

# Supplementary Materials: Photorealistic Style Transfer via Wavelet Transforms

Jaejun Yoo\*    Youngjung Uh\*    Sanghyuk Chun\*    Byeongkyu Kang    Jung-Woo Ha  
 Clova AI Research, NAVER Corp.

{jaejun.yoo, youngjung.uh, sanghyuk.c, bk.kang, jungwoo.ha}@navercorp.com

## 1. Frame-based signal reconstruction

Our proposed model WCT<sup>2</sup> is inspired by the recent theoretical advancement of frame-based signal reconstruction approaches [7, 8]. To make the paper self-contained, we provide a brief introduction to the frame theory (Section 1.1), tightness of Haar wavelets (Section 1.2) and our theoretical motivation (Section 1.3).

### 1.1. Perfect reconstruction condition

Consider an *analysis operator*  $\Phi = [\phi_1 \ \cdots \ \phi_m] \in \mathbb{R}^{n \times m}$ , where  $\{\phi_k\}_{k=1}^m$  is a family of functions in a Hilbert space  $H$ . Then,  $\{\phi_k\}_{k=1}^m$  is called a *frame* if it satisfies the following inequality:

$$\alpha \|f\|^2 \leq \|\Phi^\top f\|^2 \leq \beta \|f\|^2, \quad \forall f \in H, \quad (1)$$

where  $f \in \mathbb{R}^n$  is an input signal and  $\alpha, \beta > 0$  are called the frame bounds.

The original signal  $f$  can be exactly recovered from the frame coefficient  $z = \Phi f$  when there is the *dual frame*  $\tilde{\Phi}$  (*i.e.*, *synthesis operator*) satisfying the *perfect reconstruction (PR) condition*:  $\tilde{\Phi}\Phi^\top = I$ , since  $f = \tilde{\Phi}z = \tilde{\Phi}\Phi^\top f = f$ . Here, we call such frame *tight* (*i.e.*,  $\alpha = \beta$  in (1)) which is equivalent to  $\tilde{\Phi} = \Phi$  or  $\Phi\Phi^\top = I$ . Note that a tight frame does not amplify the power of the input and thus it has the minimum noise amplification factor. To achieve the best reconstruction performance, frame bases should satisfy another property, called energy compaction. This is particularly important to parametric models, which have to adaptively deal with varying amounts of information with a fixed number of parameters, *e.g.*, deep neural networks (DNNs). For example, singular value decomposition (SVD) provides both tight and energy compact bases given an arbitrary signal. However, SVD is data-dependent, which makes it hard to use for a large dataset.

### 1.2. Wavelet frames

Wavelets are known to compactly represent signals while maintaining important information such as edges, thus resulting in a good energy compaction [8]. Therefore, by using a tight wavelet filter-bank, we can improve the reconstruction performance of encoder-decoder type of networks with minimal noise amplification. Specifically, the non-local basis  $\Phi^T$  is now composed of a filter bank:

$$\Phi = [T_1 \cdots T_L], \quad (2)$$

where  $T_k$  denotes the  $k$ -th subband operator and the filter bank is tight, *i.e.*

$$\Phi\Phi^T = \sum_{k=1}^L T_k T_k^T = I. \quad (3)$$

In this paper, we use Haar wavelets which is one of the simplest tight filter bank frames with low and high sub-band decomposition. Here,  $T_1 \in \mathbb{R}^{\frac{n}{2} \times n}$  is the low-pass subband. This is equivalent to the average pooling:

$$T_1^\top = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & 1 & \cdots & 0 & 0 \\ \vdots & & & & & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & 1 & 1 \end{bmatrix}. \quad (4)$$

---

\* indicates equal contribution

Then,  $T_2$  is the high pass filtering given by

$$T_2^\top = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & -1 & \cdots & 0 & 0 \\ \vdots & & & & & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & 1 & -1 \end{bmatrix} \quad (5)$$

and we can easily see that

$$T_1 T_1^\top + T_2 T_2^\top = I, \quad (6)$$

so the Haar wavelet frame is tight.

### 1.3. Theoretical motivation

In the perspective of the frame-based signal reconstruction, the commonly used encoder-decoder convolution structure of deep neural networks (DNNs), such as U-net [6], can be interpreted as the data-driven way of learning the local bases  $\Psi$  (*e.g.*, convolution filters) with hand-crafted global bases  $\Phi$  (*e.g.*, max-pooling) [7]. Recently, Ye *et al.* [7] interpreted training DNNs as finding a multi-layer realization of the convolution framelets [8]:

$$Z = \Phi^T (f \circledast \Psi) \quad (7)$$

$$f = (\tilde{\Phi} Z) \circledast \tilde{\Psi}, \quad (8)$$

where  $\Phi = [\phi_1, \dots, \phi_n]$  and  $\tilde{\Phi} = [\tilde{\phi}_1, \dots, \tilde{\phi}_n] \in \mathbb{R}^{n \times n}$  (resp.  $\Psi = [\psi_1, \dots, \psi_q]$  and  $\tilde{\Psi} = [\tilde{\psi}_1, \dots, \tilde{\psi}_q] \in \mathbb{R}^{d \times q}$ ) are frames and their duals. Here,  $\circledast$  stands for the convolution operation.

Therefore, the convolutional layers of the encoder learns the signal representation with a global pooling operation. We refer to  $\Phi$  as global bases because it observes the entire image dimension  $n$  while  $\Psi$  learns local features from the data by  $d \times d$  convolution kernels of  $q$  channels. When these frames satisfy the PR condition:

$$\tilde{\Phi} \Phi^\top = I_{n \times n}, \quad \Psi \tilde{\Psi}^\top = I_{d \times d}, \quad (9)$$

the input signal  $f$  can be exactly recovered from the learned representations. Note that the encoder-decoder architectures of WCT [3] and PhotoWCT [4] cannot satisfy the perfect reconstruction condition because of the max-pooling, which does not have its exact inverse (*i.e.*, not a frame). On the other hand, our model WCT<sup>2</sup> can fully exploit the information from the encoder due to the favorable property of the wavelet decomposition and reconstruction, *i.e.*, Haar wavelet pooling and unpooling.

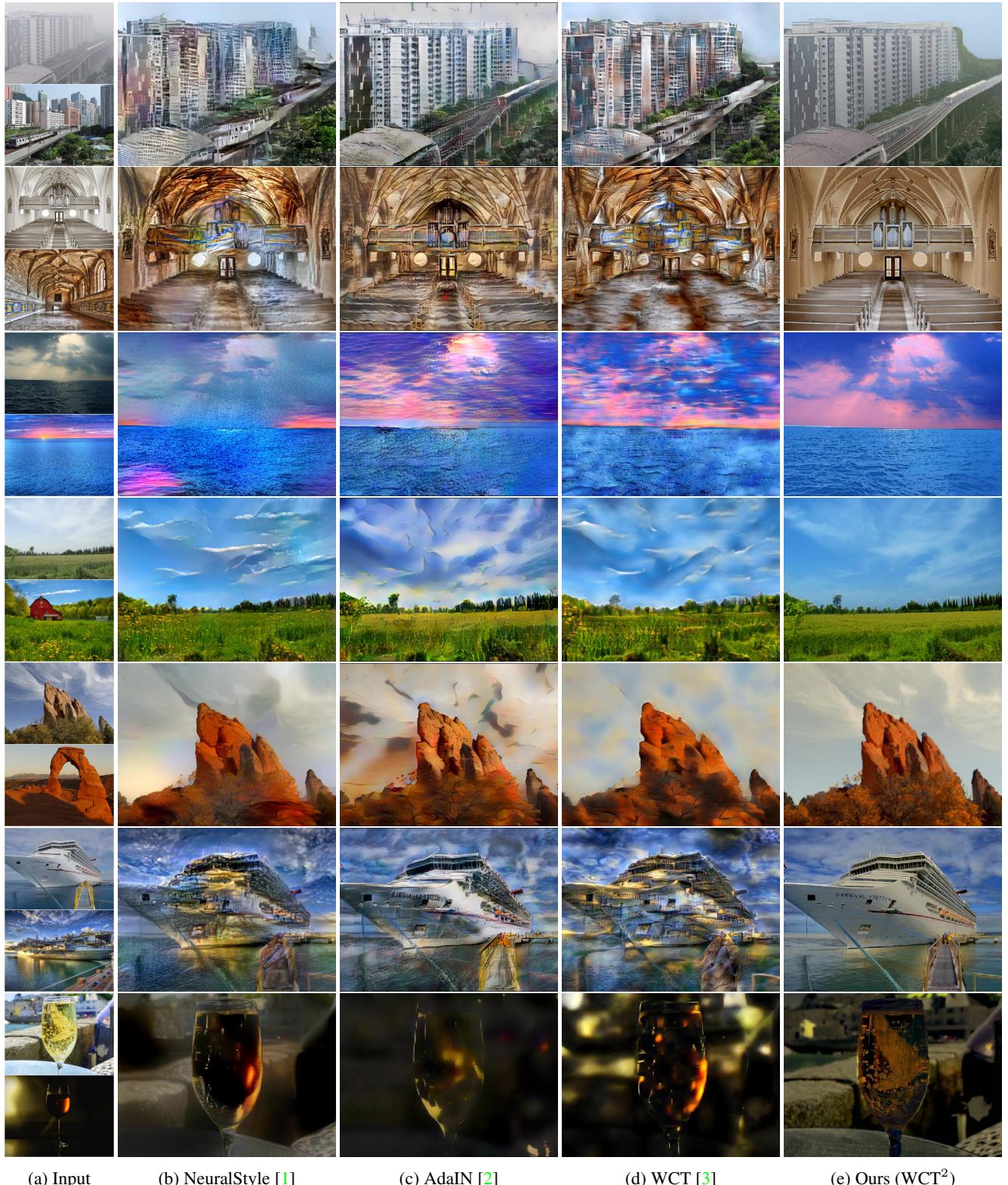
## 2. Results

### 2.1. Qualitative comparison with artistic style transfer results

We compare our proposed WCT<sup>2</sup> with popular artistic style transfer methods including NeuralStyle [1], AdaIN [2] and WCT [3] in Figure 1. To apply semantic segmentation map to the artistic style transfer methods, we followed the spatial control techniques proposed by the authors [5, 2, 3] respectively. In the figure, artistic style transfer methods generate undesired distortions and artifacts and often fail to maintain the structural information despite the spatial control with segmentation maps. In comparison, because of the proposed wavelet corrected transfer, our proposed WCT<sup>2</sup> prevents unrealistic artifacts and preserve the structure information such as edges.

### 2.2. Additional Qualitative comparison with photorealistic style transfer

Additional qualitative results using WCT<sup>2</sup> and its variants are shown in Figure 2, Figure 3 and Figure 4. The video stylization results can be found in one of the other supplementary materials.



(a) Input      (b) NeuralStyle [1]      (c) AdaIN [2]      (d) WCT [3]      (e) Ours (WCT<sup>2</sup>)

Figure 1: Qualitative comparison with artistic style transfer results. Given (a) an input pair (top: content, bottom: style), we compare the results of (b) NeuralStyle [1], (c) AdaIN [2] (d) WCT [3] and (e) ours (WCT<sup>2</sup>).

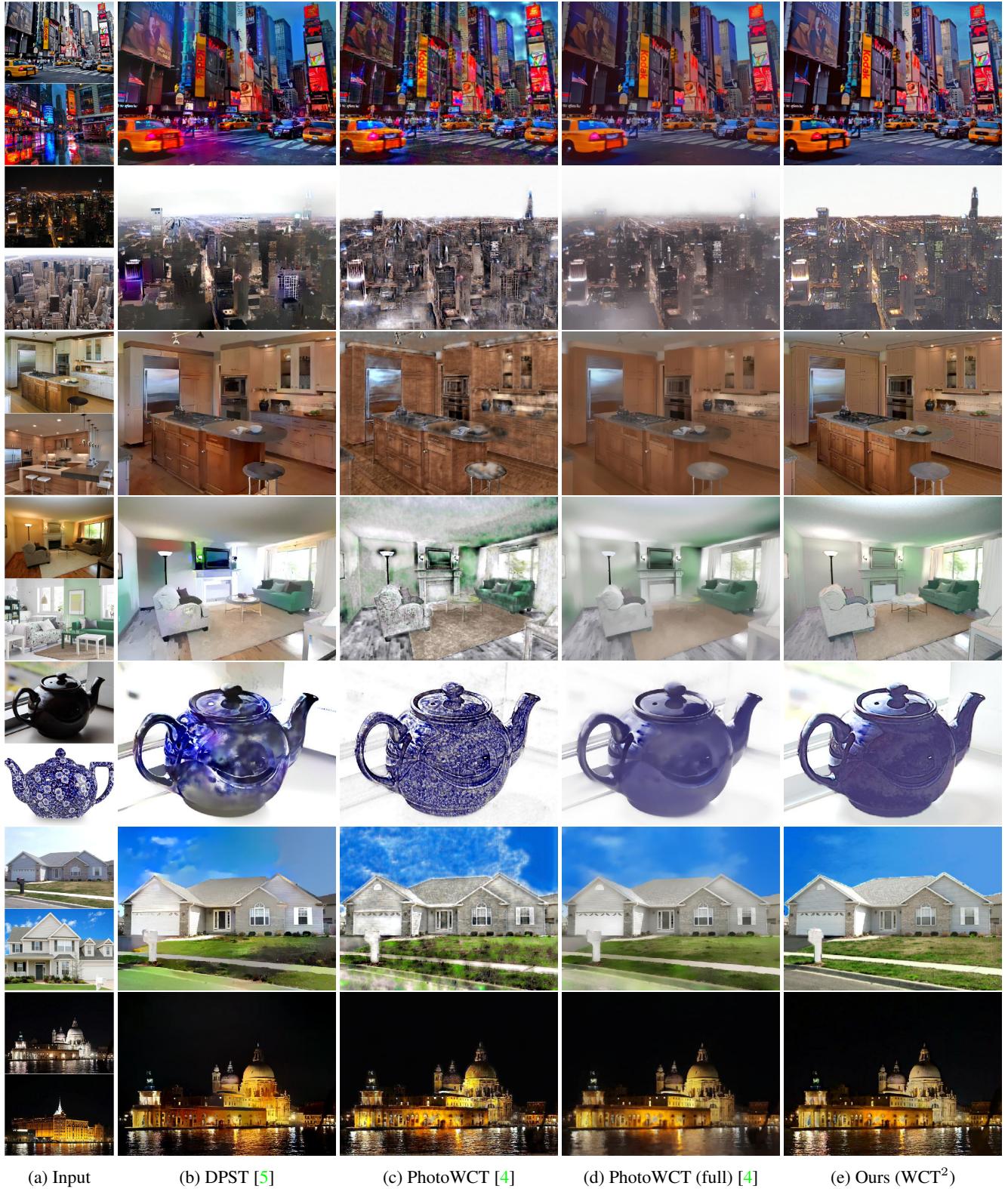


Figure 2: Photorealistic stylization results. Given (a) an input pair (top: content, bottom: style), the results of (b) deep photo style transfer (DPST) [5], (c) and (d) PhotoWCT [4], and (e) ours (WCT<sup>2</sup>) are shown. PhotoWCT (full) denotes the results after applying two post-processing steps proposed by the authors [4]. Note that WCT<sup>2</sup> **does not** need any post-processing.

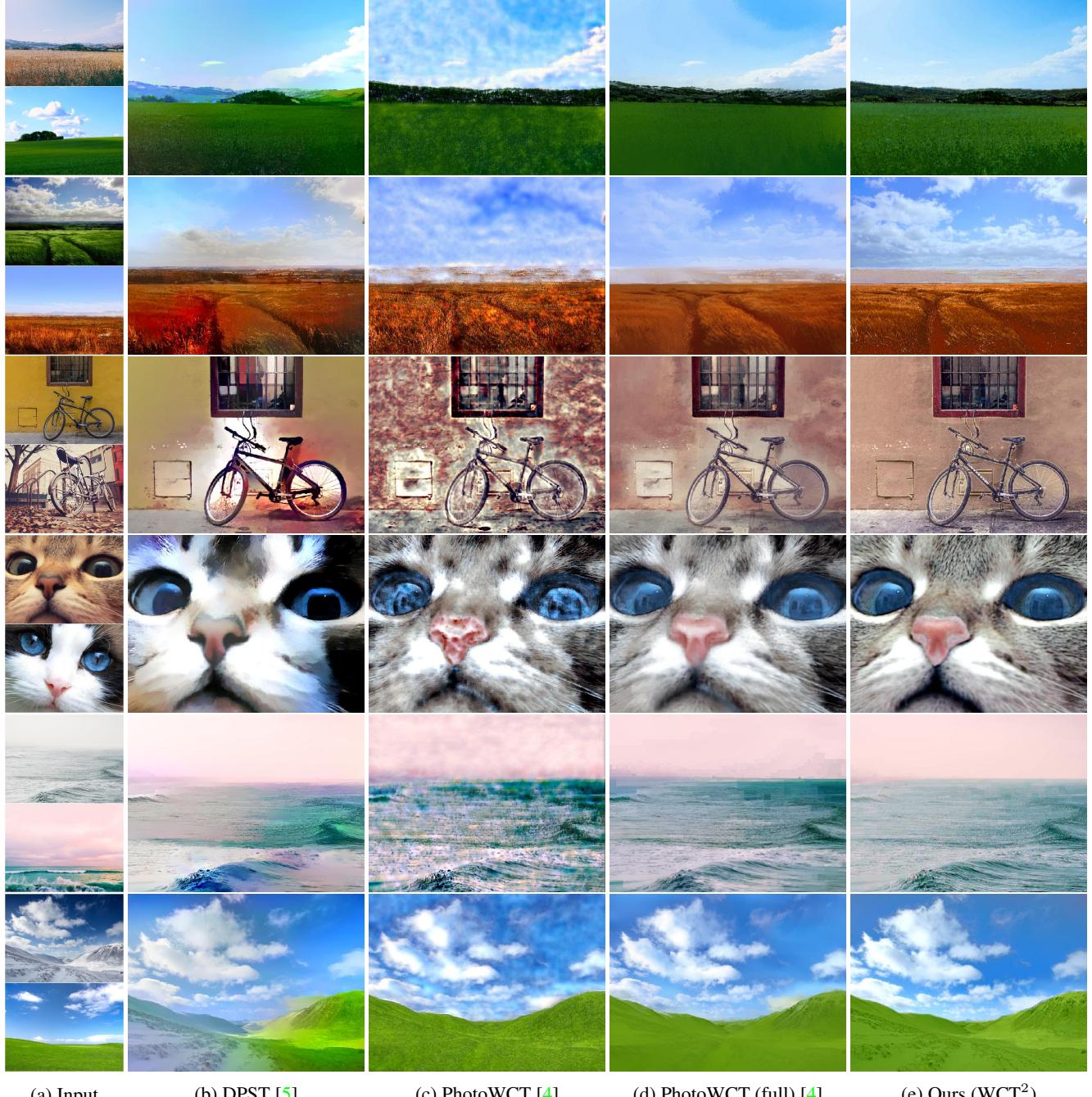
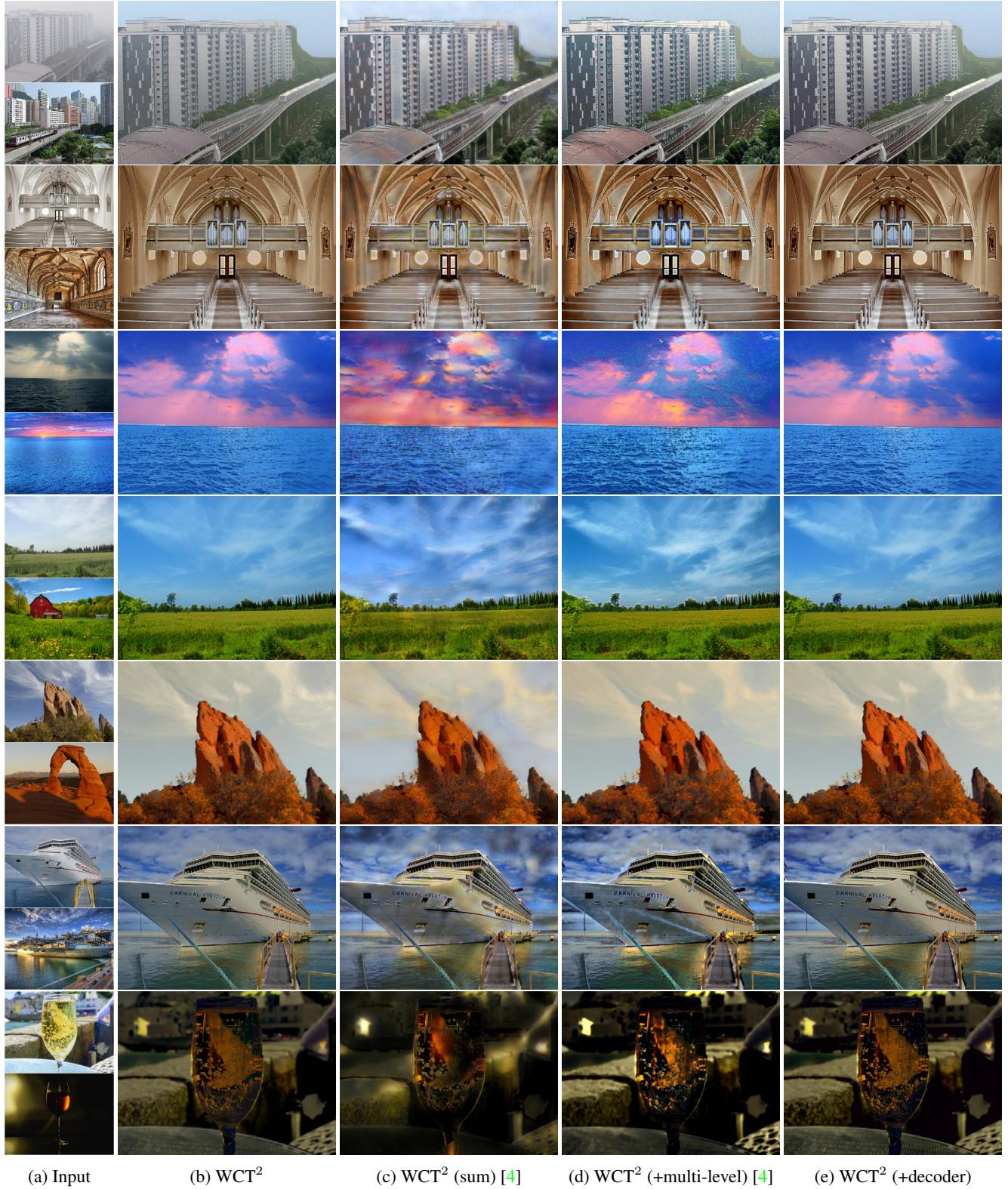


Figure 3: Photorealistic stylization results. Given (a) an input pair (top: content, bottom: style), the results of (b) deep photo style transfer (DPST) [5], (c) and (d) PhotoWCT [4], and (e) ours ( $\text{WCT}^2$ ) are shown. PhotoWCT (full) denotes the results after applying two post-processing steps proposed by the authors [4]. Note that  $\text{WCT}^2$  **does not** need any post-processing.



(a) Input                    (b)  $\text{WCT}^2$                     (c)  $\text{WCT}^2$  (sum) [4]                    (d)  $\text{WCT}^2$  (+multi-level) [4]                    (e)  $\text{WCT}^2$  (+decoder)

Figure 4: Photorealistic style transfer results. Given (a) an input pair (top: content, bottom: style), we compare the results of  $\text{WCT}^2$  and its variants, *i.e.*, (b)  $\text{WCT}^2$ , (c)  $\text{WCT}^2$  (sum) (d)  $\text{WCT}^2$  (+multi-level) and (e)  $\text{WCT}^2$  (+decoder).

### 2.3. Proposed network architecture

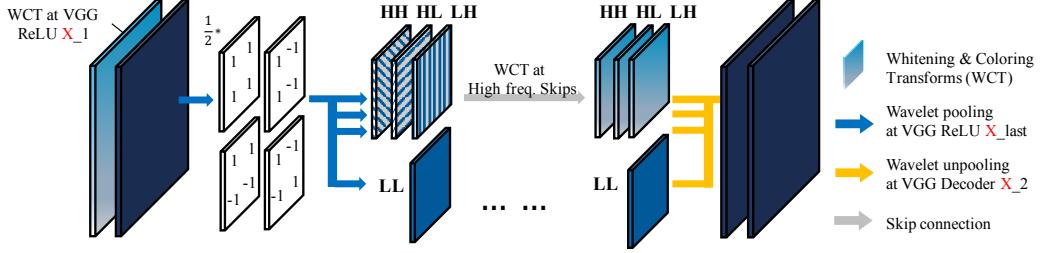


Figure 5: The proposed module using Haar wavelet pooling and unpooling. A pair of encoder and decoder at same scale are shown. WCT is performed on the output of VGG  $\text{convX}_1$  layer followed by subsequent VGG layers and wavelet pooling. Only the low component passes to the next layer and the high frequency components are directly skipped to the corresponding decoding layer. At the decoder, the components are aggregated by the wavelet unpooling.

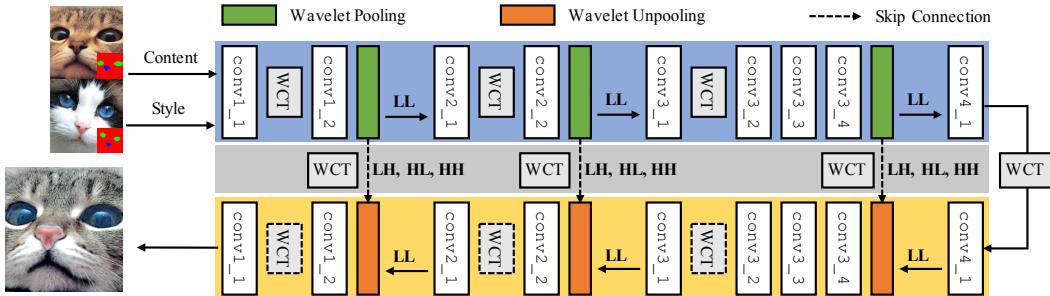


Figure 6: Overview of the proposed progressive stylization. For the encoder, we perform WCT on the output of  $\text{convX}_1$  layer and skip connections. For the decoder, we apply WCT on the output of  $\text{convX}_2$  layer, which is optional.

## References

- [1] L. A. Gatys, A. S. Ecker, and M. Bethge. Image style transfer using convolutional neural networks. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2414–2423. IEEE, 2016. [2](#), [3](#)
- [2] X. Huang and S. J. Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*, pages 1510–1519, 2017. [2](#), [3](#)
- [3] Y. Li, C. Fang, J. Yang, Z. Wang, X. Lu, and M.-H. Yang. Universal style transfer via feature transforms. In *Advances in Neural Information Processing Systems*, pages 386–396, 2017. [2](#), [3](#)
- [4] Y. Li, M.-Y. Liu, X. Li, M.-H. Yang, and J. Kautz. A closed-form solution to photorealistic image stylization. *arXiv preprint arXiv:1802.06474*, 2018. [2](#), [4](#), [5](#), [6](#)
- [5] F. Luan, S. Paris, E. Shechtman, and K. Bala. Deep photo style transfer. [2](#), [4](#), [5](#)
- [6] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. [2](#)
- [7] J. C. Ye, Y. Han, and E. Cha. Deep convolutional framelets: A general deep learning framework for inverse problems. *SIAM Journal on Imaging Sciences*, 11(2):991–1048, 2018. [1](#), [2](#)
- [8] R. Yin, T. Gao, Y. M. Lu, and I. Daubechies. A tale of two bases: Local-nonlocal regularization on image patches with convolution framelets. *SIAM Journal on Imaging Sciences*, 10(2):711–750, 2017. [1](#), [2](#)