

1 We thank all reviewers for their careful reading of the manuscript and their constructive comments.

2 **Reviewer-1: The search space of  $n$ .** In Figure 2(b),  $\log(n)$  has to be an integer, which is  $9 \leq \log(n) \leq 14$ .

3 **Reviewer-1: The latency per slot.** Instead of the latency per slot, we reported the end-to-end inference latency of a  
4 privacy-preserving neural network as the key performance metric. First, we would like to emphasize that the numbers  
5 of  $n$  and  $q$  we provided in Figure 3 and 6 are **all in the log scale**. We also would like to point out that the comment  
6 of “The range for  $q$  also shows about a 20% variation” is **inaccurate**. What we have is “The range for  $\log(q)$  also  
7 shows about a 20% variation”. Second, although we agree with the reviewer on the analysis on the latency per slot, we  
8 think the reviewer misunderstood the latency per slot on our baseline using the same  $n$  and  $q$  for all layers. Different  
9 neural network layers have different numbers of input and output channels, and different weight kernel sizes. Therefore,  
10 different layers in a neural network require different values of  $n$  to pack all weights. Our baseline selected a large  $n$  to  
11 have enough slots to pack the layer with the largest number of weights. However, for the other layers, most slots are  
12 empty since they have less weights. In this way, the “effective” latency per slot of our baseline is very long, where the  
13 effective latency means the latency per non-empty slot.

14 **Reviewer-1: The 70% latency reduction.** We would like to emphasize the fact that the numbers of  $n$  and  $q$  we  
15 provided in Figure 3 and 6 are **all in the log scale** again. We measured the latency of each HE layers on a real machine.  
16 By reducing  $n$  and  $q$ , the cache hit rate greatly increases during the NTT and CRT computations. Therefore, we  
17 observed a great latency reduction. We do NOT think simply scaling the latency per slot with  $n$  and  $q$  is a good latency  
18 estimation.

19 **Reviewer-1: Comparison against the fixed aggressive setting.** We compared AutoPrivacy against DARL [4] in Table  
20 2 and Figure 6. DARL aggressively sets the same  $n$  and  $q$  for all layers of a neural network.

21 **Reviewer-1: The models were used.** We explained the models we studied in Section 4. We studied a 7-layer CNN  
22 network used by [5] (7CNET), ResNet32 [24] (RESNET), and MobileNet-V2 [25] (MOBNET). We quantized all  
23 models with 8-bit. Due to the limited space, we cannot include the details of the network architecture in the manuscript.  
24 We will try to add the information in the next version of this manuscript.

25 **Reviewer-1: Figure 3 does not represent the runtime per slot.** Figure 3 shows the execution time of a HE multipli-  
26 cation we measured on a real machine with different values of  $n$  and  $q$ . Again, we believe the reviewer underestimated  
27 the “effective” latency per slot of our baseline.

28 **Reviewer-2: Analysis on the decryption of multiplied ciphertexts.** We believe it is difficult to do a mathematical  
29 analysis on the error rate of a neural network with different values of  $n$  and  $q$ . A mathematical error derivation is too  
30 complicated for an inference of a specific privacy-preserving neural network. This is why we propose AutoPrivacy  
31 that selects a set of  $n$  and  $q$ , and feeds them into the HE protocol and real HE-enable neural network to calculate the  
32 accuracy. The error tolerance of a privacy-preserving neural network is architecture- and application-dependent.

33 **Reviewer-2: Figure 1(d) Why  $n$  does not influence the accuracy.** The decryption error is not related to  $n$ .  $n$  decides  
34 the number of slots that can be packed in a cyphertext. Any large convolution can be broken into smaller pieces, so that  
35 we can always use a smaller  $n$  to perform the computation but with longer latency.

36 **Reviewer-3: Comparison against the other search techniques.** In this paper, we present a new and important prob-  
37 lem on the parameter selection of privacy-preserving neural networks. Presenting the problem is our first contribution.  
38 Our design is the first work to identify the fact that the mathematical error derivation is not necessary for the inferences of  
39 a privacy-preserving neural network. We can achieve better inference latency than two of the most recent state-of-the-art  
40 designs. We will compare AutoPrivacy against other search techniques in the next version of this manuscript.

41 **Reviewer-3: 19.55 seconds of inference latency on the Cifar dataset isn’t practically useful.** The inference time is  
42 architecture-dependent, but not dataset-dependent. By the architecture of 7CNNet, our fastest inference on CIFAR-10  
43 requires only 6.92 seconds, which is much faster than all existing state-of-the-art designs.