

1 We would like to thank the reviewers for their positive and constructive comments, and for finding the idea novel, and
2 the results impressive. The major concerns are addressed below. The final paper will be updated accordingly. Also, the
3 typos will be fixed, and additional references suggested by the reviewers will be cited.

4 **Universal texture synthesis claim and non-stationary textures.** We agree that the proposed method does not handle
5 well non-stationary textures, and suits better textures with repetitive patterns. It, however, is universal in a sense that it
6 generalizes to unseen textures. Note, state-of-the-art methods such as [Zhou et al'18, Shaham et al'19] train and test
7 on a single texture and do not generalize to unseen examples. Having said that, looking at Fig. 5 (paper) and Fig. 2
8 (appendix), Neural-FFT handles a good range of regular, near-regular, and even irregular and stochastic textures. In
9 order to clarify this issue, as suggested by R4, we will add a discussion in the final paper, and move appendix's Fig. 5
10 (failure cases with non-stationary textures) to the main paper as suggested by R2.

11 **R1-1. Inference time comparison with self-tuning.** Self-tuning [Kaspar et al'15] is not amenable to parallelization
12 due to its iterative optimization nature and bookkeeping mechanism for the spatial uniformity constraint. Neural-FFT,
13 however, is quite parallelizable taking great advantage of GPUs. Having said that, we ran Neural-FFT in CPU mode
14 (same setting as self-tuning), and the inference time is: 1.29 sec versus 140 sec for self-tuning showing 108x faster
15 synthesis. Note, thanks to the parallel nature of Neural-FFT, inference on GPU takes 45 msec that is 3,111x faster than
16 self-tuning. We will include this discussion in the final paper. Regarding the synthesis quality, as shown in Fig. 2, and
17 Fig. 2,3 (appendix), self-tuning tends to produce repetitive outputs, and can break the regularities, whereas texture CNN
18 uses ground-truth textures for synthesis.

19 **R1-2. Comparison with [Efros and Freeman 2002].** Texture synthesis has advanced a lot after this pioneering work.
20 As per reviewer's suggestions, we compare with image quilting (using publicly available code¹), and representative
21 examples are shown in Fig. 1. As evident from Fig. 1, image quilting performs poorly in synthesizing large-scale
22 structures and multi-scale texture details. Similar observations have been made by prior works e.g., in [Kaspar et al'15].

23 **R1-3. Evaluation metrics.** There is no gold-
24 standard image quality metric, and most previous
25 work on texture synthesis only provide visual
26 comparison. We however provide a diverse array of
27 eight metrics as each single metric could have its
28 own bias. We agree that SSIM and LPIPS are not
29 the best metrics for synthesis (where there are several
30 good solutions), and thus we provided FID,
31 c-FID, and c-LPIPS as well to compare the distribution of input/ground-truth images with the synthesized images. In
32 order to compute the FID score, we measure the Frchet distance between the Inception-v3 statistics for a set of 200
33 synthesized images (resolution: 256x256) and the corresponding set of original (ground-truth, resolution: 256x256)
34 images based on [Heusel et al'17].

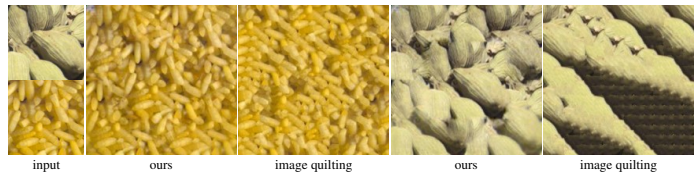


Figure 1: Representative examples for neural FFT versus image quilting.

35 **R1-4. Remark 1 and failure examples.** It seems that there is a misunderstanding. Indeed, remark 1 means that there
36 could be several good solutions for synthesis given an input example, where one can use a knob to control the trade-off
37 between structural similarity and diversity. The examples in Fig. 5 (appendix) are failures since based on the very first
38 definition of texture synthesis, the structural patterns are broken irrespective of the ground truth.

39 **R2-1. Clarifications on the PS number.** Preference score (PS) indicates the portion of workers that prefer our result
40 over the other method, averaged over 200 examples.

41 **R3-1. High resolution texture synthesis.** Regarding
42 the higher amounts of upsampling, a simple
43 heuristic is to take the trained 2x synthesis model,
44 and perform synthesis sequentially a few times
45 (fully convolutional network is invariant to the input
46 image dimension) to reach the desired resolution.
47 For example, for 4x synthesis results in Fig. 6 (main
48 paper) and Fig. 3 (appendix), we first perform 128→256 synthesis, and then 256→512 synthesis.

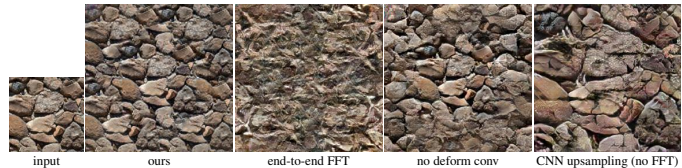


Figure 2: Illustration of ablation study for a representative example.

49 **R4-1. Qualitative results for the ablation study.** Thank you for the suggestion. A representative example is shown in
50 Fig. 1, and more examples will be added to the final version.

51 **R4-2. More discussions about the shortcomings** We discussed the shortcoming of no diverse output and possible
52 solutions in Sec. 5 of the paper. As most non-stationary textures emphasize directional effects, one possible way to
53 handle non-stationary textures could be emphasizing some specific FFT components while suppressing the others.

54 **R4-3. Alternatives to transposed convolution upsampling.** We have tested several other local upsampling methods
55 for FFT upsampling such as nearest neighbor, bilinear and trilinear interpolation. We empirically found that transposed
56 convolution works the best. We will include the results of these alternative upsampling experiments to the ablation
57 study.

¹<https://github.com/rohitrango/Image-Quilting-for-Texture-Synthesis>