# ADAPTING LEARNED IMAGE CODECS TO SCREEN CONTENT VIA ADJUSTABLE TRANSFORMATIONS

*H. Burak Dogaroglu*[1]    *A. Burakhan Koyuncu*[1,2]    *Atanas Boev*[2]
*Elena Alshina*[2]    *Eckehard Steinbach*[1,3]

[1] School of Computation, Information and Technology, TU Munich
[2] Huawei Munich Research Center
[3] Munich Institute of Robotics and Machine Intelligence
Email: burak.dogaroglu@tum.de

## ABSTRACT

As learned image codecs (LICs) become more prevalent, their low coding efficiency for out-of-distribution data becomes a bottleneck for some applications. To improve the performance of LICs for screen content (SC) images without breaking backwards compatibility, we propose to introduce parameterized and invertible linear transformations into the coding pipeline without changing the underlying baseline codec's operation flow. We design two neural networks to act as prefilters and postfilters in our setup to increase the coding efficiency and help with the recovery from coding artifacts. Our end-to-end trained solution achieves up to 10% bitrate savings on SC compression compared to the baseline LICs while introducing only 1% extra parameters.

*Index Terms*— screen content coding, learned image compression, prefilter, postfilter, deep learning

## 1. INTRODUCTION

As available compute power of consumer grade hardware increases, artificial intelligence (AI) products start to replace more classical solutions. Image compression is one of these fields where the neural network-based learned image codecs (LICs) [1, 2, 3, 4, 5, 6, 7] can produce similar level or more efficient encodings than their classical counterparts. This opens up opportunities to develop more efficient AI-based coding standards for image and video compression.

Screen content (SC) images have unique characteristics. Unlike natural images where the pixels are produced by sensors capturing natural phenomena such as light, the SC pixels are mostly artificially generated. This causes SC to have drastically different color and texture patterns [8]. They usually consist of sharp edges with reduced colorspaces. Some examples for SC images are graphical user interfaces, large bodies of letters and text, computer simulation results and so on. It is also possible to have natural images embedded in a SC context, e.g. game streaming where the user's video feed is combined with artificially generated game images.

The SC compression is often overlooked in the domain of LICs due to lack of representation in most datasets and lack of effort in the exploitation of unique redundancies. It is also very hard to update LICs after they are standardized since any change in network weights may break the backwards compatibility of the bitstreams and it's not feasible to store multiple copies of large networks for every iteration. For these reasons, we research solutions in adapting LICs to out-of-distribution data like SC without breaking the compatibility of already generated encodings.

In that regard, we develop new modules around the codec that can be disabled with a signal. This setup allows us to improve coding efficiency of LICs without retraining the codec itself. For this purpose, we develop two neural networks that perform as the prefilter and postfilter around the codec's input and output. To enhance the performance on SC compression, we also propose the use of linear forward and inverse transformations.

The prefilter model focuses on creating better compressible images through addition of spatial modulations to the source image. On the other hand, the postfilter model aims to remove the modulations and reduce the artifacts created by the forward transformation and the prefilter model. We train our models end-to-end so that they can learn to use information transmitted through the codec more efficiently.

## 2. BACKGROUND

### 2.1. Classical screen content compression

Rise of digitalization has increased the demand for SC compression. As a response, video codec standards such as HEVC [9] and VVC [10] adopted specialized extensions. Among these, most prominent ones are the Intra Block Matching (IBM) and Palette Mode Coding (PMC) [11, 12].

In IBM, for each coding unit, the intra encoder searches for similar or repeating patches in the already encoded parts

**Fig. 1**: Architecture of the residual network based on MBConv layers. This implementation depicts $L = 8, C = 32$ setup.

of the image. This allows the encoder to skip encoding large chunks of pixels and just transmit the position of the similar patch to the decoder. In turn, the decoder just copies the referenced block from already decoded coding units. This technique increases the coding efficiency especially when the source image has large chunks of repeating blocks, such as graphical user interfaces, pixel art or gaming.

On the other hand, PMC focuses on exploiting the color space redundancy of SC images. The algorithm encodes the set of unique colors in a given block and replaces the pixels with the indexes corresponding to their colors. Instead of transmitting the color channels of the image, we encode the palette and send the indexes alongside it. This way, we can get lower bitrates for source images with narrow color spaces.

### 2.2. Learned image compression

After the invention of differentiable entropy bottleneck layers and end-to-end trainable autoencoder codecs [13], learned image compression has seen growing interest from the researchers. Later, hyper-encoders and hyper-decoders were introduced to decrease the spatial redundancy on latent codes [1]. This was further improved by allowing a context model to process the reconstructed latents [2]. Designing context models with better priors were shown to reduce required bitrate to transmit images [3]. Researchers began to test transformer based context models to replace the autoregressive architectures [4]. Since transformer context models required more compute, there has been also some attempts to make them more efficient via routing attention to selected channels [5].

### 2.3. Hybrid approaches

Tao et al. [14] combined two neural network models around a classical image codec and trained them together to improve the coding efficiency of the underlying codec. To get around the non-differentiable codec's lack of backpropagation, they used auxilary loss functions to train the networks separately.

Recently, Guleryuz et al. [15] solved the same problem by switching the actual codec with a differentiable image proxy during training time. This idea was also combined with applying transformations such as downscaling or quantization to enhance codec capabilities [16].

## 3. METHODOLOGY

Training a LIC is computationally expensive. In this work, we attempt to improve the SCC efficiency of an existing LIC without retraining of the whole codec, but using pre- and post-processing operations instead. The codec is treated as a black box and its weights are frozen during training. This helps our solution to be implemented easily on top of any differentiable codec while keeping the produced bitstream compatible with the underlying software.

### 3.1. Proposed solution

Our solution requires four extra modules on top of the codec. In the encoder side, we have a linear forward transformation $T$ and Compact Representation (CR) module with a neural network. Following a similar structure, the decoder consists of a linear inverse transformation $T^{-1}$ and the Reconstruction Stage (RS) module. Fig. 2 depicts the operational flow of this solution.

### 3.2. Linear transformations

In our setup, we can't change the internal flow of the codec. That means, we must work with the same kind of inputs and outputs as the codec. This puts a requirement on the valid transformations such that they must be compatible with 3-channel RGB images. In cases where the exact inverse transformation is not possible, we can approximate it with the Moore - Penrose inverse [17] of the forward transformation.
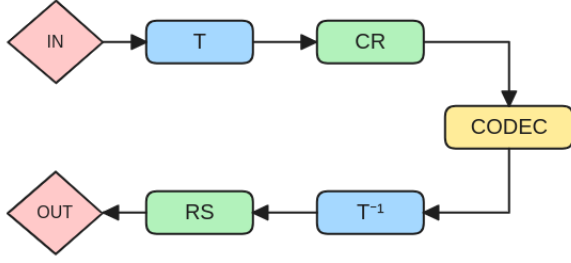
**Fig. 2**: Operational diagram of the proposed pipeline.

In this work, we consider 3 candidate transformations. Desaturation transformation takes an RGB image and linearly interpolates between the grayscale version according to the formula $\hat{x} = x + (1 - \alpha)Gray(x)$ where $x$ is the RGB input, $Gray$ is the 3-channel grayscale color conversion and $\alpha$ is the saturation level. Its inverse can easily be computed by just changing $\alpha$ to its multiplicative inverse.

The PCA downscaling transformation first converts an image to its PCA colorspace consisting of the principal channel (PC), side channel 1 (SC1) and side channel 2 (SC2) according to the order of lowest standard deviation. Then, it downscales the side channels by the parameters $d_{SC1}$ and $d_{SC2}$ respectively where the downscaled images have image dimensions divided by these values. While bicubic transformation is used during downscaling, we preferred to upscale using the nearest neighbor approximation so the color channels become the same resolution as the original image. Once the image is constructed back, we remap the colors to RGB using the inverse PCA color transformation.

Similarly, we also evaluate PCA quantization transform which uses the same color transformation as before. However, instead of downscaling, we quantize the side channel amplitudes with $Q_{SC1}$ and $Q_{SC2}$ bits by replacing them with cluster centers produced by the k-means [18] algorithm. Finally, the quantized image is remapped back to the RGB space to be consumed by the LICs.

### 3.3. Network architecture

Our CR and RS modules use the exact same architecture with same number of layers and channels. We use a ResNet architecture [19] where the main processing blocks are MBConv layers [20]. The network consists of a convolution layer with stride of 2 that maps RGB images to $C$ channels, $L$ number of MBConv layers where the last layer doesn't squeeze its channels, another convolution layer that maps the unsqueezed $4C$ channels to 12 so that we can finish with pixel shuffling [21] with a stride of 2 to get the original resolution again. To stabilize the training and let network focus on only extracting modulations, we have an end-to-end residual connection between the input and output of the module. This structure is illustrated in Fig. 1.

### 3.4. Dataset

We used the SC dataset collected for the JPEGAI [22] standardization effort, denoted as JPEGAI-SC hereinafter. It has approximately 3 thousand training images and 1 thousand validation images. We use the validation set in our analysis section since we don't use it in any part of training. The dataset consists of 5 categories of SC images: AI-generated, gaming, screenshots, 3D renderings, illustrations.

### 3.5. Training scheme

Our models were trained for 50 epochs where the codecs' weights were frozen. The batch size in our trainings was set to 8 and we worked with center cropped 256x256 images both for training and testing. We used Adam [23] optimizer with $\gamma = 10^{-4}, \beta_1 = 0.9, \beta_2 = 0.999$. The codec implementations are from the CompressAI [24]. We used the same Rate - Distortion loss with the same $\lambda$ parameters as the baseline codecs [1, 2, 3]. We optimized the Mean Squared Error (MSE) as the codec variants we use.

## 4. EXPERIMENTAL RESULTS

We performed a series of experiments to find optimal parameters and settings for various questions. In most cases, we only trained for quality points between 2 and 5 due to limited access to compute resources and 4 quality points are enough to compute Bjontegaard-Delta (BD) rate with Akima interpolation [25, 26].

### 4.1. Compressibility of transformations

Since we are interested in creating bitrate reductions through reducing the information of the images with transformations, we designed an experiment where the transformed images are passed through a LIC and we measure their compressed sizes. We used the Ballé encoder for quality points between 2 and 5 and measured averaged the bitrate differences $\Delta BPP$ of the augmented images to their baselines.

In Fig. 3, we can see that while desaturation transformation gets smaller bitrates compared to the Ballé codec, the information loss created by other transformations are not well received by the codec and does not create noticable savings even though they introduce more noise to the system. This indicates that only the desaturation transformation have any potential for performance gains for SC compression, so we focus on only desaturation transformation in the following experiments.

### 4.2. Effect of desaturation magnitude

We introduced the inverse transformation on the decoder side and measured end-to-end reconstruction performance of the pipeline without the neural network modules. From Fig. 4,
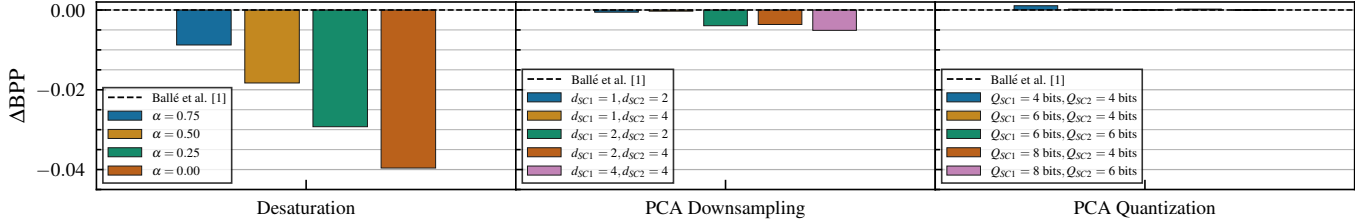
**Fig. 3**: Compressibility of desaturation, PCA downsampling and PCA quantization. Lower bitrate indicates a potential for performance gain. While desaturation clearly allows us to compress images into smaller files, the PCA downsampling and PCA quantization transformations have no visible effect on the bitrate difference.
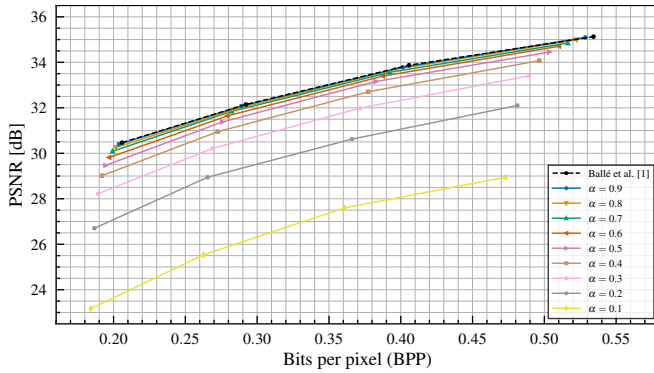


**Fig. 4**: Rate-Distortion plot of different desaturation levels on top of Ballé baseline codec. The reconstructed image quality drops higher for stronger transformations. However, weaker transformations produce almost little to none bitrate gains, indicating a sweet spot for the quality - bitrate tradeoff.



**Fig. 5**: Rate-Distortion plot for our solution with baselines as Ballé, Minnen and Cheng codecs.

it can be seen that the transformations alone hurt the performance more as the saturation level $\alpha$ gets smaller. Assuming our networks can produce 0.5 - 1 dB PSNR gain on top of the codecs, we opt to focus on a saturation level of $\alpha = 0.8$ in the following experiments.

### 4.3. Performance analysis

Finally, we activate the neural network modules and train the whole pipeline in an end-to-end manner. We repeat this for all available quality points in Ballé et al. [1], Minnen et al. [2] and Cheng et al. [3] baseline codecs. For ease of comparison, we used the same neural network models of $L = 8$ layers and $C = 32$ channels in all cases. Fig. 5 shows that in all cases, our modules improved over the baseline codecs.

### 4.4. Computational complexity

We performed an analysis to find out the relative costs and benefits of using our pipeline compared to the baseline codecs. According to Table 1, we see that it is possible to adapt learned image codecs to out-of-distribution SC data by using just %1 more parameters. We also measured the
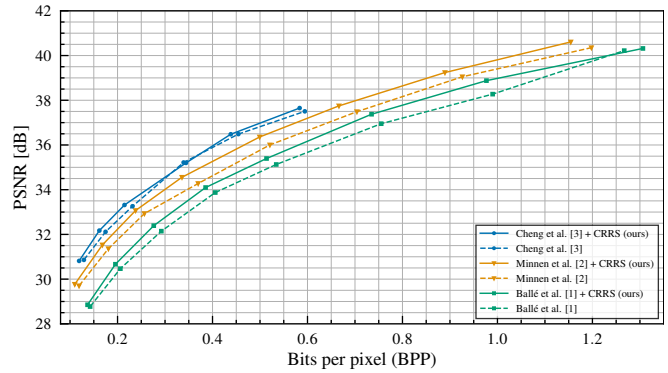
Multiply-Accumulate Operations per pixel (MAC/px) of the end-to-end pipelines and our modules introduce few extra computations to increase the coding efficieny of SC images. We follow the technique proposed in JPEGAI [22] standardization to find the MAC/px of the models. The reported values are from the highest quality point baseline codec available in the CompressAI [24] library.

### 4.5. Ablation study

Our ablation study has two steps. In the first step, we kept every module active and tested the performance of changing number of channels $C$ and number of layer $L$ for our neural network models. In the second step, we fixed the neural network architecture and only measured the performance contribution of modules separately and together. Table 2 summarizes our findings from these experiments.

From the first stage ablation experiments, in general, increasing the model complexity helps with the performance. However, increasing the number of layers from 8 to 10 for 32 channels is an exception to this observation. We hypothesize that relatively deeper models may require longer training or better normalization between submodules to be effective. We don't see this effect for higher number of channels. Since the best performing model was created with $L = 5, C = 64$, we decided to use it for the second stage of the ablation study.

| Model | # Param. | MAC/px | BD-Rate (%) ↓ | |
|---|---|---|---|---|
| Ballé et al. [1] | 11.82M | 418K | 0.0 | |
| + CRRS | 12.02M | 460K | -9.0 | |
| Minnen et al. [2] | 25.51M | 450K | -20.6 | (0.0) |
| + CRRS | 25.71M | 492K | -28.9 | (-10.5) |
| Cheng et al. [3] | 26.60M | 927K | -36.9 | (0.0) |
| + CRRS | 26.80M | 969K | -39.6 | (-4.5) |

**Table 1**: Model complexity and performance analysis relative to the Ballé baseline codec. For Minnen and Cheng pipelines, relative BD-Rate gains are reported in parantheses. Lower BD-Rate is better.

| Name | $\alpha$ | CR | RS | $L$ | $C$ | $\Delta$MAC/px | BD-Rate (%) ↓ |
|---|---|---|---|---|---|---|---|
| Ballé [1] | ✗ | ✗ | ✗ | ✗ | ✗ | 0 | 0.0 |
| | ✓ | ✓ | ✓ | 5 | 64 | 102K | **-11.5** |
| I | ✓ | ✓ | ✓ | 3 | 64 | 65K | -9.4 |
| | ✓ | ✓ | ✓ | 10 | 32 | 51K | -6.2 |
| | ✓ | ✓ | ✓ | 8 | 32 | 42K | -8.0 |
| | ✗ | ✓ | ✗ | 5 | 64 | 51K | -0.9 |
| | ✗ | ✗ | ✓ | 5 | 64 | 51K | -4.4 |
| II | ✗ | ✓ | ✓ | 5 | 64 | 102K | -7.4 |
| | ✓ | ✗ | ✗ | 5 | 64 | 0 | 2.0 |
| | ✓ | ✗ | ✓ | 5 | 64 | 51K | -2.9 |

**Table 2**: Ablation study with Ballé baseline codec on quality points 2-5. We investigate the effect of different model size in (I) and the effect of each proposed modules such as the desaturation, CR and RS in (II). We set the desaturation level to 0.8 for these experiments.

We disabled some modules of our pipeline with Ballé codec with neural networks of size $C = 64$ channels and $L = 5$ layers. We measured the performance on quality points between 2 and 5. We found out that the CR model alone doesn't produce much gain while RS model alone is reasonably performant. However, best results come when two modules are combined. The desaturation transformation hurts the performance alone but the introduction of RS model on top of this results in bitrate savings. When the full pipeline is activated, we see that every module contributes and we get the best possible performance.

In our experiments, we also tested the effect of desaturation coefficient $\alpha$ on the performance of our pipeline. From Fig. 6, it can be seen that there is an inverse log-linear correlation between the desaturation coefficient and the BD-Rate when no neural networks are activated in the pipeline. We investigated the sudden drop of the curve around $\alpha = 0.9$
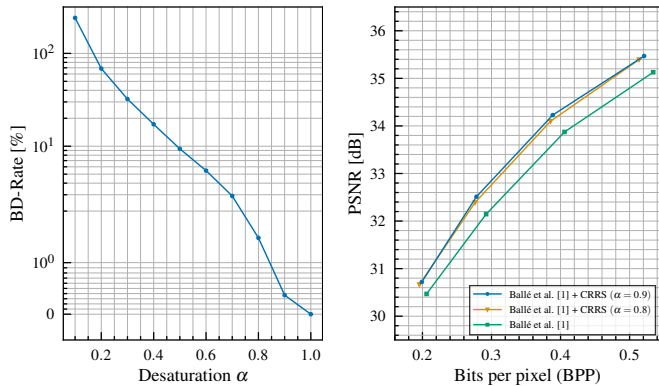


**Fig. 6**: Effect of desaturation coefficient $\alpha$ on top of Ballé baseline codec. Left image shows the BD-Rate without the neural network modules and the right figure depicts the Rate-Distortion plot of $\alpha = 0.8$ and $\alpha = 0.9$ with the neural network modules.

by training our neural pipeline to see whether it can produce better gain than its $\alpha = 0.8$ counterpart. We observed that the model with $\alpha = 0.9$ reaches a slightly higher performance gain of 1.1% compared to the model with $\alpha = 0.8$. This shows that the desaturation coefficient can be further optimized for a higher performance.

## 5. CONCLUSION

In this paper, we tried to solve the adaptation of LICs to SC data without breaking the codec's self consistency. To do so, we proposed two neural networks, CR and RS models to act as the prefilter and the postfilter. We also demonstrated that introduction of some forward and inverse transformations help with the codec efficiency. Specifically, the desaturation transformation can reduce the bitrate required to compress an image via destroying some information. However, the neural networks are able to recover from this transformation for reasonable magnitude. It can be clearly seen that with a little bit of extra compute cost, it is possible to adapt LICs to unseen domains without retraining the codecs themselves.

While the performance improvement shown here is promising, LICs still require lots of compute power compared to the classical codecs. This creates a barrier of adaptation for standardized LICs and introducing new networks may not be favorable in some use cases. Creating more efficient and faster LICs is research direction with lots of real world applications. Another interesting approach could be in adapting this approach to other special domains such as the medical, microscopy, multispectral or space imaging where more efficient codecs are still in need. These new domains may require new specialized transformations and architectures. Finally, it is also a promising research direction to find optimal transformation coefficients per image to get even more specialized gains.

# 6. REFERENCES

[1] Johannes Ballé, David Minnen, Saurabh Singh, Sung Jin Hwang, and Nick Johnston, "Variational image compression with a scale hyperprior," *arXiv preprint arXiv:1802.01436*, 2018.

[2] David Minnen, Johannes Ballé, and George D Toderici, "Joint autoregressive and hierarchical priors for learned image compression," *Advances in neural information processing systems*, vol. 31, 2018.

[3] Zhengxue Cheng, Heming Sun, Masaru Takeuchi, and Jiro Katto, "Learned image compression with discretized gaussian mixture likelihoods and attention modules," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 7939–7948.

[4] A. Burakhan Koyuncu, Han Gao, Atanas Boev, Georgii Gaikov, Elena Alshina, and Eckehard Steinbach, "Contextformer: A transformer with spatio-channel attention for context modeling in learned image compression," in *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XIX*. Springer, 2022, pp. 447–463.

[5] A. Burakhan Koyuncu, Panqi Jia, Atanas Boev, Elena Alshina, and Eckehard Steinbach, "Efficient contextformer: Spatio-channel window attention for fast context modeling in learned image compression," *arXiv preprint arXiv:2306.14287*, 2023.

[6] Dailan He, Ziming Yang, Weikun Peng, Rui Ma, Hongwei Qin, and Yan Wang, "Elic: Efficient learned image compression with unevenly grouped space-channel contextual adaptive coding," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 5718–5727.

[7] Jinming Liu, Heming Sun, and Jiro Katto, "Learned image compression with mixed transformer-cnn architectures," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 14388–14397.

[8] Yingbin Wang, Xin Zhao, Xiaozhong Xu, Shan Liu, Zhijun Lei, Mariana Afonso, Andrey Norkin, and Thomas Daede, "An open video dataset for screen content coding," in *2022 Picture Coding Symposium (PCS)*. Dec. 2022, IEEE.

[9] Gary J. Sullivan, Jens-Rainer Ohm, Woo-Jin Han, and Thomas Wiegand, "Overview of the high efficiency video coding (hevc) standard," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 22, no. 12, pp. 1649–1668, 2012.

[10] Benjamin Bross, Ye-Kui Wang, Yan Ye, Shan Liu, Jianle Chen, Gary J. Sullivan, and Jens-Rainer Ohm, "Overview of the versatile video coding (vvc) standard and its applications," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 10, pp. 3736–3764, 2021.

[11] Shan Liu, Xiaozhong Xu, Shawmin Lei, and Kevin Jou, "Overview of hevc extensions on screen content coding," *APSIPA Transactions on Signal and Information Processing*, vol. 4, pp. e10, 2015.

[12] Xiaozhong Xu and Shan Liu, "Overview of screen content coding in recently developed video coding standards," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 2, pp. 839–852, 2022.

[13] Johannes Ballé, Valero Laparra, and Eero P Simoncelli, "End-to-end optimized image compression," *arXiv preprint arXiv:1611.01704*, 2016.

[14] Wen Tao, Feng Jiang, Shengping Zhang, Jie Ren, Wuzhen Shi, Wangmeng Zuo, Xun Guo, and Debin Zhao, "An end-to-end compression framework based on convolutional neural networks," in *2017 Data Compression Conference (DCC)*, 2017, pp. 463–463.

[15] Onur G. Guleryuz, Philip A. Chou, Hugues Hoppe, Danhang Tang, Ruofei Du, Philip Davidson, and Sean Fanello, "Sandwiched image compression: Wrapping neural networks around a standard codec," in *2021 IEEE International Conference on Image Processing (ICIP)*, 2021, pp. 3757–3761.

[16] Onur G. Guleryuz, Philip A. Chou, Hugues Hoppe, Danhang "Danny" Tang, Ruofei Du, Philip Davidson, and Sean Fanello, "Sandwiched image compression: Increasing the resolution and dynamic range of standard codecs," in *2022 Picture Coding Symposium (PCS)*, 2022, Best Paper Finalist.

[17] Roger Penrose, "A generalized inverse for matrices," *Mathematical Proceedings of the Cambridge Philosophical Society*, vol. 51, no. 3, pp. 406–413, July 1955.

[18] Stuart Lloyd, "Least squares quantization in pcm," *IEEE Transactions on Information Theory*, vol. 28, no. 2, pp. 129–137, Mar. 1982.

[19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2016, IEEE.

[20] Mingxing Tan and Quoc Le, "EfficientNet: Rethinking model scaling for convolutional neural networks," in *Proceedings of the 36th International Conference*

*on Machine Learning*, Kamalika Chaudhuri and Ruslan Salakhutdinov, Eds. 09–15 Jun 2019, vol. 97 of *Proceedings of Machine Learning Research*, pp. 6105–6114, PMLR.

[21] Wenzhe Shi, Jose Caballero, Ferenc Huszár, Johannes Totz, Andrew P Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang, "Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 1874–1883.

[22] João Ascenso, Elena Alshina, and Touradj Ebrahimi, "The jpeg ai standard: Providing efficient human and machine visual data consumption," *IEEE MultiMedia*, vol. 30, no. 1, pp. 100–111, 2023.

[23] Diederik P Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[24] Jean Bégaint, Fabien Racapé, Simon Feltman, and Akshay Pushparaja, "Compressai: a pytorch library and evaluation platform for end-to-end compression research," *arXiv preprint arXiv:2011.03029*, 2020.

[25] Gisle Bjøntegaard, "Calculation of average psnr differences between rd-curves," 2001.

[26] Christian Herglotz, Matthias Kränzler, Ruben Mons, and André Kaup, "Beyond bjøntegaard: Limits of video compression performance comparisons," in *2022 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2022, pp. 46–50.