

EE1-1.3 related: Lightweight and Efficient CNN In-loop Filter

Hao Zhang¹, Cheolkon Jung¹, Yang Liu², and Ming Li²

¹Xidian University, China

²OPPO, China

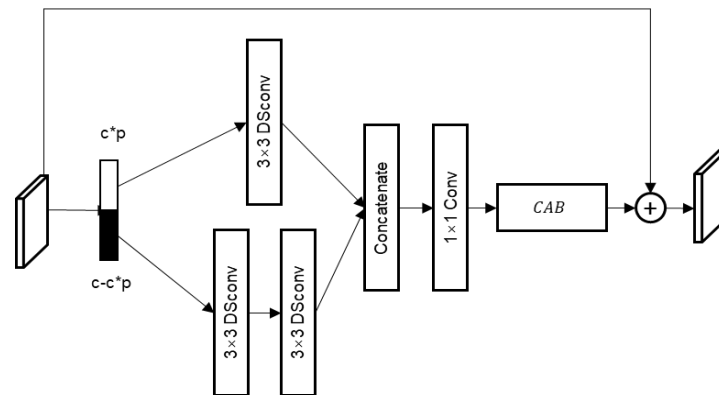
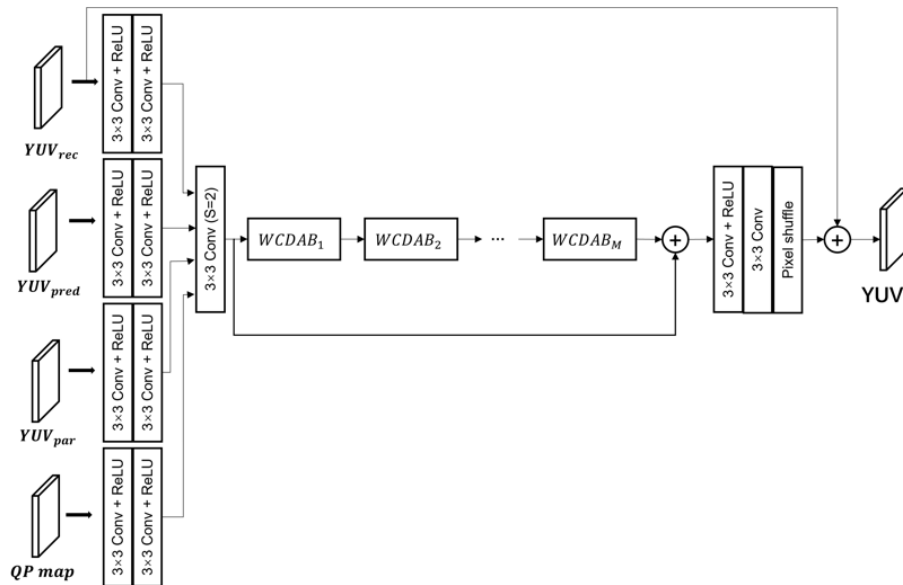
Introduction

In this contribution, the CNN-based in-loop filter proposed in JVET-AA0074 (AA0074 filter) is made lightweight and efficient.

- The number of skip connections, convolution layers and channels are reduced in each residual block, and correspondingly increased the number of residual blocks to ensure good learning ability
- A multi-stage training strategy is proposed to train the proposed network taking advantage of progressive learning(In AI configuration)
- Verify the performance of the improved network using AA0074 training strategy (In RA configuration)

Proposed Network in EE1-1.3

3



Proposed network

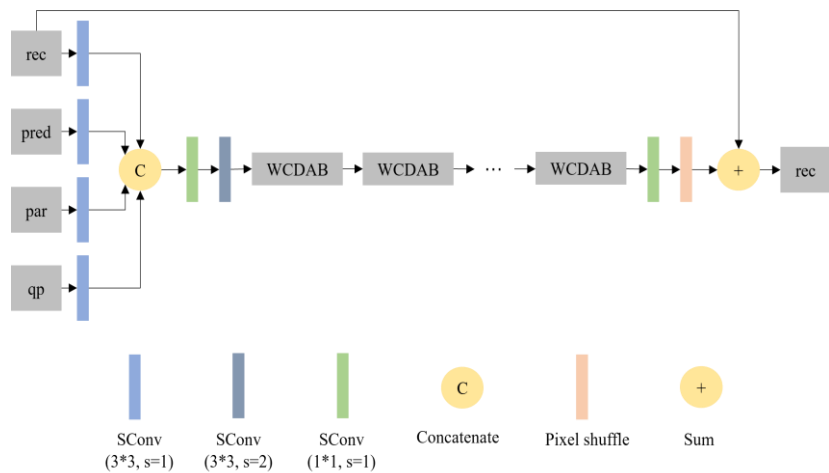
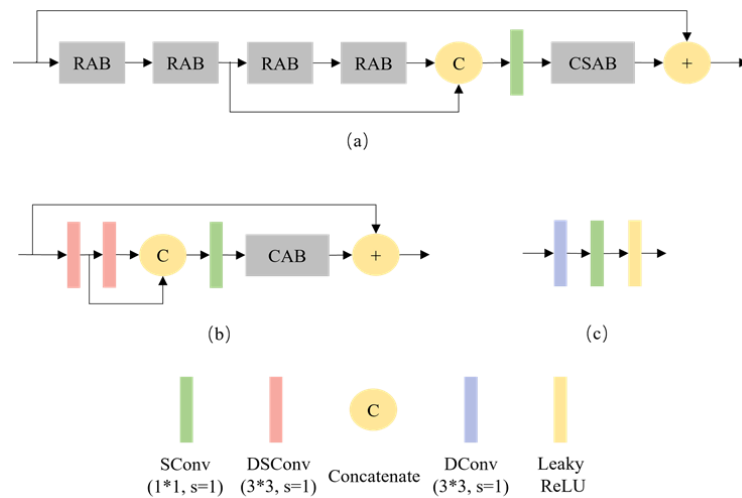


Illustration of the proposed NN-based in-loop filter. (SConv is the standard convolution)



(a) The architecture of WCDAB. (b) The architecture of RAB. (c) The architecture of Depthwise separable convolution. (SConv is standard convolution, DSConv is depthwise separable convolution, and DConv is depthwise convolution)

Training Detail

Dataset: DIV2K and BVI-DVC

Loss function: weighted L1 and L2 loss

$$L_{total} = 8 * L(Y_{out}, Y_{label}) + L(U_{out}, U_{label}) + L(V_{out}, V_{label})$$

Training data acquisition: Obtain reconstruction map, prediction map and partition map after LMCS, and obtain labels after ALF

Training strategy

In this contribution, a multi-stage training strategy is proposed to maximize the performance of the model during training.

- A parameter qp_dis is used, which represents the QP difference between the network input and the label. First, a smaller qp_dis is used to train the network, and then gradually increase qp_dis after the network converges while continuing to train the network.
- Since our loss function is multi-stage loss, we combine the loss function with the training strategy, i.e. a multi-stage training strategy.

Training strategy

Network performance at each training stage in terms of BD-rate (Class D in AI configuration).

| Stage | Y-PSNR | U-PSNR | V-PSNR |
|-----------------------------|--------|---------|---------|
| <i>qp_dis</i> =5 & L1 loss | -4.61% | -14.09% | -14.59% |
| <i>qp_dis</i> =5 & L2 loss | -5.60% | -14.84% | -17.12% |
| <i>qp_dis</i> =10 & L1 loss | -6.12% | -14.39% | -17.23% |
| <i>qp_dis</i> =10 & L2 loss | -6.63% | -14.65% | -16.90% |
| <i>qp_dis</i> =15 & L2 loss | -6.91% | -13.44% | -15.53% |

For specific *qp_dis*, the difference of image quality corresponding to Y channel is far greater than that of UV channel. This results in that the proposed training strategy cannot train both luminance and chrominance components at the same time.

Therefore, in actual training, *qp_dis* is only used for the Y channel. For the UV channel, its label uses the uncompressed image.

This problem can be solved by training luminance and chrominance filters respectively.

Network Information in Training Stage

| Network Information in Training Stage | | |
|---------------------------------------|--|---|
| Mandatory | GPU Type | GPU: GeForce RTX 3090 |
| | Framework: | PyTorch v1.9.0 |
| | Number of GPUs per Task | 1 |
| | | |
| | Epoch: | 800 |
| | Batch size: | 32 |
| | Training time: | ~200h/model |
| | Training data information: | BVI-DVC, DIV2K |
| | Training configurations for generating compressed training data (if different to VTM CTC): | VTM-11.0, QP {12,17,22, 27, 32, 37, 42} |
| Optional | Loss function: | Weighted L1 and L2 |
| | | |
| | Number of iterations | 3100 |
| | Patch size | 144x144 |
| | Learning rate: | 1e-4 |
| | Optimizer: | ADAM |
| | Preprocessing: | random cropping |
| | Other information: | |

Network Information in Inference Stage

| Network Information in Inference Stage | | |
|--|--|---|
| Mandatory | HW environment: | |
| | GPU Type | CPU only |
| | Framework: | Libtorch v1.9.0 |
| | Number of GPUs per Task | 0 |
| | | |
| | Number of Parameters (Each Model) | 0.78M |
| | Total Number of Parameters (All Models) | 0.78M |
| | Parameter Precision (Bits) | 32 |
| | Memory Parameter (MB) | 3.12M |
| | Multiply Accumulate (MAC)/pixel | 200K |
| Optional | | |
| | Total Conv. Layers | 64 Depthwise separable convolution 56 Common convolution |
| | Total FC Layers | 72 |
| | Total Memory (MB) | |
| | Batch size: | 1 |
| | Patch size | 144x144 |
| | Changes to network configuration or weights required to generate rate points | |
| | Peak Memory Usage (Total) | |
| | Peak Memory Usage (per Model) | |
| | Border handling | |
| | Other information: | |
| | | |

Performance of the proposed method in AI configuration

[illegible]

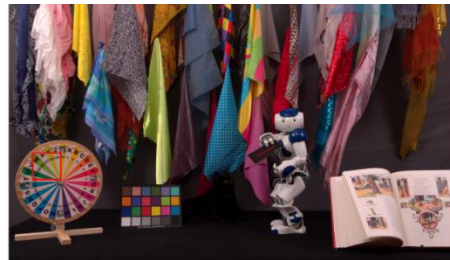
Performance of the proposed method in RA configuration

[illegible]

Visual comparison



(a)



(a)



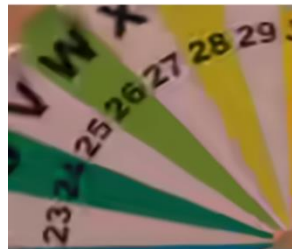
(b)



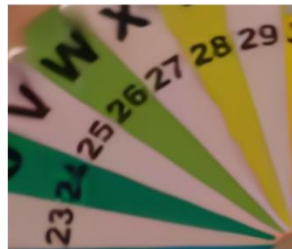
(c)



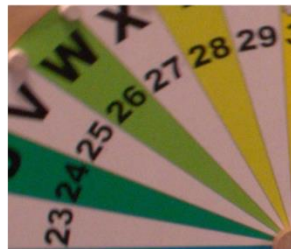
(d)



(b)



(c)



(d)

Visual quality comparison on the test datasets. (a) Video frame of Class A2-CatRobot. (b) Result of VTM-11.0_NNVC-2.0. (c) Result of the proposed filter. (d) Ground truth.

Conclusions

- Fine-tune the structure of AA0074 filter, make it more lighter.
- Propose a multi-stage training strategy to take advantage of progressive learning in the network training.
- The proposed CNN filter achieve average {7.08%, 12.46%, 12.75%} and {5.43%, 15.34%, 14.79%} BD-rate reductions over VTM-11.0-NNVC-2.0 for {Y, Cb, Cr} channels in AI and RA configurations.
- We recommend further study on the proposed CNN filter by including it in EE1.



THANK YOU!

