order statistics.

Evidence 1 Figure 4 demonstrates that incorporating second-order statistics ($R_e = 32.16$ dB) indeed enhances image quality and captures more intricate details compared to first-order ($R_d = 32.12$ dB).

Table 1. Effects of different modules. We report the best PSNR (dB) values on Set5 (4×) in 5.6×10^5 iterations.

	R_a	R_b	R_c	R_d	R_e	R_f	R_g	R_h	R_i
RL-NL(before SSRG)	✓					✓	✓	✓	✓
RL-NL(after SSRG)		✓				✓	✓	✓	✓
share-source skip connection (SSC)			✓			✓	✓	✓	✓
First-order channel attention (FOCA)				✓				✓	
Second-order channel attention (SOCA)					✓				✓
	32.00	32.04	32.06	32.07	32.12	32.16	32.08	32.10	32.14
									32.20

Figure 4. Table 1 is taken from page 7, section 4.2 of the original paper (Dai et al., 2019) as a screenshot.

220
 221
 222
 223
 224
 225
 226
 227
 228
 229
 230
 231
 232
 233
 234
 235
 236
 237
 238
 239
 240
 241
 242
 243
 244
 245
 246
 247
 248
 249
 250
 251
 252
 253
 254
 255
 256
 257
 258
 259
 260
 261
 262
 263
 264
 265
 266
 267
 268
 269
 270
 271
 272
 273
 274
 275
 276
 277
 278
 279
 280
 281
 282
 283
 284
 285
 286
 287
 288
 289
 290
 291
 292
 293
 294
 295
 296
 297
 298
 299
 300
 301
 302
 303
 304
 305
 306
 307
 308
 309
 310
 311
 312
 313
 314
 315
 316
 317
 318
 319
 320
 321
 322
 323
 324
 325
 326
 327
 328
 329
 330
 331
 332
 333
 334
 335
 336
 337
 338
 339
 340
 341
 342
 343
 344
 345
 346
 347
 348
 349
 350
 351
 352
 353
 354
 355
 356
 357
 358
 359
 360
 361
 362
 363
 364
 365
 366
 367
 368
 369
 370
 371
 372
 373
 374
 375
 376
 377
 378
 379
 380
 381
 382
 383
 384
 385
 386
 387
 388
 389
 390
 391
 392
 393
 394
 395
 396
 397
 398
 399
 400
 401
 402
 403
 404
 405
 406
 407
 408
 409
 410
 411
 412
 413
 414
 415
 416
 417
 418
 419
 420
 421
 422
 423
 424
 425
 426
 427
 428
 429
 430
 431
 432
 433
 434
 435
 436
 437
 438
 439
 440
 441
 442
 443
 444
 445
 446
 447
 448
 449
 450
 451
 452
 453
 454
 455
 456
 457
 458
 459
 460
 461
 462
 463
 464
 465
 466
 467
 468
 469
 470
 471
 472
 473
 474
 475
 476
 477
 478
 479
 480
 481
 482
 483
 484
 485
 486
 487
 488
 489
 490
 491
 492
 493
 494
 495
 496
 497
 498
 499
 500
 501
 502
 503
 504
 505
 506
 507
 508
 509
 510
 511
 512
 513
 514
 515
 516
 517
 518
 519
 520
 521
 522
 523
 524
 525
 526
 527
 528
 529
 530
 531
 532
 533
 534
 535
 536
 537
 538
 539
 540
 541
 542
 543
 544
 545
 546
 547
 548
 549
 550
 551
 552
 553
 554
 555
 556
 557
 558
 559
 560
 561
 562
 563
 564
 565
 566
 567
 568
 569
 570
 571
 572
 573
 574
 575
 576
 577
 578
 579
 580
 581
 582
 583
 584
 585
 586
 587
 588
 589
 590
 591
 592
 593
 594
 595
 596
 597
 598
 599
 600
 601
 602
 603
 604
 605
 606
 607
 608
 609
 610
 611
 612
 613
 614
 615
 616
 617
 618
 619
 620
 621
 622
 623
 624
 625
 626
 627
 628
 629
 630
 631
 632
 633
 634
 635
 636
 637
 638
 639
 640
 641
 642
 643
 644
 645
 646
 647
 648
 649
 650
 651
 652
 653
 654
 655
 656
 657
 658
 659
 660
 661
 662
 663
 664
 665
 666
 667
 668
 669
 670
 671
 672
 673
 674
 675
 676
 677
 678
 679
 680
 681
 682
 683
 684
 685
 686
 687
 688
 689
 690
 691
 692
 693
 694
 695
 696
 697
 698
 699
 700
 701
 702
 703
 704
 705
 706
 707
 708
 709
 710
 711
 712
 713
 714
 715
 716
 717
 718
 719
 720
 721
 722
 723
 724
 725
 726
 727
 728
 729
 730
 731
 732
 733
 734
 735
 736
 737
 738
 739
 740
 741
 742
 743
 744
 745
 746
 747
 748
 749
 750
 751
 752
 753
 754
 755
 756
 757
 758
 759
 760
 761
 762
 763
 764
 765
 766
 767
 768
 769
 770
 771
 772
 773
 774
 775
 776
 777
 778
 779
 780
 781
 782
 783
 784
 785
 786
 787
 788
 789
 790
 791
 792
 793
 794
 795
 796
 797
 798
 799
 800
 801
 802
 803
 804
 805
 806
 807
 808
 809
 810
 811
 812
 813
 814
 815
 816
 817
 818
 819
 820
 821
 822
 823
 824
 825
 826
 827
 828
 829
 830
 831
 832
 833
 834
 835
 836
 837
 838
 839
 840
 841
 842
 843
 844
 845
 846
 847
 848
 849
 850
 851
 852
 853
 854
 855
 856
 857
 858
 859
 860
 861
 862
 863
 864
 865
 866
 867
 868
 869
 870
 871
 872
 873
 874
 875
 876
 877
 878
 879
 880
 881
 882
 883
 884
 885
 886
 887
 888
 889
 890
 891
 892
 893
 894
 895
 896
 897
 898
 899
 900
 901
 902
 903
 904
 905
 906
 907
 908
 909
 910
 911
 912
 913
 914
 915
 916
 917
 918
 919
 920
 921
 922
 923
 924
 925
 926
 927
 928
 929
 930
 931
 932
 933
 934
 935
 936
 937
 938
 939
 940
 941
 942
 943
 944
 945
 946
 947
 948
 949
 950
 951
 952
 953
 954
 955
 956
 957
 958
 959
 960
 961
 962
 963
 964
 965
 966
 967
 968
 969
 970
 971
 972
 973
 974
 975
 976
 977
 978
 979
 980
 981
 982
 983
 984
 985
 986
 987
 988
 989
 990
 991
 992
 993
 994
 995
 996
 997
 998
 999
 1000
 1001
 1002
 1003
 1004
 1005
 1006
 1007
 1008
 1009
 1010
 1011
 1012
 1013
 1014
 1015
 1016
 1017
 1018
 1019
 1020
 1021
 1022
 1023
 1024
 1025
 1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079
 1080
 1081
 1082
 1083
 1084
 1085
 1086
 1087
 1088
 1089
 1090
 1091
 1092
 1093
 1094
 1095
 1096
 1097
 1098
 1099
 1100
 1101
 1102
 1103
 1104
 1105
 1106
 1107
 1108
 1109
 1110
 1111
 1112
 1113
 1114
 1115
 1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133
 1134
 1135
 1136
 1137
 1138
 1139
 1140
 1141
 1142
 1143
 1144
 1145
 1146
 1147
 1148
 1149
 1150
 1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172
 1173
 1174
 1175
 1176
 1177
 1178
 1179
 1180
 1181
 1182
 1183
 1184
 1185
 1186
 1187
 1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241
 1242
 1243
 1244
 1245
 1246
 1247
 1248
 1249
 1250
 1251
 1252
 1253
 1254
 1255
 1256
 1257
 1258
 1259
 1260
 1261
 1262
 1263
 1264
 1265
 1266
 1267
 1268
 1269
 1270
 1271
 1272
 1273
 1274
 1275
 1276
 1277
 1278
 1279
 1280
 1281
 1282
 1283
 1284
 1285
 1286
 1287
 1288
 1289
 1290
 1291
 1292
 1293
 1294
 1295
 1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319
 1320
 1321
 1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349
 1350
 1351
 1352
 1353
 1354
 1355
 1356
 1357
 1358
 1359
 1360
 1361
 1362
 1363
 1364
 1365
 1366
 1367
 1368
 1369
 1370
 1371
 1372
 1373
 1374
 1375
 1376
 1377
 1378
 1379
 1380
 1381
 1382
 1383
 1384
 1385
 1386
 1387
 1388
 1389
 1390
 1391
 1392
 1393
 1394
 1395
 1396
 1397
 1398
 1399
 1400
 1401
 1402
 1403
 1404
 1405
 1406
 1407
 1408
 1409
 1410
 1411
 1412
 1413
 1414
 1415
 1416
 1417
 1418
 1419
 1420
 1421
 1422
 1423
 1424
 1425
 1426
 1427
 1428
 1429
 1430
 1431
 1432
 1433
 1434
 1435
 1436
 1437
 1438
 1439
 1440
 1441
 1442
 1443
 1444
 1445
 1446
 1447
 1448
 1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457
 1458
 1459
 1460
 1461
 1462
 1463
 1464
 1465
 1466
 1467
 1468
 1469
 1470
 1471
 1472
 1473
 1474
 1475
 1476
 1477
 1478
 1479
 1480
 1481
 1482
 1483
 1484
 1485
 1486
 1487
 1488
 1489
 1490
 1491
 1492
 1493
 1494
 1495
 1496
 1497
 1498
 1499
 1500
 1501
 1502
 1503
 1504
 1505
 1506
 1507
 1508
 1509
 1510
 1511
 1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565
 1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1600
 1601
 1602
 1603
 1604
 1605
 1606
 1607
 1608
 1609
 1610
 1611
 1612
 1613
 1614
 1615
 1616
 1617
 1618
 1619
 1620
 1621
 1622
 1623
 1624
 1625
 1626
 1627
 1628
 1629
 1630
 1631
 1632
 1633
 1634
 1635
 1636
 1637
 1638
 1639
 1640
 1641
 1642
 1643
 1644
 1645
 1646
 1647
 1648
 1649
 1650
 1651
 1652
 1653
 1654
 1655
 1656
 1657
 1658
 1659
 1660
 1661
 1662
 1663
 1664
 1665
 1666
 1667
 1668
 1669
 1670
 1671
 1672
 1673
 1674
 1675
 1676
 1677
 1678
 1679
 1680
 1681
 1682
 1683
 1684
 1685
 1686
 1687
 1688
 1689
 1690
 1691
 1692
 1693
 1694
 1695
 1696
 1697
 1698
 1699
 1700
 1701
 1702
 1703
 1704
 1705
 1706
 1707
 1708
 1709
 17010
 17011
 17012
 17013
 17014
 17015
 17016
 17017
 17018
 17019
 17020
 17021
 17022
 17023
 17024
 17025
 17026
 17027
 17028
 17029
 1

275 term dependencies within the images. Despite the superior
 276 performance of vision transformers in single-image super-
 277 resolution compared to methods based on CNNs, their high
 278 memory consumption has limited their popularity.
 279

280 **Prior work on this problem** Although deep learning net-
 281 works have been used in single-image super-resolution, the
 282 concept of vision transformers is comparatively new. Some
 283 Recurrent Neural Network (RNN) based architectures were
 284 also deployed for this task but the reduced network pa-
 285 rameters led them to perform poorly such as DRCN (Kim
 286 et al., 2016), SRRFN (Li et al., 2019), and IMDN (Hui
 287 et al., 2019). The transformer network came from the
 288 NLP literature, especially the element “self-attention” in
 289 the transformer which captures the long-term dependen-
 290 cies in sequence. However, in vision, it is represented as
 291 long-term dependencies among image patches (Dosovit-
 292 skiy et al., 2020). A significant advancement occurred with
 293 SwinIR (Liu et al., 2021), employing Swin Transformers
 294 for single-image super-resolution with remarkable pre-
 295 cision, surpassing the performance of state-of-the-art CNN
 296 techniques.

297 **Research gap** CNN-based methods incur substantial GPU
 298 costs, prompting researchers to explore RNN-based alterna-
 299 tives. Additionally, CNNs can leverage local features, but
 300 they do not comprehensively analyze global features. How-
 301 ever, RNN methods perform notably worse than CNN-based
 302 approaches due to their shallower network depth. Although
 303 Swin Transformer significantly improves reconstruction per-
 304 formance it suffers from high memory consumption.
 305

306 **Contributions** In this paper (Lu et al., 2022), the author
 307 proposed an efficient transformer network for image super-
 308 resolution known as ESRT (Efficient Super-Resolution
 309 Transformer). Specifically, the author proposed a hybrid
 310 model that consists of a lightweight CNN backbone (LCB)
 311 and a lightweight transformer backbone (LTB). LCB utilizes
 312 high-preserving blocks (HPBs) to adjust the feature maps
 313 with low-cost computation and extract texture details before
 314 feeding into the LTB. The LTB utilizes efficient multihead
 315 attention (EMHA) mechanisms to efficiently capture the
 316 dependencies among image patches with reduced cost.
 317
 318

4.2. Proposed Solution

319 Efficient Super-resolution transformer consists of four parts.
 320 They are shallow feature extraction, Lightweight CNN back-
 321 bone (LCB), Lightweight transformer backbone (LTB), and
 322 reconstruction block. The image input is I_{LR} and the output
 323 I_{SR} , the shallow feature extractor is defined as a convolution
 324 operation (f_s).
 325

$$F_0 = f_s(I_{LR}) \quad (12)$$

The extracted shallow feature, F_0 is passed into the LCB which consists of several High Preserving Blocks (HPBs) denoted as ζ .

$$F_n = \zeta^n(\zeta^{n-1}(\dots(\zeta^2(\zeta^1(F_0))))) \quad (13)$$

Here, ζ^n represents the n^{th} HPB block and F_n denotes the output of n^{th} HPB. The outputs are combined and sent to several efficient transformers (denoted as ϕ) to merge intermediate features.

$$F_d = \phi^n(\phi^{n-1}(\dots(\phi^2(\phi^1([F_1, F_2, \dots, F_n]))))) \quad (14)$$

Here, F_d denotes the output of LTB. Finally, in the recon-
 struction module, the shallow feature F_0 and the output
 of the LTB, F_d are used jointly to reconstruct the high-
 resolution image.

$$I_{SR} = f(f_p(f(F_d))) + f(f_p(F_0)) \quad (15)$$

Here, f denotes the convolution layer, and f_p denotes the
 pixel shuffle layer.

4.3. Claims-Evidence

Claim 1 There are similar patches in a single image with
 the same type of pixels which means we can count them as
 ground truth. Therefore, for each patch, we can have mul-
 tiple ground-truth images. The author claimed that similar
 patches can be used for ground truth. Therefore, certain tex-
 ture details can be resolved using multiple reference patches.
 This operates in a manner akin to the way transformers
 function.

Evidence 1 Figure 8 demonstrates that without the trans-
 former the PSNR drops from 32.18 dB to 31.96 dB. There-
 fore, the transformers are able to exploit the patches in the
 image and improve the reconstruction.

Case	PSNR(dB)	Param.(K)	GPU Memory
w/o TR	31.96	554	1931M
Original TR [38]	32.14	971	16057M
1 ET	32.18	751	4191M
2 ET	32.25	949	6499M

Table 4. Study of Efficient Transformer (ET) on Set5 ($\times 4$). The GPU memory here refers to the cost of the model during training, which patch_size = 48*48 and batch_size=16.

Figure 8. Table 4 is taken from page 8, section 4.5.2 of the original paper (Lu et al., 2022) as a screenshot.

Claim 2 Because the Fourier Transform requires a sig-
 nificant amount of memory, the authors suggested using a
 lightweight, differentiable high-frequency filtering module
 (HFM) to extract high-frequency information from low-
 resolution input in a more efficient manner.

330
 331 **Evidence 2** In Figure 9, we can see that for cases 1, 2,
 332 and 3 the HFM module improves the performance with less
 333 amount of parameters.
 334

Case Index	1	2	3	4
HFM		✓	✓	✓
CA	✓	✓		✓
ARFB	✓	✓	✓	
RB				✓
Parameters	658K	751K	724K	972K
PSNR	32.02dB	32.19dB	32.08dB	32.20dB

Table 3. Study of each component in HPB on Set5 ($\times 4$).

341
 342 *Figure 9.* Table 3 is taken from page 7, section 4.5.1 of the original
 343 paper (Lu et al., 2022) as a screenshot.
 344
 345

346
 347 **Claim 3** While Swin Transformer (SwinIR) is effective
 348 for image super-resolution, the author asserted that it con-
 349 sumes more GPU memory. The proposed approach, how-
 350 ever, presents a superior solution by enhancing image qual-
 351 ity with fewer parameters, addressing the issue of high GPU
 352 memory consumption associated with SwinIR.
 353
 354

355 **Evidence 3** As we can see from Figure 10, the number of
 356 parameters for SwinIR is higher than the proposed ESRT
 357 but the difference between PSNR is 0.01 dB. Therefore,
 358 the proposed ESRT can fully exploit the advantages of a
 359 transformer with fewer parameters.
 360
 361

Method	Parame.	GPU Memory	BSD100	Manga109
SwinIR	886K	6966M	29.20/0.8082	33.98/0.9478
ESRT	770K	4191M	29.15/0.8063	33.95/0.9455

Table 7. A detailed comparison of SwinIR and ESRT ($\times 4$).

362 *Figure 10.* Table 7 is taken from page 8, section 4.7 of the original
 363 paper (Lu et al., 2022) as a screenshot.
 364
 365

4.4. Critique and Discussion

366 The writing in this paper is undeniably exceptional, partic-
 367 ularly in the introduction section. What fascinated me the
 368 most was the explanation of utilizing transformers, demon-
 369 strated in a figure that illustrates how multiple patches cor-
 370 respond to the same image. Therefore, the rationale under-
 371 lying the transformer was evident and comprehensible.
 372
 373

374 In the related works section, the authors provided a brief
 375 overview of CNN-based research. It is recommended that
 376 the authors expand their coverage of prior studies to enhance
 377 the reader’s comprehension of the field’s historical context
 378 and development. Nevertheless, the sections on attention
 379 and transformers were presented in a general manner, al-
 380
 381

382 lowing the reader to discern the underlying trend and the
 383 author’s intended message.
 384

385 While certain implementation specifics were missing, the
 386 paper provided a thorough examination of each network
 387 block utilized, rendering it comprehensive yet accessible to
 388 a broader audience.
 389

5. Implementation

5.1. Implementation motivation

390 The authors of the paper did not share any code, prompting
 391 me to undertake the project and reimplement their work
 392 using PyTorch. The paper lacks crucial details, such as the
 393 difference image between input and output, and the rationale
 394 behind the use of l_1 loss instead of the more commonly used
 395 l_2 loss. Additionally, there is a lack of explanation regarding
 396 the utilization of RCAG. It remains unclear how attention
 397 blocks collectively work to preserve high-frequency infor-
 398 mation. The paper vaguely states the use of 15 Residual
 399 Channel Attention Blocks (RCABs) in each RCAG. The se-
 400 lection of this specific number of RCABs should be substi-
 401 tuted with several experiments to validate its effectiveness.
 402 Through the reimplementation, I will seek to answer these
 403 questions.
 404

405 Additionally, I have one particular idea to implement inside
 406 the network - a sharpness enhancer feature block. This
 407 block would utilize phase images (which inherently contain
 408 edges) in the reconstruction process, employing RCAGs to
 409 incorporate them by addition or multiplication before the
 410 final concatenation. Given that many networks face issues
 411 with blurriness, I believe incorporating a sharpness enhancer
 412 could mitigate blurriness and enhance sharp edges.
 413

414 If it does not work, I will try to enhance sharpness by addi-
 415 tional constraints within the model loss as shown below:
 416

$$L_{\Theta'} = L_{\Theta} + \|edge(I_{HR}) - edge(I_{SR})\|_1 \quad (16)$$

417 The edge can be extracted whether by reconstructing with
 418 phase only or through traditional algorithms e.g., Canny
 419 edge detection. If the aforementioned methods prove in-
 420 effective, I will attempt to integrate the ‘Sobel’ kernel for
 421 edge extraction, which is differentiable. The image ¹ in Fig-
 422 ure 11 somewhat explains the intuition behind using phase
 423 reconstruction as a sharpness enhancer.
 424

425 Among phase and magnitude, phase determines the maxi-
 426 mum information of an image, including the edges. There-
 427 fore, I am expecting to see improvements in sharpness after
 428 incorporating the sharpness feature enhancer.
 429

430 ¹<https://math.stackexchange.com/questions/849382/image-reconstructionphase-vs-magnitude>
 431

385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439



Figure 11. Both the image and the caption are taken from the website mentioned in the footnote.

5.2. Implementation plan and setup

In the process of reimplementing, I plan to utilize the DIV2K dataset (Agustsson & Timofte, 2017; Timofte et al., 2017) and PyTorch framework for developing the network. First, the DIV2K training data set (consisting of 800 images) and the test set (consisting of 100 images) are used to create the baseline of the project. The setup is implemented on an 8 GB RTX 4060 GPU. In the initial phase of my implementation, I will preprocess the dataset and execute a basic super-resolution using a simple CNN model along with MSE loss between the reconstructed image and the ground truth image. This step will help me to familiarize myself with the dataset and the super-resolution process. After that, I will try to implement a super-resolution network only with the Residual Channel Attention Block (RCAB) to fully understand the necessity of residual channel attention blocks for super-resolution. Afterward, I will undertake the task of implementing the proposed TSN, omitting the second stage, which is responsible for transforming the high-resolution image into a low-resolution version with high frequency (the second component of the loss function). The authors of the paper also mentioned it in their “Ablation Study” and will check whether my results are consistent with the paper. From that, I will try to implement the original TSN network with consistent parameters. However, parameter study will not be a priority for this project rather I will focus on implementing the network architecture. The evaluation metric considered for this project is PSNR (Peak Signal-to-Noise

Ratio) (Horé & Ziou, 2010).

To search for concerns raised in Section 5.1, I will include a difference image between I_{SR} and I_{HR} to fully visualize the missing pixels in the reconstruction process. Following that, I will conduct multiple experiments, varying the number of RCABs (1, 15, 20) to comprehensively assess the impact of attention blocks in the context of super-resolution. Additionally, I will incorporate the l_2 loss instead of the l_1 loss to provide justification for using the first one.

To validate my hypothesis of the edge enhancer block, I will attempt to incorporate a sharpness enhancer block and assess whether it enhances overall performance.

5.3. Implementation details

The original network’s code is not accessible to the public, necessitating the creation of the code from scratch. However, I did take motivation from some of the previous works. The channel attention block is first introduced by Zhang et al., (Zhang et al., 2018b). In their work, they stacked residual channel attention blocks to increase the super-resolution performance. I took the channel attention layer code directly from them.

The data loader was written by me but took some inspiration from the book, “Deep Learning with Pytorch” (Stevens et al., 2020). I initially trained the network in a dual-loop setting like the other papers did. But later I found out that the whole dataset is small enough to fit within the GPU along with the model. So, I came up with a different approach. At first, I loaded the GPU with all the possible batches of data at once. During training, a random number is called in each epoch that takes one random batch of images to calculate the gradient and backpropagate. It saved more time than expected. Initially, in a dual-loop setting, it took 15 hours to train for 500 epochs. Later, in a single-loop setting, I was able to train the model for 4000 epochs and it took 2 hours only.

At first, I attempted to utilize the Mean Squared Error (MSE) loss, but subsequently transitioned to the l_1 loss, which aligns more coherently with both theoretical principles and experimental observations. In a super-resolution task, the goal is to match as much as pixels possible. The l_1 loss tends to sparsify the difference, implying that the gap between the low-resolution image and the predicted high-resolution image should approach zero. That’s why it is better to use l_1 loss instead of l_2 in super-resolution tasks. The PSNR was calculated in a different way. Traditionally, the PSNR is expected to be based on the calculation of the MSE error from the RGB images. However, in Python, the values are significantly lower than Matlab values. Nonetheless, it is a common practice to convert the RGB images to YCbCr

² and take the Y plane only to measure the PSNR. In my implementation, I have also followed the same convention to match the PSNR values with the paper.

5.4. Results and interpretation

Initial experiments focused more on setting up a fundamental Convolutional Neural Network (CNN) architecture for super-resolution which includes the construction of the workflow and fitting the data into the network. The first experiment is summarized below:

- a simple super-resolution network was built with convolutional neural networks only.
 - For this experiment, a simple MSE loss was used. For further experiments, MSE loss was used unless specified.
 - The parameters (batch size, epochs) do not exactly follow the original paper, as the parameter study is not the priority.

A basic CNN should theoretically be capable of reconstructing high-resolution images due to its inherent averaging property in convolution, where neighboring pixels are summed together. However, there should be a significant blurring. We can observe this blurring in Figure 12.

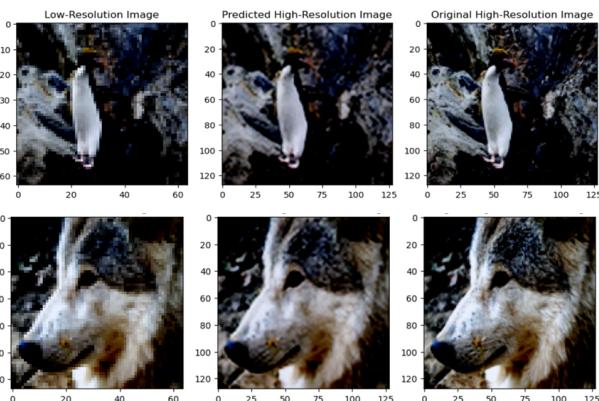


Figure 12. The input image is scaled by 0.5. All the images are normalized between values 0 to 1. This simple CNN model was trained for 100 epochs.

To fully understand the necessity of residual channel attention blocks a network is trained only using Residual Channel Attention Blocks with some basic convolutions. In this case, the MSE loss is used and the training persisted for 100 epochs. The predicted images are shown in Figure 13.

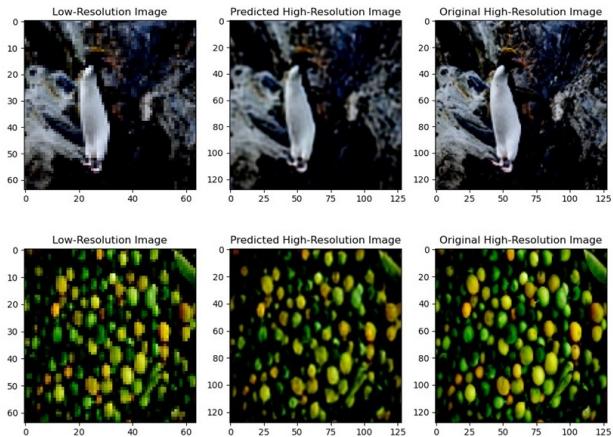


Figure 13. The input image is scaled by 0.5. All the images are normalized between values 0 to 1. In this model, a single RCAB block is utilized along with some convolution operations before and after the RCAB block.

As we can see, the images are much sharper compared to using CNNs alone. It is evident that the residual channel attention block plays a crucial role in enhancing image sharpness and removing blurriness.

The paper (Han et al., 2021) mentioned that they used 15 Residual Channel Attention Blocks (RCABs) per Residual Channel Attention Group (RCAG). We conducted an experiment where one RCAG (consisting of 15 RCABs) was utilized. The predicted images are shown in Figure 14.

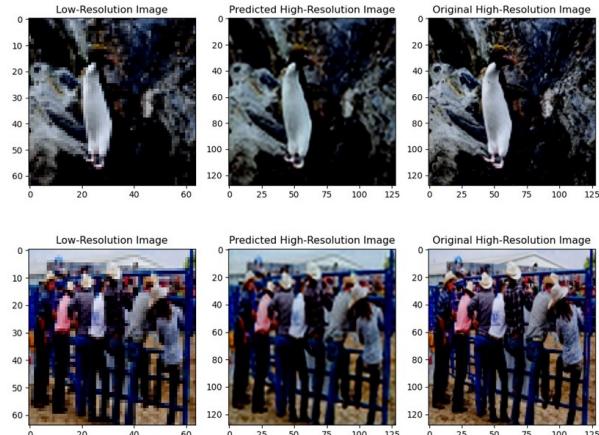


Figure 14. The input image is scaled by 0.5. All the images are normalized between values 0 to 1. In this model, a single RCAG block is utilized along with some convolution operations before and after the RCAG block.

Finally, I implemented the original network without the second stage (the second term in the loss function in Equation (6)). In this experiment, 8 RCAG blocks were used where each RCAG consists of 15 RCAB blocks. The train-

²<https://github.com/krasserm/super-resolution/issues/19>

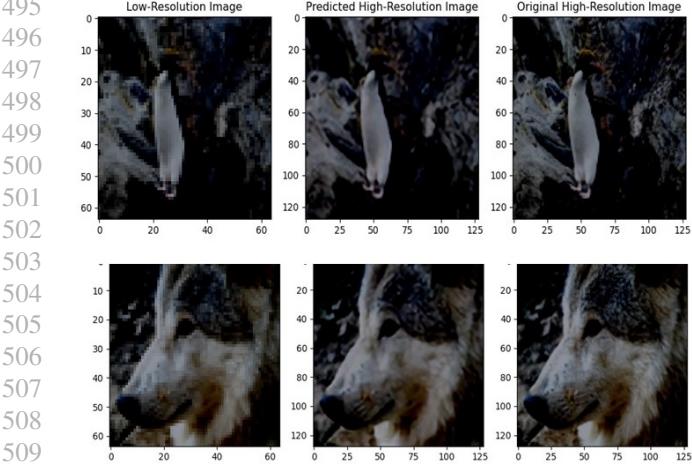


Figure 15. The input image is scaled by 0.5. All the images are normalized between values 0 to 1. In this model, eight consecutive RCAGs were used to reconstruct the image.

ing was stopped at 135 epochs because it required significantly more time compared to previous experiments. This time I used the l_1 loss, keeping consistency with the original paper. The output images as well as the difference images are shown in Figure 15, and Figure 16.

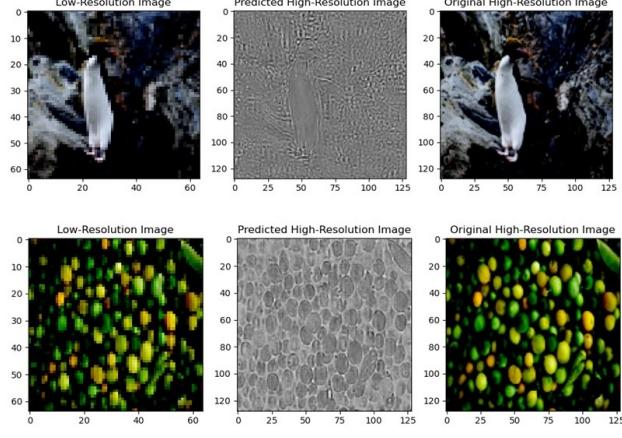


Figure 16. To generate the difference image, both the ground truth and the predicted image were transformed into grayscale. In this process, subtraction was applied between the two images to calculate the difference.

A difference image can visualize the performance of good reconstruction. Therefore, from Figure 16, it is evident that the reconstruction preserved most of the information, although it did miss certain pixel values.

It is expected that incorporating more RCAB blocks can lead to good reconstruction performance. From Table 1, it is evident that using 15 RCAB instead of one improves the performance. However, as we can see after implementing

Table 1. PSNR values on the test set for different methods.

DATA SET	MODEL	AVERAGE PSNR
DIV2K	1 RCAB	26.74 dB
DIV2K	1 RCAG (EACH INCLUDES 15 RCAB)	27.23 dB
DIV2K	TSN (8 RCAG)	26.07 dB

the TSN network the performance drops. This is due to the fact that the training was stopped at epoch 135. Moreover, the PSNR values slightly differ from the values mentioned in the paper. This is a common issue faced by the Python community. I used Python to calculate PSNR. Python "imresize" function is different from the Matlab "imresize" and this leads to some errors while calculating the PSNR. Therefore, I utilized the YCbCr planes to calculate PSNR (NOD, 2007) in the next experiments as mentioned in the "Implementation details" section.

In the first checkpoint, the first stage of the two-stage network was implemented. In particular, I implemented a comprehensive deep network, but the second-stage network, which served as an additional loss component, was not implemented. In the second checkpoint, I implemented the second network, through which the high-resolution input goes through multiple residual channel attention groups and is then compared with a corresponding layer from the initial network. The resulting loss is combined with the original reconstruction loss. If I recall from the previous checkpoint, I implemented the second portion of the loss shown in Equation 17.

$$L_{\Theta} = \frac{1}{N} \sum_{n=1}^N \lambda \|I_{HR} - I_{SR}\|_1 + \|I_{HF} - I_{LR'}\|_1 \quad (17)$$

Furthermore, I initially intended to integrate a phase loss to emphasize edge enhancement within the network. Unfortunately, this plan had to be abandoned because the Fourier Transform method required significantly more computa-

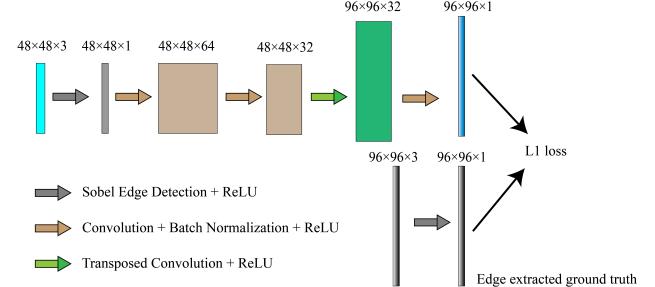


Figure 17. In the third stage network, a CNN initialized with a "Sobel" kernel is integrated to compute the edge loss.

Table 2. PSNR values on the test set for different methods.

DATA SET	MODEL	AVERAGE PSNR
DIV2K	1 RCAB	26.74 dB
DIV2K	1 RCAG (EACH INCLUDES 15 RCAB)	27.23 dB
DIV2K	TSN (8 RCAG)	26.07 dB
DIV2K	TSN	30.39 dB
DIV2K	TSN + SOBEL LOSS	30.87 dB

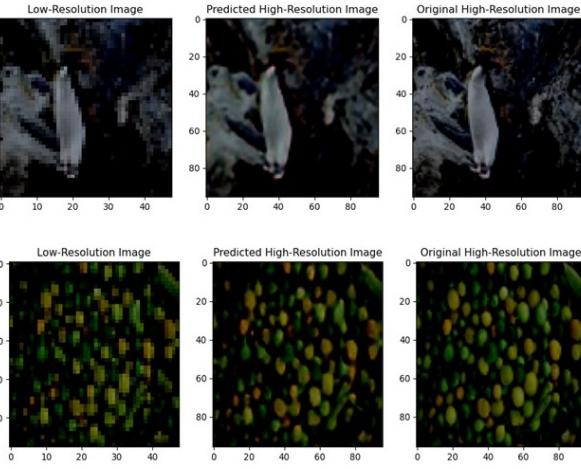


Figure 18. Predicted high-resolution images for the two-stage network (TSN).

tional power than my current setup could handle. Consequently, I opted for an alternative approach. I incorporated a shallow convolution neural network using a “Sobel” kernel at the network’s input and included it as another loss term. The results of this approach, as well as the previous one (two-stage network), are shown in Table 2. The third stage that I added is shown in Figure 17. For the TSN method, the predicted high resolution images as well as the difference images are shown in Figures 18, and 19.

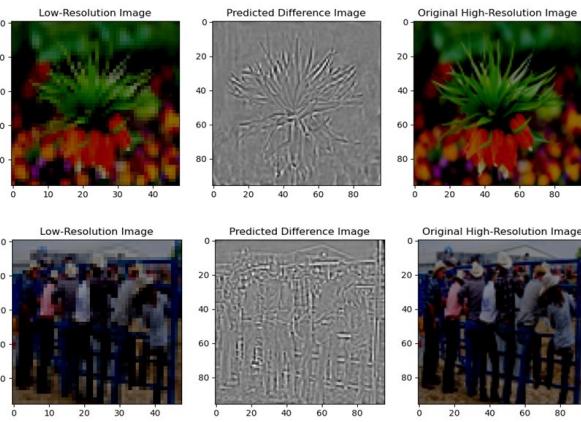


Figure 19. Difference images for the two-stage network (TSN).

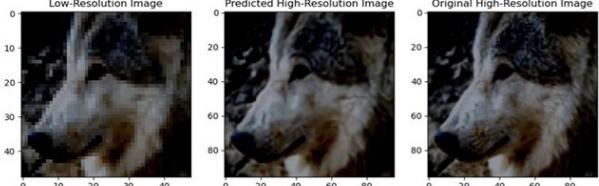


Figure 20. Predicted high-resolution images for the network where Sobel loss is incorporated as an additional loss with the TSN network.

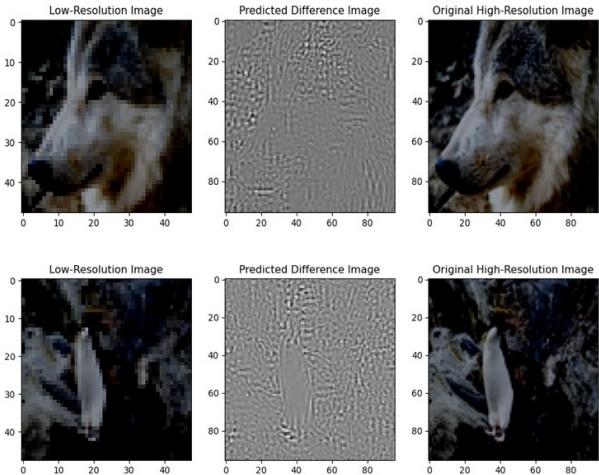


Figure 21. Difference images for the network where Sobel loss is incorporated as an additional loss with the TSN network.

As we can see from the difference images in Figure 19, the edges are not fully resolved. However, after incorporating Sobel loss this edge differences are slightly improved as shown in Figures 20, and 21. Nevertheless, the discrepancy is not substantial, but I attribute this largely to employing a relatively shallow network for enhancing edges. For future studies, I will try to pass the edge extracted image (by “Sobel” kernel) through residual channel attention blocks so that the network can pull more low-level details from the edges.

6. Conclusion and Discussion

The proposed approach integrates residual channel attention blocks to preserve high-frequency information, consequently enhancing the performance of the super-resolution

network. The dual-stage loss functions facilitated the collaborative optimization of the network, guiding it to preserve high-frequency details. While implementing the network, I discovered various details that were not explicitly outlined in the original paper.

Initially, the original paper does not clearly identify the specific block responsible for preserving high-frequency information. However, based on my implementation, it becomes apparent that the residual channel attention block is the layer responsible for retaining high-frequency information. In each channel attention block, the input is passed through several layers and the final output is multiplied with the input. In my understanding, the elementwise multiplication boosts the pixels, thereby preserving high-frequency information. Therefore, in theory, increasing the attention blocks should lead to improved performance. From the data presented in Table 2, it is observed that integrating several RCABs within a single RCAG resulted in an approximately 1 dB increase in PSNR.

Secondly, the paper lacks clarity on the reasoning behind the number of RCAGs. Therefore, I experimented with 4, 8, and 15 RCABs within a single RCAG. It is confirmed that augmenting the number of RCABs marginally boosts the PSNR. Unfortunately, I couldn't explore further due to hardware limitations.

Thirdly, the idea behind incorporating an edge-enhanced layer was to maintain the edges of the low-resolution image. As evident in the difference image shown in Figure 21, the edges are preserved more effectively compared to previous methods. However, it is noteworthy that the edge-enhance layer was relatively shallow. While deepening the layer might potentially enhance performance, I did not explore this aspect in the current project, reserving it for future investigation.

During the implementation of the paper, I discovered a method that enabled me to train my network for 1000 epochs in just 20 minutes. I loaded all my data onto the GPU in a single loop. Subsequently, for each epoch, I randomly accessed the stored data, computed gradients, and updated them. This approach significantly saved time, considering that in the traditional double-loop format, data loading consumes a considerable amount of time.

I am presently engaged in an MRI project, wherein the task involves reconstructing a high-resolution image from an aliased image. In this context, I can seamlessly integrate the method discussed. The input needs to be adjusted within the MRI reconstruction pipeline for effective implementation.

References

- Colorization in ycbr color space and its application to jpeg images. *Pattern Recognition*, 40(12):3714–3720, 2007. ISSN 0031-3203. doi: <https://doi.org/10.1016/j.patcog.2007.04.005>. URL <https://www.sciencedirect.com/science/article/pii/S0031320307001793>.
- Agustsson, E. and Timofte, R. Ntire 2017 challenge on single image super-resolution: Dataset and study. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, July 2017.
- Bevilacqua, M., Roumy, A., Guillemot, C., and Alberi Morel, M.-L. Neighbor embedding based single-image super-resolution using semi-nonnegative matrix factorization. In *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1289–1292, 2012. doi: 10.1109/ICASSP.2012.6288125.
- Chakrabarti, A., Rajagopalan, A. N., and Chellappa, R. Super-resolution of face images using kernel pca-based prior. *IEEE Transactions on Multimedia*, 9(4):888–892, 2007. doi: 10.1109/TMM.2007.893346.
- Chang, H., Yeung, D.-Y., and Xiong, Y. Super-resolution through neighbor embedding. In *Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2004. CVPR 2004.*, volume 1, pp. I–I, 2004. doi: 10.1109/CVPR.2004.1315043.
- Chen, C., Xiong, Z., Tian, X., Zha, Z.-J., and Wu, F. Camera lens super-resolution. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, June 2019. doi: 10.1109/cvpr.2019.00175. URL <https://doi.org/10.1109/cvpr.2019.00175>.
- Dai, T., Cai, J., Zhang, Y., Xia, S.-T., and Zhang, L. Second-order attention network for single image super-resolution. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 11057–11066, 2019. doi: 10.1109/CVPR.2019.01132.
- Dong, C., Loy, C. C., He, K., and Tang, X. Image super-resolution using deep convolutional networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(2):295–307, 2016. doi: 10.1109/TPAMI.2015.2439281.
- Dong, W., Zhang, L., Shi, G., and Wu, X. Image deblurring and super-resolution by adaptive sparse domain selection and adaptive regularization. *IEEE Transactions on Image Processing*, 20(7):1838–1857, 2011. doi: 10.1109/TIP.2011.2108306.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer,

- 660 M., Heigold, G., Gelly, S., Uszkoreit, J., and Houlsby,
 661 N. An image is worth 16x16 words: Transformers
 662 for image recognition at scale, 2020. URL <https://arxiv.org/abs/2010.11929>.
- 663
- 664 Freedman, G. and Fattal, R. Image and video upscaling
 665 from local self-examples. *ACM Transactions on*
666 Graphics, 30(2):1–11, April 2011. doi: 10.1145/
 667 1944846.1944852. URL <https://doi.org/10.1145/1944846.1944852>.
- 668
- 669 Glasner, D., Bagon, S., and Irani, M. Super-resolution
 670 from a single image. In *2009 IEEE 12th International*
671 Conference on Computer Vision, pp. 349–356, 2009. doi:
 672 10.1109/ICCV.2009.5459271.
- 673
- 674 Han, Y., Du, X., and Yang, Z. Two-stage network for sin-
 675 gle image super-resolution. In *2021 IEEE/CVF Con-*
676 ference on Computer Vision and Pattern Recog-
677 nition Workshops (CVPRW), pp. 880–887, 2021. doi: 10.1109/
 678 CVPRW53098.2021.00098.
- 679
- 680 Horé, A. and Ziou, D. Image quality metrics: Psnr vs.
 681 ssim. In *2010 20th International Conference on Pattern*
682 Recognition, pp. 2366–2369, 2010. doi: 10.1109/ICPR.
 683 2010.579.
- 684
- 685 Huang, J.-B., Singh, A., and Ahuja, N. Single image
 686 super-resolution from transformed self-exemplars. In
 687 *2015 IEEE Conference on Computer Vision and Pat-*
688 tern Recognition (CVPR), pp. 5197–5206, 2015. doi:
 689 10.1109/CVPR.2015.7299156.
- 690
- 691 Hui, Z., Gao, X., Yang, Y., and Wang, X. Lightweight
 692 image super-resolution with information multi-distillation
 693 network. In *Proceedings of the 27th ACM Interna-*
694 tional Conference on Multimedia. ACM, October 2019. doi: 10.
 695 1145/3343031.3351084. URL <https://doi.org/10.1145/3343031.3351084>.
- 696
- 697 Johnson, J., Alahi, A., and Fei-Fei, L. Perceptual
 698 losses for real-time style transfer and super-resolution.
 699 In *Computer Vision – ECCV 2016*, pp. 694–711.
 700 Springer International Publishing, 2016. doi: 10.1007/
 701 978-3-319-46475-6_43. URL https://doi.org/10.1007/978-3-319-46475-6_43.
- 702
- 703 Kim, J., Lee, J. K., and Lee, K. M. Deeply-recursive
 704 convolutional network for image super-resolution. In
 705 *2016 IEEE Conference on Computer Vision and Pat-*
706 tern Recognition (CVPR), pp. 1637–1645, 2016. doi:
 707 10.1109/CVPR.2016.181.
- 708
- 709 Lai, W.-S., Huang, J.-B., Ahuja, N., and Yang, M.-H. Deep
 710 laplacian pyramid networks for fast and accurate super-
 711 resolution. In *2017 IEEE Conference on Computer Vision*
712 and Pattern Recognition (CVPR), pp. 5835–5843, 2017.
 713 doi: 10.1109/CVPR.2017.618.
- 714
- Ledig, C., Theis, L., Huszar, F., Caballero, J., Cunningham,
 715 A., Acosta, A., Aitken, A., Tejani, A., Totz, J., Wang, Z.,
 716 and Shi, W. Photo-realistic single image super-resolution
 717 using a generative adversarial network. pp. 105–114, 07
 718 2017. doi: 10.1109/CVPR.2017.19.
- 719
- Li, J., Yuan, Y., Mei, K., and Fang, F. Lightweight and accu-
 720 rate recursive fractal network for image super-resolution.
 721 In *2019 IEEE/CVF International Conference on Com-*
722 puter Vision Workshop (ICCVW), pp. 3814–3823, 2019.
 723 doi: 10.1109/ICCVW.2019.00474.
- 724
- 725 Lim, B., Son, S., Kim, H., Nah, S., and Lee, K. M. Enhanced
 726 deep residual networks for single image super-resolution.
 727 In *2017 IEEE Conference on Computer Vision and Pat-*
728 tern Recognition Workshops (CVPRW), pp. 1132–1140,
 729 2017. doi: 10.1109/CVPRW.2017.151.
- 730
- 731 Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin,
 732 S., and Guo, B. Swin transformer: Hierarchical vision
 733 transformer using shifted windows. In *2021 IEEE/CVF*
734 International Conference on Computer Vision (ICCV),
 735 pp. 9992–10002, 2021. doi: 10.1109/ICCV48922.2021.
 736 00986.
- 737
- 738 Lu, Z., Li, J., Liu, H., Huang, C., Zhang, L., and Zeng,
 739 T. Transformer for single image super-resolution. In
 740 *2022 IEEE/CVF Conference on Computer Vision and*
741 Pattern Recognition Workshops (CVPRW), pp. 456–465,
 742 Los Alamitos, CA, USA, jun 2022. IEEE Computer
 743 Society. doi: 10.1109/CVPRW56347.2022.00061. URL
<https://doi.ieee.org/10.1109/CVPRW56347.2022.00061>.
- 743
- 744 Stevens, E., Antiga, L., and Viehmann, T. *Deep learning*
 745 with PyTorch. Manning Publications, 2020.
- 746
- Tai, Y., Yang, J., and Liu, X. Image super-resolution via
 747 deep recursive residual network. In *2017 IEEE Con-*
748 ference on Computer Vision and Pattern Recog-
749 nation (CVPR), pp. 2790–2798, 2017. doi: 10.1109/CVPR.2017.
 750 298.
- 751
- 752 Timofte, R., Agustsson, E., Van Gool, L., Yang, M.-H.,
 753 Zhang, L., Lim, B., et al. Ntire 2017 challenge on sin-
 754 gle image super-resolution: Methods and results. In *The*
755 IEEE Conference on Computer Vision and Pattern Recog-
756 nation (CVPR) Workshops, July 2017.
- 757
- 758 Yang, G., Yu, S., Dong, H., Slabaugh, G., Dragotti, P. L.,
 759 Ye, X., Liu, F., Arridge, S., Keegan, J., Guo, Y., and
 760 Firmin, D. Dagan: Deep de-aliasing generative adver-
 761 sarial networks for fast compressed sensing mri recon-
 762 struction. *IEEE Transactions on Medical Imaging*, 37(6):
 763 1310–1321, 2018. doi: 10.1109/TMI.2017.2785879.
- 764

715 Yang, J., Wright, J., Huang, T. S., and Ma, Y. Image super-
716 resolution via sparse representation. *IEEE Transactions*
717 *on Image Processing*, 19(11):2861–2873, 2010. doi: 10.
718 1109/TIP.2010.2050625.

719
720 Yang, W., Zhang, X., Tian, Y., Wang, W., Xue, J.-H., and
721 Liao, Q. Deep learning for single image super-resolution:
722 A brief review. *IEEE Transactions on Multimedia*, 21(12):
723 3106–3121, 2019. doi: 10.1109/TMM.2019.2919431.

724 Zhang, K., Gao, X., Tao, D., and Li, X. Single image
725 super-resolution with non-local means and steering kernel
726 regression. *IEEE Transactions on Image Processing*, 21
727 (11):4544–4556, 2012. doi: 10.1109/TIP.2012.2208977.

728
729 Zhang, K., Zuo, W., and Zhang, L. Learning a single con-
730 volutional super-resolution network for multiple degra-
731 dations. In *2018 IEEE/CVF Conference on Computer*
732 *Vision and Pattern Recognition*, pp. 3262–3271, 2018a.
733 doi: 10.1109/CVPR.2018.00344.

734
735 Zhang, Y., Li, K., Li, K., Wang, L., Zhong, B., and Fu, Y.
736 Image super-resolution using very deep residual channel
737 attention networks. In *Proceedings of the European Con-*
738 *ference on Computer Vision (ECCV)*, September 2018b.

739
740 Zhang, Y., Tian, Y., Kong, Y., Zhong, B., and Fu, Y. Resid-
741 ual dense network for image super-resolution. In *2018*
742 *IEEE/CVF Conference on Computer Vision and Pattern*
743 *Recognition*, pp. 2472–2481, 2018c. doi: 10.1109/CVPR.
744 2018.00262.

745
746 Zhang, Y., Li, K., Li, K., Zhong, B., and Fu, Y. Residual
747 non-local attention networks for image restoration, 2019.
748 URL <https://arxiv.org/abs/1903.10082>.

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770 **A. Appendix**

771
772 The paper used ChatGPT as a paraphrasing tool. It is important to note that all content was originally written by me and
773 later applied chatGPT to paraphrase them, ensuring that no sentences were directly copied from the original paper or any
774 other source.

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824