

An Image Coder With CNN Optimizations

Yundong Zhang^{*1}, Jianhua Hu^{*2}, Ming Li^{*2}, Changsheng Xia^{*2}, Jinwen Zan¹
Zhangming Huang², Weiwan Liao², Dekai Chen², Jing Nie²

¹Vimicro AI Chip Technology Corporation
The National Key Laboratory of Digital Multimedia Chip Technology
Room607,6/F, Shining Tower, No.35 Xueyuan Road, Haidian District, Beijing 100191 China

²Guangdong Vimicro Microelectronics Corporation
Building 16,Hengqin Financial District, Hengqin New Area, Zhuhai City, Guangdong Province, P.R.C
hu.jianhua@zxelec.com

Abstract

Convolutional neural networks (CNNs) has achieved great success in image processing and computer vision, especially in high level vision applications, such as classification and image compression. In this paper, CNN based optimizations have been proposed to improve the performance of an open source image coder, and the coding gain mainly comes from three modules: firstly, a classification CNN is employed to generate a region of interest (ROI) map, highlighting the part of the image containing more visual information that might be more sensitive to coding loss than other part, and thus guiding the bit allocation; secondly, a remedy CNN is introduced on the reconstructed YUV image, to learn and compensate for the coding loss; thirdly, adaptive loop filter(ALF) algorithm is applied to carry out color space conversion, and to minimize the color information loss during conversion. The improvement of the proposed optimizations, both objectively and subjectively, has been demonstrated on the CLIC validation data set.

1. Introduction

In recent years, image compression attracts increasing interest in image processing and computer vision due to its potential applications in many vision systems. Many image compression methods have been developed to efficiently compress the image such as JPEG, WebP, H.265 and H266. But compressed images and videos often suffer from block and ringing artifacts for areas with rich texture and sharp edges, especially when the bit-rate is relatively low. Human vision naturally focuses on familiar objects, and is particularly sensitive to distortions of these objects as compared to distortions of background details [1], so in this paper, we have improved the subjective quality of region of interesting by using a higher bit rate to encode, and lowering the

bit rate in background region to guarantee high coding efficiency. Also we proposed an approach to reduce artifacts after reconstructing image by applying deep CNNs in-loop filter. Moreover, due to the lossy of yuv-to-rgb translation, we proposed a adaptive loop filter filter in RGB color space.

2. The Proposed Compression Methods

In the proposed approach, we develop our codec based on the VTM/H266 platform. Three algorithms have been proposed to improve the performance of H266 codec: CNN based in-loop filter (CNNIF) after reconstruction frame; CNN based regions of interesting control different quantization parameters(QP); ALF to enhance the coding performance.

Our image compression framework is shown in Fig 1. Each image is split into block-shaped regions, and coded using intra prediction and other coding modules. The residual signal of intra prediction is transformed by a linear spatial transformation. The transform coefficients are then scaled, quantized and coded with entropy coding. Moreover, we apply different QP values according to the ROI mapping which is generated using CNN network to improving the subjection quality. The raw input image is RGB fromat, we transfer it into yuv420 format, and then encode to bitstreams with H266 encoder. After reconstructing the yuv420 data by the h266 decoder and interpolating chroma u/v by cubic interpolation, we get yuv444 data, and then a CNN network is used to filter luma y and chroma u/v data to get better image quality. We transform yuv444 into the RGB data. Finally we apply a RGB color space based adaptive loop filter to filter the RGB data and pack it to PNG format.

2.1. CNN ROI Use Different QP(ROIDQP)

In a traditional classification cnn, e.g.VGG-16, there are two fully-connected (non-convolutional) layers as the final layers of the network. The final layer has one neuron for every class in the training data, and the final step in the in-

*These authors share first-authorship

unit(CU) may have their own QP value, so we can set a QP offset (max_qp_offset) in H266 encode relative to frame QP. This is to say, when the ROI mapping value is high, we set a small QP value, otherwise a big QP value. Each of cu we can calculate a max QP offset (cur_cu_qp_offset) in the H265 encode is shown in Fig 2:

```

cu_delta = cu_average_value - global_average_value
max_diff = global_max_value - global_min_value
left_diff = global_average_value - global_min_value
right_diff = global_max_value - global_average_value
if abs(cu_delta) < k*0.5*max_diff
    cur_cu_qp_offset = 0
if (cu_average_value < global_average_value)
    cur_cu_qp_offset = cu_delta/left_diff
if (cu_average_value > global_average_value)
    cur_cu_qp_offset = cu_delta/right_diff

```

Where cu_average_value means the average ROI mapping value of current CU and global_average_value means the average ROI mapping value of all CU, k is response coefficient. what's more, global_max_value is max value of ROI mapping and global_min_value is the min value of ROI mapping.

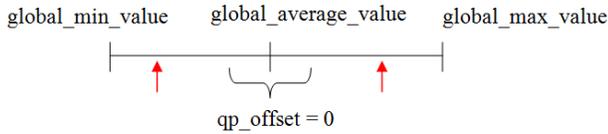


Figure 2. The calculation of CU QP offset

2.2. Reconstruction Frame Filtering With CNNs

In-loop filtering is an important technique in current mainstream codec for improving the quality of compressed image, so we design a novel CNN architecture and further boost the performance of in-loop filtering. The input of the network is the reconstruction frame in YUV color space, let's denote lumina of the H266 decoded image as Y. The proposed CNN model focuses on learning the residuals between the decoded Y and the ground truth lumina X of source image. Our goal is to fit a mapping function $X \approx F(Y)+Y$ that reverses image degradation due to H266 codec as much as possible. The whole architecture of our network is shown in Fig 3. We wish to learn the F by training a CNN, which conceptually consists of two operations: the feature extraction and image detail's reconstruction. The filter reconstruction frame model is a fully CNN network that consists of a set of convolution layers and non-linear layers cascades. To extract both the local and the global image features, all outputs of the hidden layers are concatenated at the end of feature extraction as skip connections from different layer domains. After concatenating all of the features, a simple reconstruction net is used to reconstruct the image details. Input Y is fed into the network,

residual is output from the second last layer, finally adding Y to form a $F(Y)+Y$ function.

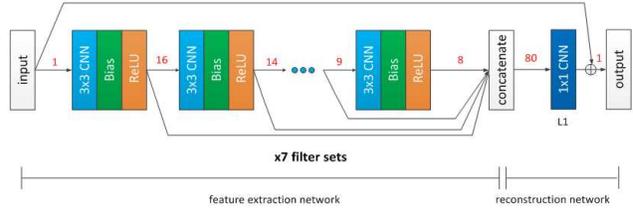


Figure 3. CNN filter structure

Although the reconstruction network is trained only on the luminance channel, there is great improvement in chroma uv reconstruction, which using the trained net of luminance. The model is not specifically designed to be an end-to-end solution. On the contrary, the proposed optimizes an end-to-end mapping. It is faster at speed because of less layers and channels.

The feature extraction part is responsible for extracting hidden features of the H266's reconstructed image. It consists of 7 consecutive 2d-convolution layers. We use Relu as activation function for each layer. We optimize the number of filters at each layer. The 7 filters with output feature num N are as follows: 16,14,12,11,10,9,8. All filter kernel size is k=3. The concatenation concatenates all layers' outputs, therefore the channels num in concatenation is 80.

The detail reconstruction part is responsible for outputting the residue of the image. Because of all of the hidden features are concatenated at the input layer of the reconstruction network, the dimension of input data is rather large. So we only use a 1x1 CNN filter as final mapping to generate output residual pixels data, not only reduces the dimensions of the previous layer for faster computation with less information loss, but also adds more nonlinearity to enhance the potential representation of the network [3]. The output residual is then added with the original input Y.

The training inputs are Y patches decoded by H266 decoder using qp36. We trained our model with BSD200 dataset[4]. Patch size is 64x64. Using MSE as the loss function favors a high PSNR.

3. ALF In RGB Color Space

Because there is a loss of YUV420 to YUV444 and YUV444 color to RGB color space, we apply ALF method to reduce loss during translation in RGB color space and R, G, B each have their own filter coefficients and master switch. The ALF was designed to minimize the error between the original frame and the reconstructed frame. It use different method in reconstructed image to divide the pixels into different classification. The pixels of the same category will share the same filter coefficients. ALF switch for each Block and its taps are decided based on the rate distortion

References

- [1] M. Jiang, S. Huang, J. Duan, and Q. Zhao *Perceptual adaptation of objective video quality metrics*, in Proc. Ninth International Workshop on Video Processing and Quality Metrics (VPQM), 2015.
- [2] G.Griffin, A.Holub, and P Perona, *Caltech-256 object category dataset*, 2007.
- [3] Jin Yamanaka¹, Shigesumi Kuwashima¹ and Takio Kurita². *Fast and Accurate Image Super Resolution by Deep CNN with Skip Connection and Network in Network.*, 24th International Conference On Neural Information Processing (ICONIP 2017), 2017. 9
- [4] Arbelaez, P., Maire, M., Fowlkes, C., Malik, J. *Contour Detection and Hierarchical Image Segmentation*, IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 33, no. 5, pp. 898–916 (2011)
- [5] Xinfeng Zhang, Ruiqin Xiong, Siwei Ma, Wen Gao *Adaptive loop filter with temporal prediction*, Picture Coding Symposium, May 7-9, 2012.