



MOTION IMAGING JOURNAL

Covering Emerging Technologies for the Global Media Community



Understanding Banding—Perceptual Modeling and Machine Learning Approaches

By Hojatollah Yeganeh, Kai Zeng, and Zhou Wang

Introdução:

O efeito *Banding* é um artefato visual irritante que frequentemente aparece em várias etapas ao longo da cadeia de aquisição de vídeo, produção, distribuição e exibição. Com a crescente popularidade do conteúdo de Ultra Definição (UHD), esse efeito tem recebido atenção cada vez maior devido ao seu forte impacto negativo na experiência do espectador. Neste artigo vamos ver dois tipos de estruturas para detectar o *Banding*, o primeiro é orientado pelo conhecimento e é construído com base em modelos computacionais que levam em consideração as características do sistema visual humano (SVH), já o segundo é orientado por dados e baseia-se em métodos de aprendizado de máquina através do treinamento de Redes Neurais Profundas (DNNs). Confesso que gostei muito mais da primeira abordagem, não é por que esta na moda, que vamos abandonar o certo pelo duvidoso. Mas e você, qual sistema prefere? Leia o artigo e vamos debater mais sobre esse assunto, escreva-me: tvdigitalbr@gmail.com

Tom Jones Moreira

Abstract

Banding is an annoying visual artifact that frequently appears at various stages along the chain of video acquisition, production, distribution, and display. With the thriving popularity of ultrahigh definition (UHD), high-dynamic range (HDR), wide-color-gamut (WCG) content, and the increasing user expectations that follow, the banding effect has been attracting increased attention for its strong negative impact on viewer experience in visual content that would otherwise have nearly perfect quality. Here, we present two different types of frameworks to detect the banding artifact. The first is knowledge-driven and is built upon computational models that account for the characteristics of the human visual system (HVS), the content acquisition, production, distribution, and display processes, and the interplay between them. The second is data-driven and is based on machine learning methods, by training deep neural networks (DNNs) with large-scale datasets.

Keywords

Banding impairment, contouring impairment, human visual system (HVS), perceived video quality

Introduction

Banding is a visual artifact that appears frequently at many stages in video acquisition, production, distribution, and display systems. Banding typically appears as perceived discontinuities or false contours in large and smooth image regions of slow color or intensity gradients. An example is given in **Fig. 1**, where a severe banding effect is observed in the sky. Although

heavy video compression is a potential source of banding, banding may also occur in the absence of any lossy compression and may create annoying visual quality degradations in otherwise pristine quality images or video content.

What often frustrates many industrial practitioners is that simply increasing the bit-depth or bitrate of a video does not necessarily lead to the removal or even reduction of banding. Indeed, with the recent accelerated growth of ultrahigh-definition (UHD), high dynamic range (HDR),

wide-color-gamut (WCG) content production, distribution services, and consumer display devices, severe banding occurs even more frequently than before and the visual effect is often much stronger. This is because UHD/HDR/WCG content typically covers a wider range of luminance levels and color variations than those of the traditional standard dynamic range (SDR) content. This, together with the limited and varying capabilities of display devices, creates major challenges to maintain smooth visual transitions simultaneously across all luminance levels and color variations. Significant effort has been made over the years on removing or reducing banding effects in video distributions. Depending on where these banding

reduction techniques are applied, they may be classified into preprocessing, post-processing, and banding-aware encoding methods. However, without having a reliable objective measure to detect banding, improving dithering or debanding efforts are quite cumbersome and directionless. Therefore, the industry is in urgent need of innovative approaches that are able to detect, control, and remove/reduce banding in an automated fashion.

Significant effort has been made over the years on removing or reducing banding effects in video distributions. Depending on where these banding reduction techniques are applied, they may be classified into preprocessing, postprocessing, and banding-aware encoding methods. However, without having a reliable objective measure to detect banding, improving dithering or debanding efforts are quite cumbersome and directionless. Therefore, the industry is in urgent need of innovative approaches that are able to detect, control, and remove/reduce banding in an automated fashion.

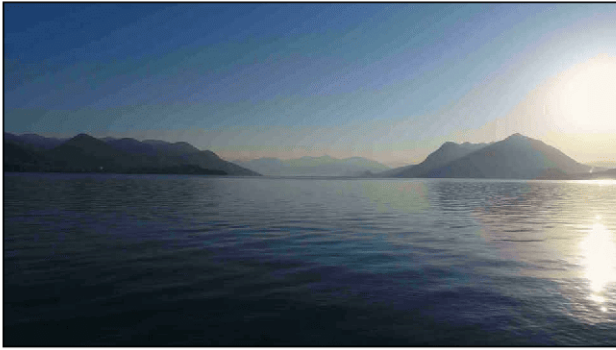


FIGURE 1. Sample image of visual banding effect (in the sky region).

Automatic or objective image/video quality assessment (IQA/VQA) has been a highly active topic in the past two decades. However, popular IQA/VQA methods such as peak signal to noise ratio (PSNR),¹ structural similarity index (SSIM),² multi-scale SSIM (MS-SSIM),³ SSIMPLUS,⁴ and Video Multimethod Assessment Fusion (VMAF),⁵ often require full access to a pristine reference when assessing a test image or video—and pristine references are rarely available in real-world testing environments. Moreover, although the quality maps created by SSIM types of approaches^{2–4} often successfully capture local banding artifacts, the overall assessment of these approaches mixes all visual distortion types together and there is no simple mechanism to single out the banding effect. Therefore, it is desirable to develop novel IQA/VQA methods dedicated to detecting and assessing banding without access to a pristine-quality original image/video as a reference.

Recently, two substantially different types of approaches have shown notable success at banding detection. The first is based on the domain knowledge gained through deep and thorough understanding of the human visual system (HVS) and the video acquisition, production, distribution, and display processes. Computational models derived from such domain knowledge are then combined to construct an overall banding detection and assessment model.^{6,7}

In contrast to the first type of domain knowledge-driven methods, the second type of approaches are data-driven, with no or little domain knowledge assumed. Instead, a large number of images/videos and their ground-truth labels (with or without banding) are collected, and machine learning methods are then used to train black-box models such as the deep neural networks (DNNs) using the image/video dataset, so that the learned model may make good banding predictions on unseen image/video content.

Knowledge-driven Method

There is rich literature on computational modeling of HVS characteristics and the individual components in the sophisticated video acquisition, production,

distribution, and display processes. Knowledge-driven banding detection methods select relevant models and combine them in a systematic way, so as to produce a prediction of the perceived banding effect.

A knowledge-driven method⁸ is illustrated in **Fig. 2**, where the input is an image or a video frame at the pixel level, and the outputs are a banding score together with a banding map. The banding score denotes the overall level of perceived banding by considering two important factors: banding spread and banding strength. The spread of banding impairment impacts viewing experience but does not solely represent banding annoyance. The contrast sensitivity of the HVS varies based on multiple signal components and viewing conditions, and thus not every banded edge or false contour is perceived equally. The severity of banding is captured by the banding strength component. In other words, the goal here is to detect abrupt local activities in smooth image regions and then analyze the visibility of such activities from the perspectives of HVS characteristics.

Pixels that correspond to abrupt activities deemed visible as banding are then marked, which collectively constitute a banding map of the image or video frame. The banding map illustrates the presence of banding impairment in an image or a video frame and does not reflect the banding strength. Examples of such banding maps are shown in **Fig. 3** (right), where the banding artifacts are highlighted by the white pixels. It appears that this knowledge-driven approach not only detects banding, but also precisely localizes the banding impairment at pixel precision.

Banding regions may be determined by pixels in the image or video frame that have significant local signal activity, while the signal activity in a majority of its surrounding regions is not significant. Therefore, classifying pixels into significant and nonsignificant categories is the first step in detecting banding impairment. To classify pixels, a significance threshold is determined based on the characteristics of HVS as well as a series of signal and system properties. **Figure 4** depicts the workflow of marking pixels based on their signal activity.

Determining the significance threshold in **Fig. 4** is the key to detecting banding impairment. It requires a deep understanding of the HVS properties such as the contrast sensitivity function (CSF)⁹ and various visual masking effects.¹⁰

Figure 5 shows how the significance threshold is generated, starting with HVS modeling, and followed by adjustments based on important video workflow and display factors. Banding is a local activity in smooth image regions that is visible under certain conditions and viewing environments. Therefore, not only signal properties, but also display devices and viewing conditions affect perceived banding.

Modeling CSF and visual masking of the HVS provides a starting point in determining the significance

threshold, which would need to be tuned further to determine precisely the visibility of banding as shown in **Fig. 5**. The CSFs are typically derived based on psychophysical studies on the visibility of patterns with varying luminance levels and spatial frequency.⁹ Chroma component has a significant impact on the visibility of a signal contrast,¹¹ and may be used to adjust the initial significance threshold. SDR is often interpreted as the ratio of the brightest to the darkest luminance. Different opto-electrical transfer functions (OETFs) and electro-optical transfer functions (EOTFs) are designed to accommodate signals with standard and HDR using certain bit depths, that is, 8 and 10 for SDR videos and 10 and 12 for HDR videos. The bit depth of the content determines signal precision and in

conjunction with transfer functions impacts the visibility of banding impairment and suggests adjustments to the significance threshold. Further adjustments to the significance threshold may also be required based on the maximum and minimum values of contents and the capabilities of display devices in producing an adequate range of luminance that is needed to avoid banding.

When all these content, perceptual, chroma, bit-depth, transfer function, and display factors are properly modeled, a precise prediction of visual banding may be achieved. Recently, technologies following this path have emerged and enjoyed growing adoption. The advantages of such knowledge-driven approaches are

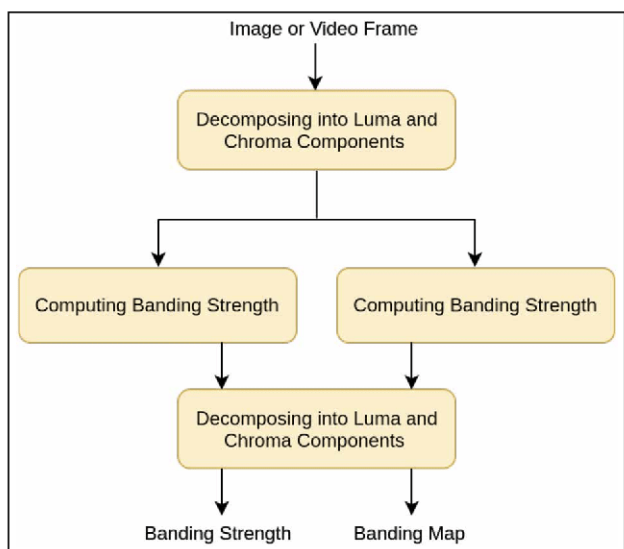


FIGURE 2. Diagram of a knowledge-driven method.



FIGURE 3. Sample images and banding maps created by knowledge-driven method.

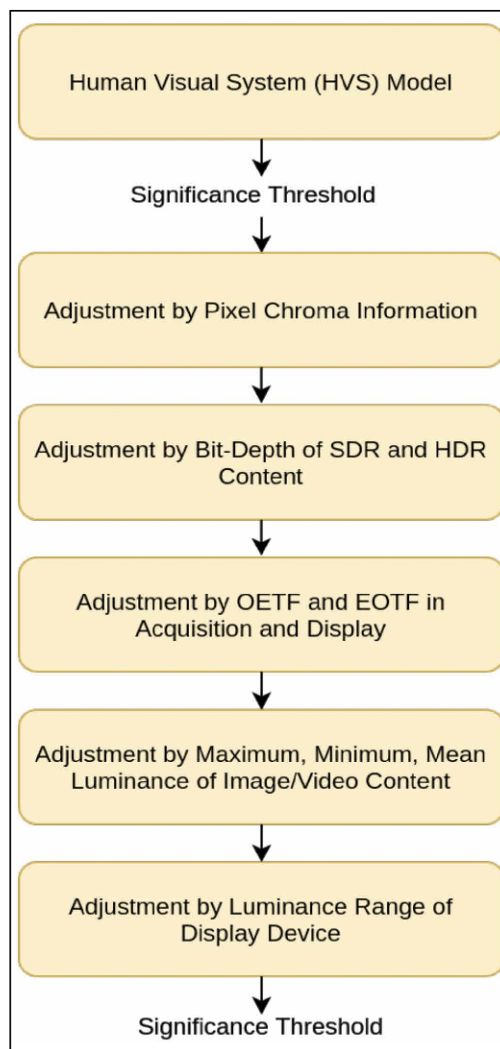


FIGURE 5. Computation and adjustment of significance threshold factors.

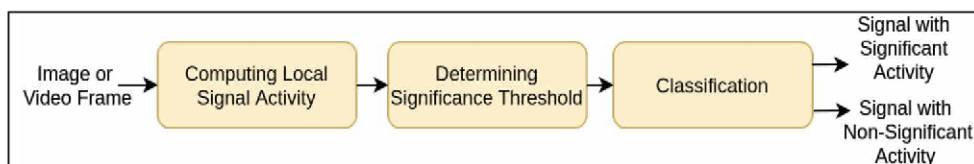


FIGURE 4. Classification of local signals based on an activity threshold.

not only the high pixel-precision accuracy (that allows for the creation of banding maps) and low computational complexity, but also the high explainability—meaning that when banding happens, a deep investigation is plausible to find the cause of banding and then localize the problem to be fixed. The disadvantage of this approach is mainly in the difficulty of the modeling process itself, as precise models of the contributing factors are difficult to develop and parameters of such models are hard to calibrate.

Data-driven/Machine Learning Method

Machine learning, and especially deep learning approaches, have attracted a great deal of attention in recent years and have achieved remarkable success in many application areas. These approaches are generally data-driven, with black-box models being trained by data samples. When these data samples are sufficiently representative of the real-world data distribution, the trained model may be strong enough to make good predictions on novel data samples unseen in the training dataset. The data-driven approach becomes a desirable option in the case of visual banding detection because it avoids the difficulty in developing and calibrating knowledge-driven models.

The first step in building a deep learning-based method is to obtain “big data,” that is, to construct a large-scale dataset for training, validation, and testing. Fortunately, such datasets have emerged recently. A dataset has been constructed,¹² which is composed of nearly 17,000 image patches, together with their ground-truth labels, that is, each image patch has been labeled to be either containing or not containing banding. This allows us to train a DNN, more precisely a convolutional neural network (CNN), to classify a given image patch as either with or without banding in an end-to-end manner. Such a method is end-to-end because the DNN takes a raw image patch of pixels as input and directly produces a classification result as output. As such, the feature extraction process and the classifier are trained or optimized altogether (as opposed to being constructed separately in traditional image classification methods) by learning from data samples. By doing so, a model is constructed with no

explicit domain knowledge. In other words, all knowledge is learned from data and stored in the weights of the trained CNN.

Figure 6 shows a diagram of how a DNN-based banding patch classifier may be applied to assess the visual banding of a given image or video frame.¹³ Image patches are first extracted from the test image using a sliding window that moves pixel by pixel (or a larger stripe when the computational cost is a concern) across the image space. The extracted patches are then fed into the DNN-based banding patch classifier.

The DNN is typically implemented using a CNN structure, which contains multiple convolutional layers, followed by a fully connected neural network. In each convolutional layer, there are multiple linear kernels that convolve with the input signal, followed by activation and pooling processes. The outcome of the layer is considered intermediate features that are subsequently taken as the input to the next layer. Due to the pooling process, the number of features reduces over the layers. At the starting point of the fully connected layer, the features are aligned into a vector, which is fed into a neural network of multiple layers of nodes and with the full connection between all nodes of adjacent layers. The final output of the fully connected layers produces the classification results. Once the CNN architecture is constructed, all that remains is to determine the parameters (including kernel weights in the convolutional layers and the connection weights in the fully connected layers) at each layer. This is achieved by a training process, in which ground-truth outputs are compared with the CNN output and the errors are back-propagated and used to adjust the weights in the CNN. When we have sufficient data samples for training, the CNN will converge to a stage that may make accurate classification results. The classification results obtained for individual local image patches through the CNN classifier are laid out spatially. They are then aggregated with predefined smoothness constraints into two outcomes, as shown in **Fig. 6**. The first is a scalar frame-level banding score for the whole test image and the second is a pixel-level banding map.

Sample images and their corresponding DNN-based banding scores and banding maps are shown in **Fig. 7**. **Figure 7(a)**, **(c)**, and **(e)** have ascending levels of visual banding, which are well reflected by the banding scores.

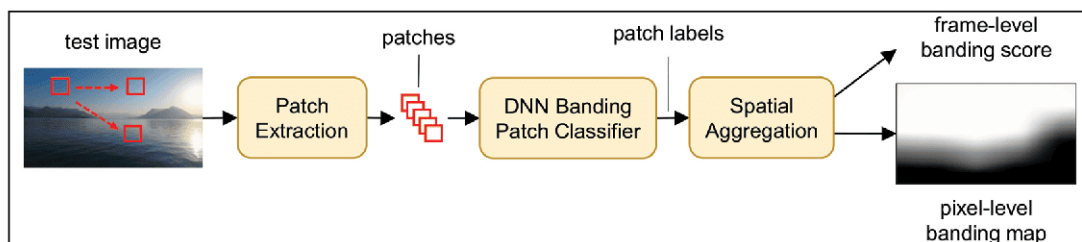


FIGURE 6. Diagram of DNN-based visual banding assessment.

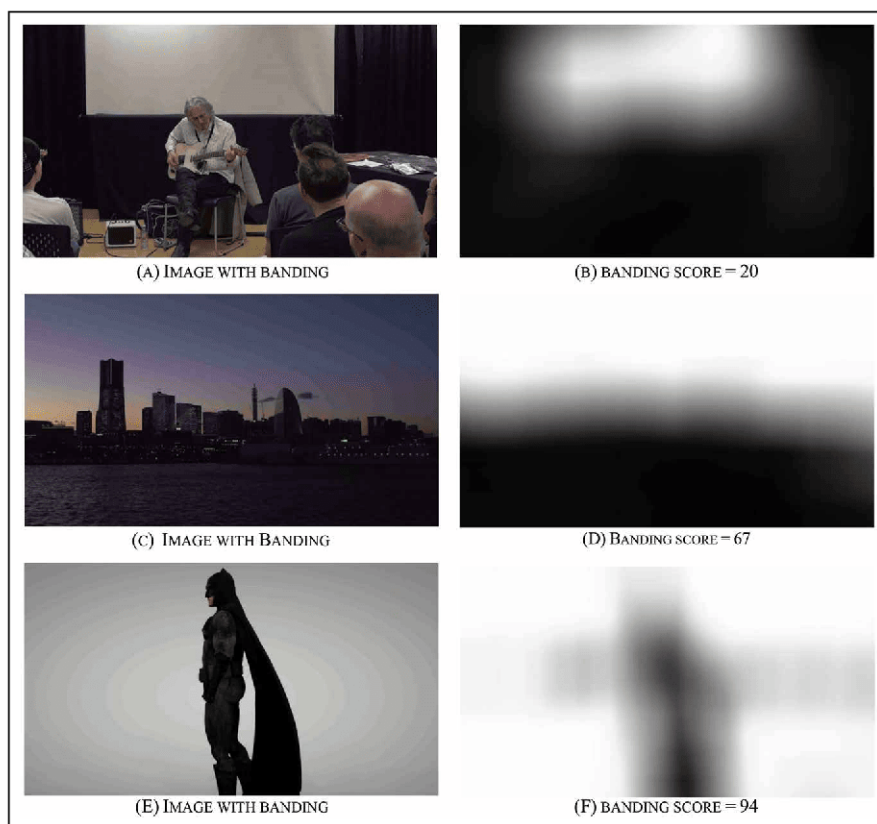


FIGURE 7. Sample images and banding maps created by DNN method. (a) Image with banding. (b) Banding score = 20. (c) Image with banding. (d) Banding score = 67. (e) Image with banding. (f) Banding score = 94.

The banding maps created by the smoothed local patch level banding assessment are given in **Fig. 7(b)**, **(d)**, and **(f)**. It is interesting to compare the banding maps created by the knowledge-driven approach and the DNN-based data-driven approach, as exemplified by **Figs. 3** and **7**, respectively. Both maps offer good predictions on the existence and the spatial locations of banding. The spatial information is very important, especially for subsequent methods that may be used to fix or reduce the banding problem. Comparatively, the maps created by the knowledge-driven method give much more precise localization of banding at the pixel level.

The advantage of the data-driven approach is mainly in the alleviation of the necessity to fully comprehend the domain knowledge, which is sophisticated and evolving over time. On the other hand, its disadvantages are manifold. First, it relies heavily on the quality and quantity of training data. The training sample images need to cover all test cases and the labels assigned to these images need to be very accurate. In general, the behavior of CNN models is difficult to predict on novel data samples, especially when the data samples contain novel content, distortion type, or distortion combinations. This becomes a major issue when the setup of the video acquisition, production, distribution, and display workflow changes, in which case we often need to collect new data and retrain the DNN

model. Second, the quality prediction results may not be as precise as knowledge-driven models, as demonstrated by the comparison of banding maps between **Figs. 3** and **7**. Third, the deployment of DNN models in practical systems may demand high computational and memory resources, especially when the use of the model requires precise localization and diagnosis, such as generating the banding maps as shown in **Fig. 7**, in which the DNN patch classifier needs to be applied to all sliding windows. Nevertheless, significant progress has been made in the past decade on accelerating DNN performance through advanced hardware and software design. Although the training process is usually time-consuming, once the model is trained, the application of the trained model in real-world testing may be made very fast, especially when localized assessment such as the production of banding maps is not required.

Challenges and Future Works

Despite the exciting progress made in the past decades on banding detection and reduction, there are still significant gaps in practice. Some of the root causes of these gaps are summarized as follows.

- First, there is no definite way to differentiate the real contours in the video content and the false contours of banding. Contours of various types are exhibited in real-world videos: some are from camera acquisi-

tion of the visual world, and some are artificial, for example, in animation and computer screen content. How to reliably and efficiently differentiate them remains a challenging problem.

- Second, while it is proved that dithering does reduce banding, there is a dilemma between banding reduction and preserving fine details in an image, because dithering involves adding noise to an image, and noise reduces the visibility of fine texture details. This dilemma becomes even stronger with UHD/HDR/WCG content, which is supposed to bring in finer details than SDR content of the lower resolution and smaller color gamut, but the details embedded in deeper depth bit-planes are even more sensitive to noise contamination. Furthermore, the noise in the dithered content makes video encoding more difficult, as the noisy pixels consume a large number of bits to encode, leaving much fewer bits for true fine details of the original content.
- Third, there has been a debate about the objective of video distribution—whether we should aim for the preservation of the creative intent of the content producers or for creating appealing visual results for the end viewers. It is important to be aware that these two goals may not always align. If preserving the creative intent is the goal, then techniques such as dithering or preprocessing-based banding reduction are problematic, because they purposely change the content.

All of these are intertwined with the never-ending progress of camera and display technologies, and the recent development of scene-adaptive and device-adaptive HDR/WCG processing in the new HDR standards such as HDR10+ and Dolby vision. Consequently, banding detection and reduction will remain an open problem in the future and will evolve with the technology front of the digital video industry.

Conclusion

In this article, we focus on the banding effect, an annoying visual artifact that appears in all stages of the life cycle of digital videos and that has been drawing an increasing amount of attention with the recent growing popularity of UHD/HDR/WCG video content. We discuss the technical details of two promising but substantially different types of approaches for banding detection—knowledge-driven approaches that are built upon deep understandings of the HVS and each component in the video acquisition, production, distribution, and display workflows, and data-driven approaches that learn to detect banding by training DNN models with big data of labeled image samples. Our experiments and analysis demonstrate promising results and predictions using both the mentioned frameworks.

References

1. Z. Wang and A. C. Bovik, “Mean Squared Error: Love It or Leave It? – A New Look at Signal Fidelity Measures,” *IEEE Signal Process. Mag.*, 26(1): 98–117, Jan. 2009.
2. Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, “Image Quality Assessment: From Error Visibility to Structural Similarity,” *IEEE Trans. Image Process.*, 13(4):600–612, Apr. 2004.
3. Z. Wang, E. P. Simoncelli, and A. C. Bovik, “Multi-Scale Structural Similarity for Image Quality Assessment.” *IEEE Asilomar Con. Signals, Syst. Comput.*, Nov. 2003.
4. A. Rehman, K. Zeng, and Z. Wang, “Display Device-Adapted Video Quality-of-Experience Assessment,” *IS&T/SPIE Electronic Imaging: Human Vision & Electronic Imaging*, Feb. 2015.
5. Z. Li, A. Aaron, I. Katsavounidis, A. Moorthy, and M. Manohara, “Toward a Practical Perceptual Video Quality Metric,” *Netflix TechBlog*, Jun. 2016.
6. Y. Wang, S. Kum, C. Chen, and A. Kokaram, “A Perceptual Visibility Metric for Banding Artifacts.” *Proc. 2016 IEEE Int. Conf. Image Process.*, pp. 2067–2071, Sep. 2016.
7. Z. Tu, J. Lin, Y. Wang, B. Adsumilli, and A. C. Bovik, “BBAND Index: A No-Reference Banding Artifact Predictor,” *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, pp. 2712–2716, May 2020.
8. Z. Wang, H. Yeganeh, A. Badr, and K. Zeng, “Image and Video Banding Assessment,” *U.S. Patent Application 17/225, 808*, Oct. 2021.
9. P. G. J. Barten, “Formula for the Contrast Sensitivity of the Human Eye” *Proc. SPIE-IS&T*, 5294:231–238, Jan. 2004.
10. Z. Wang and A. C. Bovik, *Modern Image Quality Assessment*, Morgan & Claypool Publishers, Mar. 2006.
11. G. Denes, G. Ash, H. Fang, R. Mantiuk, “A Visual Model for Predicting Chromatic Banding Artifacts,” *Human Vision and Electronic Imaging (HVEI)*, Jan. 2019.
12. A. Kapoor, J. Sapra, and Z. Wang, “HD Images Dataset With Banded and Nonbanded Region Information,” [Online]. Available: <https://zenodo.org/record/4513740>, 2021.
13. A. Kapoor, J. Sapra, and Z. Wang, “Capturing Banding in Images: Database Construction and Objective Assessment,” *Proc. IEEE Int. Conf. Acoust., Speech Signal Process.*, Jun. 2021.

About the Authors



Hojatollah Yeganeh received a PhD from the University of Waterloo, Waterloo, ON, Canada, in 2014. He is the principal video architect with SSIMWAVE Inc., Waterloo, where he leads the Research and Development (R&D) Team. He has designed and implemented many algo-

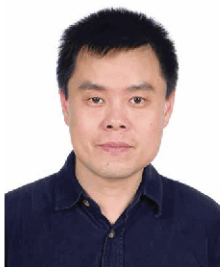
rithms in the field of image and video quality assessment and has published number of journals and conference papers and holds several U.S. patents. He is a Senior Member of IEEE and is one of the recipients of the 2020 Technology and Engineering Emmy Award, which is given by the National Academy of Television Arts and Sciences (NATAS) for outstanding achievement in the development of perceptual metrics for video encoding optimization.



Kai Zeng received a PhD in electrical and computer engineering from the University of Waterloo, Waterloo, ON, Canada, in 2013. He is the Co-Founder and a Lead Researcher at SSIMWAVE Inc., Waterloo, with research interests including computational image and video processing, computer vision,

multimedia communications, with emphasis on image and video quality assessment. He was one of the recipients of the 2020 Technology and Engineering Emmy Award which is given by the National Academy of Television Arts and Sciences (NATAS) for outstanding achievement in the development of perceptual metrics for video encoding optimization.

video quality assessment and optimization, machine learning, computational vision, and multimedia coding research. He is a co-founder and chief science officer of SSIMWAVE Inc., Waterloo, directing video quality-of-experience measurement and optimization solutions for industry. He received an Engineering Emmy Award for contributions to video quality assessment theories and technologies. He is a Fellow of the IEEE, the Royal Society of Canada, Ottawa, ON, Canada, and the Canadian Academy of Engineering, Ottawa. He has more than 200 publications with over 70,000 citations. His awards include IEEE Signal Processing Society Best Paper, IEEE Signal Processing Magazine Best Paper, IEEE-SPS Sustained Impact Paper, University of Waterloo Faculty of Engineering Research Excellence Award, and Natural Sciences and Engineering Research Council Steacie Memorial Fellowship Award.



Zhou Wang received a PhD from the University of Texas at Austin, Austin, TX, USA. He is a professor at the University of Waterloo, Waterloo, ON, Canada, holds the position of Canada Research Chair in Multimedia QoE, and directs the Image and Video Computing Laboratory focused on image/



Keep Up with Emerging Tech

SMPTE Webcasts are a great way to learn about the latest developments

Find new ways to advance your work and your career. SMPTE's live, interactive webcasts help you stay on top of what's happening in the digital media industry, from the newest technologies and standards, to critical business and management strategies.

Our monthly **Technology Series** features industry-recognized experts presenting on the latest tech challenges. In our **Thought Leadership** webcasts, global industry leaders illuminate special topics and innovative solutions in their own specialties. And our **Powered by SMPTE** webcasts offer SMPTE's deep educational expertise to create custom programs that expand your company's resources.



Explore our upcoming offerings, and stream previous webcasts on-demand at [smpte.org/webcasts](https://www.smpte.org/webcasts)