

1 Thanks for the careful and valuable comments from reviewers. We are glad to see the agreements that our paper
 2 proposes a new approach (R1), and the idea of error-correction mechanism is intuitive (R1), novel (R2) and smart
 3 (R3). Compared to previous methods, impressive visual improvement is achieved (R1) and our results look better
 4 both quantitatively and qualitatively (R2). Reviewers have some concerns about the motivation of affinity matrices in
 5 non-local blocks (R2), the relation and comparison with iterative back-propagation method (R3) and some other details.
 6 However, we believe all these comments can be addressed through this rebuttal and a minor revision, which will include
 7 some of the explanations below as well as spelling corrections.

8 **Q1. Is any special feature operation applied in ETN? & Does a larger K help? (R1)**

9 A1. Please note that, except employing the learnable matrices to fuse error features, our model is built based on some
 10 conventional operations, like convolutions. And as presented in Line 148 of submission, subtle improvement is
 11 achieved with larger K for images smaller than $2K \times 1K$.

12 **Q2. The motivation to compute affinity matrices & How to achieve the error diffusion. (R2)**

13 A2. We respectfully point out that both stylized image features and corresponding full error features encode information
 14 from the same content-style image pairs and can definitely correlate with each other. If we simply apply full error
 15 features to update stylized images, similar to Gatys et al. [7], it is easy to encounter local minima and hurt semantic
 16 structures. Thus as mentioned in Line 35 of the submission, our goal is to extract more compatible error features
 17 for refinements. We assume that when the error feature is more correlative to the feature of stylized result, they are
 18 more compatible to each other. Then for a stylized image, an affinity matrix can be used to measure its similarities
 19 to error features at different pixels. Thus similar to [23, 11], to better preserve the long-range dependency between
 20 pixels, we diffuse errors to the whole image through a matrix multiplication between a full error feature and an
 21 affinity matrix. Please see Fig. 5 in submission for example. Without joint analysis strategy, the model stylizes
 22 images with unseen white bumps and blurs the structure of clouds more. And as shown in Fig. 3 in supp file, our
 23 full model equipped with non-local blocks can produce results with clearer outlines (e.g., the edges between the
 24 sky and mountain) and better respect the semantic information while the model without blocks fails to preserve
 25 the orientation of the road in land. Moreover, as shown in the Table 1 in supp file, the non-local blocks can also
 26 improve the perceptual content & style metrics, making the diffusion mechanism effective.

27 **Q3. Why not simply concatenate the error feature of layer i with error feature of layer i-1? (R2)**

28 A3. Feeding ΔE^i into convolution layers, we denote $\Delta \hat{E}^i \in R^{N \times C^{i-1}}$ as the processed feature, $\Delta \hat{E}^i = \Delta E^i \otimes \phi_t$.
 29 N and C^{i-1} indicate the number of pixels and the channel number of $\Delta \hat{E}^i$ respectively. Thus we clarify that we
 30 apply the learnable matrix to merge $\Delta \hat{E}^i$ with $\Delta E_s^{i-1} \in R^{C^{i-1} \times C^{i-1}}$ due to their different dimensions.

31 **Q4. Performance issues, including increased training burden and running time. (R2)**

32 A4. Thanks for pointing out the mistake in real-time stylization, which will be corrected in revision. Our method
 33 indeed slightly increases training burdens, but achieves remarkable improvements with comparable running time.

34 **Q5. Relation and comparison with iterative back-propagation method. (R3)**

35 A5. We clarify that we share the same stylization goal and content & style definition as Gatys et al. [7]. However, one
 36 key difference is that [7] back-propagates the error features while our method achieves the error-correction in a
 37 feed-forward manner. Thus we can explicitly perform joint analysis between errors and a synthesized image for
 38 a better residual image and apply the error diffusion to capture the long-range dependency in pixels. Moreover,
 39 the developed technology can be extended to other fields more easily thanks to its efficiency and effectiveness.
 40 Fig. 1(a) shows that, [7] often gets stuck into local minimum and fails to capture the salient style patterns. In
 41 comparison, our approach, shown in Figs. 6-7 of the submission, yields more visually pleasing results.

42 **Q6. Comparison with ResNet like structure as an alternative to the structure in Figure 2. (R3)**

43 A6. Laplacian pyramid can facilitate performance in various fields [1, 9, 5]. In fact, we have built a baseline model
 44 by keeping the same resolution during stylization and demonstrate the results in Fig. 2 and Table 1 of supp file
 45 and Fig. 1(b). With error-correction mechanism, the baseline model is successful in preserving high-level content
 46 structure and receiving presented style patterns, like brush strokes, color distributions and block-wise patterns. But
 47 the pyramid strategy enables us to better capture large-scale style information and communicate different scales,
 48 thus exhibiting more abundant patterns and superior details. The usage of pyramid also improves the style loss.

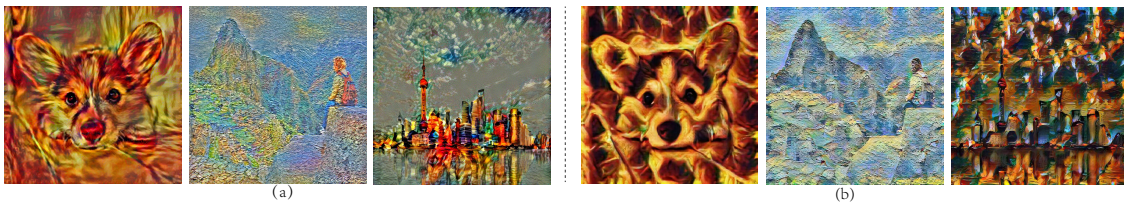


Figure 1: Stylization results of Gatys et al. (a) and ResNet-like baseline (b) for tests shown in Figs. 6-7 in submission.