

1 We thank all the reviewers for providing constructive feedback. First, we hope to bring a new theoretical result to the
2 reviewers’ attention. We can now prove that **Algorithm 1 is guaranteed to converge locally** whenever $0 < \alpha < 2$.
3 We will include this new result in the revision.

4 **To Reviewer #1**

5 **Q1:** *Difference from emerging convolutions.*

6 We thank R1 for pointing out a missing reference. We developed our approach without knowing the work of emerging
7 convolutions. We will include a discussion of it in Section 4.3. We believe that our approach has *critical advantages*
8 *over emerging convolutions*, which can be detailed as follows:

- 9 • **The inversion is much faster to compute.** Emerging convolutions are inverted using forward/back substi-
10 tutions designed for inverting triangular matrices, which requires the same number of iterations as the input
11 dimension. In stark contrast, we use a fixed-point iteration method (Algorithm 1) for inversion, which is similar
12 to i-ResNet and requires substantially fewer iterations than the input dimension. For example, our inversion
13 takes 120 iterations to converge on CIFAR-10, while emerging convolutions will need 3072 iterations. In other
14 words, our inversion is roughly 25 times faster than emerging convolutions on powerful GPUs.
- 15 • **No bottleneck in the architecture.** The input and output of an emerging convolution must have the same
16 dimensions, and therefore emerging convolutions cannot be used to expand the hidden dimensions of a neural
17 network. As a result, if a flow model is built using only emerging convolutions, the largest number of hidden
18 units for a layer will be the input dimension, which becomes a bottleneck in expressive power. In stark contrast,
19 we can use masked convolutions to expand the hidden dimensions (See Eq. (2)) and eliminate the bottleneck.
20 This is why we can build MintNet *using only masked convolutions*, while the emerging convolution paper
21 modified the Glow architecture by replacing 1×1 convolutions with emerging convolutions.

22 **Q2:** *Strengths over i-ResNet, ResNet, Flow++, etc.*

23 We would like to remind the reviewer that we use Section 4.3 to compare MintNet with many popular paradigms of flow
24 models, *e.g.*, free-form invertible models (including i-ResNet) and dimension partitioning models (including Flow++).
25 We will clarify Section 4.3 more to incorporate the following information.

26 Unlike ResNet, MintNet is an invertible architecture that can be used directly as a flow generative model for density
27 estimation and sample generation. Unlike i-ResNet, MintNet provides exact likelihood values and better empirical
28 performance (see Table 1), while having a comparable cost of sample generation. Unlike Flow++, MintNet is a
29 free-form invertible model that does not rely on dimension partitioning and coupling layers. We did not compare results
30 with Flow++, as it uses variational dequantization, while all models we compare (Glow, i-ResNet, and other models in
31 Table 1) use uniform dequantization. For the same model, the bpd values resulting from variational dequantization are
32 **lower bounds** to the bpd values of uniform dequantization, and therefore will always appear to be “better” (hence not
33 comparable).

34 **To Reviewer #2**

35 **Q1:** *Code and experiments on ImageNet 64×64 .*

36 We will put the code and checkpoints online once the paper is accepted. We agree that results on more datasets are
37 helpful. However, our resources are too limited to finish experiments on ImageNet 64×64 , and we believe that our
38 state-of-the-art results across all three datasets (MNIST, CIFAR-10, and ImageNet 32×32) are sufficient to demonstrate
39 the advantages of MintNet.

40 **To Reviewer #3**

41 **Q1:** *The initial choice of \mathbf{x}_0 .*

42 When inverting Mint layers, we always choose $\mathbf{x}_0 = \mathbf{z} \odot \frac{1}{t}$ because we empirically observe $\mathbf{z} = \mathcal{L}(\mathbf{x}) \approx \mathbf{t} \odot \mathbf{x}$ for
43 deeper layers after training. We will elaborate more on this in Section 4.1.

44 **Q2:** *More explanations on Paired Mint layers.*

45 A Paired Mint layer is a lower triangular Mint layer followed by an upper triangular Mint layer. The results of the lower
46 triangular layer are fed to the upper one as the input. There are clearer descriptions on Paired Mint layers in the network
47 architecture tables in Appendix D (see, *e.g.*, Table 3). We will clarify this more in Section 4.2.