

Table 1: (a) Image synthesis speed on CIFAR10. Glow re-implemented in PyTorch is masked with †. ‡ denotes results shown in Hooeboom et al. [2019]. (b) Image synthesis speed of MaCow on different datasets.

(a)			(b)		
CIFAR10	time/sample (ms)	Slow-down	Dataset	image size	time/sample (ms)
Glow‡	5	1.0	CIFAR10	32 × 32	38.7
MAF‡	3000	600.0	ImageNet-64	64 × 64	104.7
Emerging‡	1800	360.0	LSUN-128	128 × 128	267.9
			CelebA-HQ-256	256 × 256	434.2
Glow†	5.3	1.0			
MaCow	38.7	7.3			

1 We thank for the valuable feedback. We address the questions below and will revise our paper accordingly.

2 [R1 & R2 & R3 Generation speed] We appreciate your suggestions on adding experiments on image generation speed  
 3 comparison and plan to add a section in experiments to discuss it in details. Table 1a and Table 1b provides the  
 4 preliminary results. Following Hooeboom et al. [2019], we measure the time to sample a datapoint when computed in  
 5 mini-batches with size 100. For fair comparison, we re-implemented Glow using PyTorch and show that the speed of  
 6 Glow† and Glow‡ is comparable. On CIFAR-10, MaCow is 7.3 times slower than Glow, much faster than Emerging  
 7 Convolution and MAF, whose factors are 360 and 600 respectively. The generation speed of MaCow on different  
 8 datasets is shown in Table 1b. We see that the time of generation increases linearly with the the image resolution.

9 [R2 Related work and Originality] Thanks for pointing out the related work. We will cite the Emerging Convolution  
 10 paper in our final version and discuss the detailed relation and difference with MaCow. There are two main differences  
 11 between MaCow and Emerging Convolution. i) the pattern of the mask is different. With the mask in MaCow (Figure  
 12 1 in the paper), MaCow achieves significantly more efficient inference and sampling (shown in Table 1a and 1b) by  
 13 reducing the complexity from  $O(h \times w)$  to  $O(h)$  or  $O(w)$ , without sacrificing the receptive field. ii) the Emerging  
 14 Convolutional Flow [Hooeboom et al., 2019] is basically a linear transformation with masked convolutional kernels,  
 15 which does not introduce “nonlinearity” to the random variables. This flow is proposed as an alternative to the  $1 \times 1$   
 16 conv flow in Glow. MaCow, in contrast, is able to introduce such nonlinearity similar to the coupling layer, and is  
 17 proposed to replace both the coupling layers and the  $1 \times 1$  conv layers in Glow. We would like to point out that your  
 18 comment “the masked conv layers in MaCow are linear” is a misconception. Actually, regardless of the different  
 19 patterns of masks, Emerging Conv is a special (linear) case of MaCow, by specifying the  $s(x_{<t})$  and  $b(x_{<t})$  in Eq (4)  
 20 as linear functions.

21 [R2 & R3 Novelty and significance of improvements] First, we would like to argue that the proposed new pattern of  
 22 mask in MaCow is simple and effective, but not trivial or incremental. It produces a semi-autoregressive model, which  
 23 significantly reduces the inference time of MAF and obtains better density estimation performance.

24 Second, the improvements of MaCow on bits/dim, especially the contribution of fine-grained multi-scale architecture,  
 25 are not neglectable. Emerging Convolution [Hooeboom et al., 2019] obtained 0.02 improvement on bits/dim by  
 26 increasing both the number of parameters and inference time. Our fine-grained architecture reduces the number of  
 27 parameters (compared with the original multi-scale architecture) and obtains 0.03 bit/dim improvements on both  
 28 CIFAR-10 and ImageNet-64. The overall improvements of MaCow over Glow, without Variational dequantization, are  
 29 0.07 on CIFAR-10, 0.06 on ImageNet-64, 0.04 on LSUN-bedroom and 0.08 on CelebA-HQ.

30 [R3 Flow++ v.s. MaCow] From results shown in Table 1 of our submission, in terms of density estimation, we could find  
 31 that on CIFAR-10 with uniform dequantization, MaCow (3.28) performs better than Flow++ (3.29) and on ImageNet-64,  
 32 with variational dequantization, MaCow (3.66) outperforms Flow++ (3.69). The only exception is on CIFAR-10 with  
 33 variational dequantization, Flow++ (3.09) achieve better performance than Macow (3.16). But we have to mention that  
 34 even with similar number of parameters, Flow++ is slower and consumes much more memory than Glow and MaCow,  
 35 preventing us from evaluating it on high-resolution images.

36 [R3 Sampling with temperature] The temperature trick is only applied to LSUN and CelebA-HQ 5-bits images, where  
 37 MaCow adopts additive coupling layers. For CIFAR-10 and ImageNet 8-bits images, we sample with temperature 1.0.  
 38 [R1 non-numerical measures of modifications] We appreciate your suggestion on evaluating the difference between Ma-  
 39 Cow w./w.o. these modications in non-numerical ways. From the image samples w./w.o the variational dequantization,  
 40 we have not observed significant difference. We will consider some other non-numerical metrics.

41 [R2 & R3 Error bar] We have performed experiments on CIFAR-10 with multiple runs (more than 3) with different  
 42 random seeds and the standard deviation is less than 0.005. On other datasets, same with Glow and Flow++, we only  
 43 performed single run on each dataset due to the limits of computational resources, unfortunately.

44 [R2 statement of stability] We appreciate your comment about our irrigorous statement of stability in line 128. We will  
 45 amend this claim in the final version.

46 [R1 & R2 & R3 Paper writing] We thank for all your suggestions on revising the paper to improve the writing. We will  
 47 elaborate the Masked Convolution and Fine-grained architecture sections to explain them more clearly.

## 48 References

49 Emiel Hooeboom, Rianne V. Berg, and Max Welling. Emerging convolutions for generative normalizing flows. In *ICML*, 2019.