

## 1 A Appendix

### 2 A.1 Toy problem

3 In order to further understand how running statistics and current statistics with Gaussian noise as  
4 input compare against real data, we consider the following toy problem of binary classification of 2D  
5 points distributed as concentric circles (real data). We take a simple Multilayer perceptron (MLP)  
6 with 2 hidden layers and 2 BatchNorm layers. This is done to restrict the input and embedding space  
7 (output from penultimate layer) to 2 dimensions each which can be directly visualized, rather than  
8 looking at their distributions. Figure 1 shows the real data on which the MLP is trained, the input  
9 Gaussian noise samples, and the embeddings in each of the three cases, respectively from left to right.  
10 We see that, although the distribution of Gaussian noise is considerably different from real data, the  
11 embeddings are close to the embeddings of real data while using current statistics as compared to  
12 using running statistics.

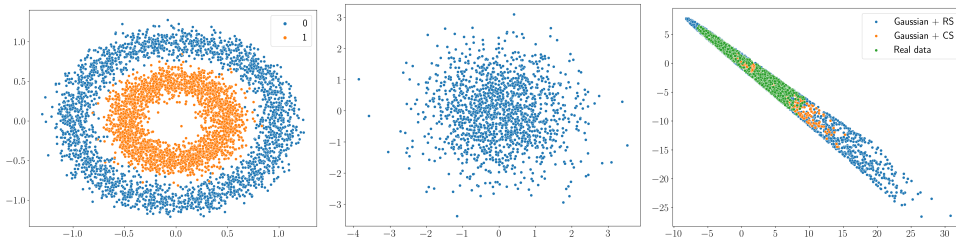


Figure 1: (Left) Circles data on which the MLP is trained. (Middle) Gaussian noise used as input to the trained MLP. (Right) Scatter plot for embeddings in different cases.

### 13 A.2 Handling different number of BatchNorm layers

14 Here we consider CIFAR10 dataset on which teacher is trained and the teacher-student pair of  
15 ResNet34-ResNet18. While performing distillation using Gaussian noise, we randomly choose a  
16 given percentage of BatchNorm layers ( $\mathcal{P}$ ) from the teacher network and restrict those layers to  
17 use running statistics, whereas the remaining BatchNorm layers use current statistics. We vary the  
18 percentage of such layers and report the result in Table 1. We observe that as the percentage of  
19 number of BatchNorm layers using running statistics increases, the student accuracy decreases. Note  
20 that, the number of BatchNorm layers in the teacher network remains same.

Table 1: Varying percentage of BatchNorm layers ( $\mathcal{P}$ ) that uses running statistics. As the percentage of number of BatchNorm layers using running statistics increases, the student accuracy decreases.

| $\mathcal{P}$            | Student accuracy |
|--------------------------|------------------|
| 100 (Running statistics) | 13.49            |
| 90                       | 18.26            |
| 75                       | 57.73            |
| 50                       | 79.54            |
| 25                       | 82.24            |
| 0 (Current statistics)   | 89.4             |

### 21 A.3 Adjusting the student

22 From Figure 2, we observe that the student activations for original data with running statistics (which  
23 are adapted to Gaussian noise) follows a very different distribution compared to other two, and  
24 does not output meaningful information. Neurons in that case either explode or rarely activate. As  
25 described in Section 4.4 of the paper, the remedy is to use current statistics while inference or adjust

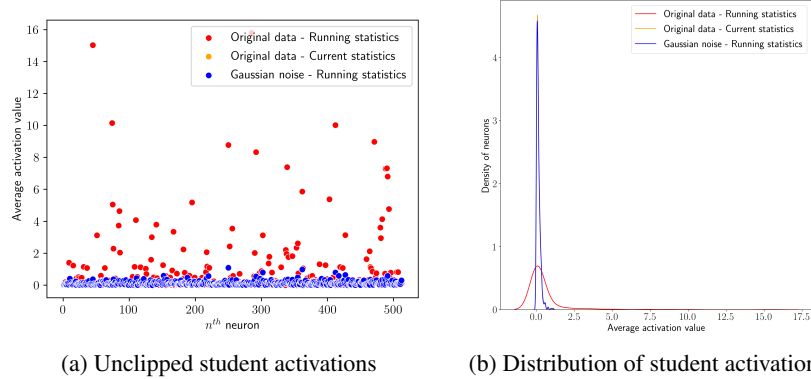


Figure 2: Unclipped scatter plot (linked to Figure 3 (Right) of the paper) and accompanying distribution plot for ‘avgpool’ layer of the *student network* trained for CIFAR10 using our approach.

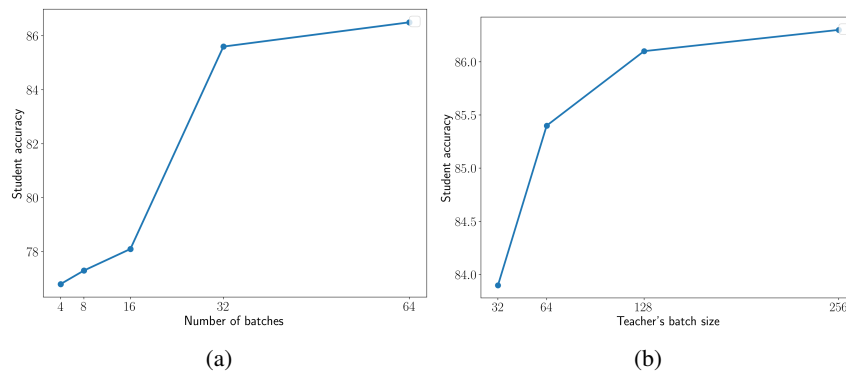


Figure 3: (a) Varying number of batches to adjust student’s running statistics. We found that the number of batches should be sufficient enough to get good student accuracy. (b) Varying teacher’s batch size during distillation. We observe that the larger the batch size, the more accurate the estimates, hence better distillation. The batch size does not need to be too large. Standard batch size for the given dataset should result in good distillation.

26 the running statistics to make them adapt to real data distribution using a small subset of evaluation  
 27 data. Figure 3a further shows the number of batches of evaluation data used to adjust the running  
 28 statistics vs the student’s accuracy when those adjusted running statistics are used. Note that the  
 29 batch size is kept constant at 16. We found that, similar to how the batch size needs to be just large  
 30 enough while using current statistics, here as well, the number of batches or data in general needs to  
 31 be sufficient enough to get a good performance out of the student model.

#### 32 A.4 Significance of Batch size during distillation

33 As stated in Section 5.1 of the paper, as opposed to the traditional knowledge distillation, the teacher  
 34 in our approach needs to rely on the current statistics of the input. In order to get the correct estimates  
 35 in the BatchNorm layer, the batch size during distillation should be large enough. A comparison is  
 36 shown in Figure 3b for the case of CIFAR10, where teacher is ResNet-34 and student is ResNet-18.  
 37 Here, the distillation is performed with the given batch size for both teacher and student models, and  
 38 the same batch size is used while inference to calculate current statistics.

#### 39 A.5 Experiments with Dead Leaves

40 We conduct additional experiments to understand the potential of the proposed approach. In particular  
 41 we notice the missing spatial consistency among pixels in the samples randomly drawn from a  
 42 Gaussian distribution. We, therefore, consider random samples obtained using the Dead leaves



Figure 4: Dead leaf samples used for KD during the experiment with CIFAR10

Table 2: CIFAR10 distillation in different cases (input fed to the networks and state of BatchNorm in teacher network) across various Student network architectures. The numbers are accuracy obtained on the test data. The teacher network here is a ResNet-34, which has an accuracy of 93.29%. The BatchNorm layers in the student model use current statistics during evaluation. Note that RS is running statistics and CS is current statistics.

| Student                     | ResNet34 | ResNet18 | MobileNetV2 | WRN-28-10 | WRN-16-8 |
|-----------------------------|----------|----------|-------------|-----------|----------|
| Original data + RS (Oracle) | 92.74    | 92.44    | 90.57       | 92.41     | 91.32    |
| Gaussian noise + RS         | 13.18    | 13.49    | 12.43       | 14.56     | 14.35    |
| Gaussian noise + CS (Ours)  | 87.11    | 85.98    | 82.47       | 88.12     | 88.76    |
| Dead leaves + RS            | 42.45    | 37.04    | 31.53       | 37.59     | 37.14    |
| Dead leaves + CS (Ours)     | 89.7     | 89.4     | 86.94       | 90.75     | 89.96    |

43 (Shapes)<sup>1</sup> model. Figure 4 shows some of the Dead leaf samples. We repeat the distillation  
 44 experiments on CIFAR10 dataset with Dead leaves samples. We use the same hyperparameter  
 45 values as the Gaussian noise distillation experiments. In addition to ResNet-34, ResNet-18, and  
 46 MobileNetV2, we consider WRN-28-10 and WRN-16-8 as student networks. Table 2 shows the  
 47 obtained results. We observe a considerable improvement in the performance of the proposed  
 48 approach with Dead leaf samples. In particular, for MobileNetV2, test accuracy improves more  
 49 than 4%. Furthermore, for WRN-28-10, the proposed noise-based KD (90.75%) is very close to  
 50 the real data-based KD (92.41%). Moreover, we consider batch size varying from 32 to 512 and  
 51 distill ResNet-34, pretrained on CIFAR10, into ResNet-18 using the proposed approach with Dead  
 52 leaves samples. We observe test accuracy of the student distilled with batch size of 32, 64, 128, 256,  
 53 and 512 as 88.09%, 88.21%, 89.35%, 89.4%, and 89.19%, respectively. This shows that the further  
 54 improvements can be obtained by adjusting the nature of random samples.

55 For further evaluation, we consider DenseNet169-DenseNet121 as teacher-student combination for  
 56 CIFAR10 dataset, and perform distillation using the proposed approach with Dead leaves samples.  
 57 Table 3 shows the obtained results where we observe similar improvements as with the previous  
 58 experiments.

<sup>1</sup>[https://mbaradad.github.io/learning\\_with\\_noise/](https://mbaradad.github.io/learning_with_noise/)

Table 3: Results on CIFAR10 teacher - student pair of DenseNet169-DenseNet121.

| Dataset                     | CIFAR10     |
|-----------------------------|-------------|
| Teacher                     | DenseNet169 |
| Student                     | DenseNet121 |
| Teacher supervised          | 86.23       |
| Original data + RS (Oracle) | 86.23       |
| Gaussian noise + RS         | 10.79       |
| Gaussian noise + CS (Ours)  | 76.57       |
| Dead leaves + RS            | 23.5        |
| Dead leaves + CS (Ours)     | 78.33       |

59 **A.6 Activation distributions**

60 This section contains the distributions of activations of the teacher networks with respect to the model  
61 architectures and datasets that they were trained on.

62 **A.6.1 CIFAR10**

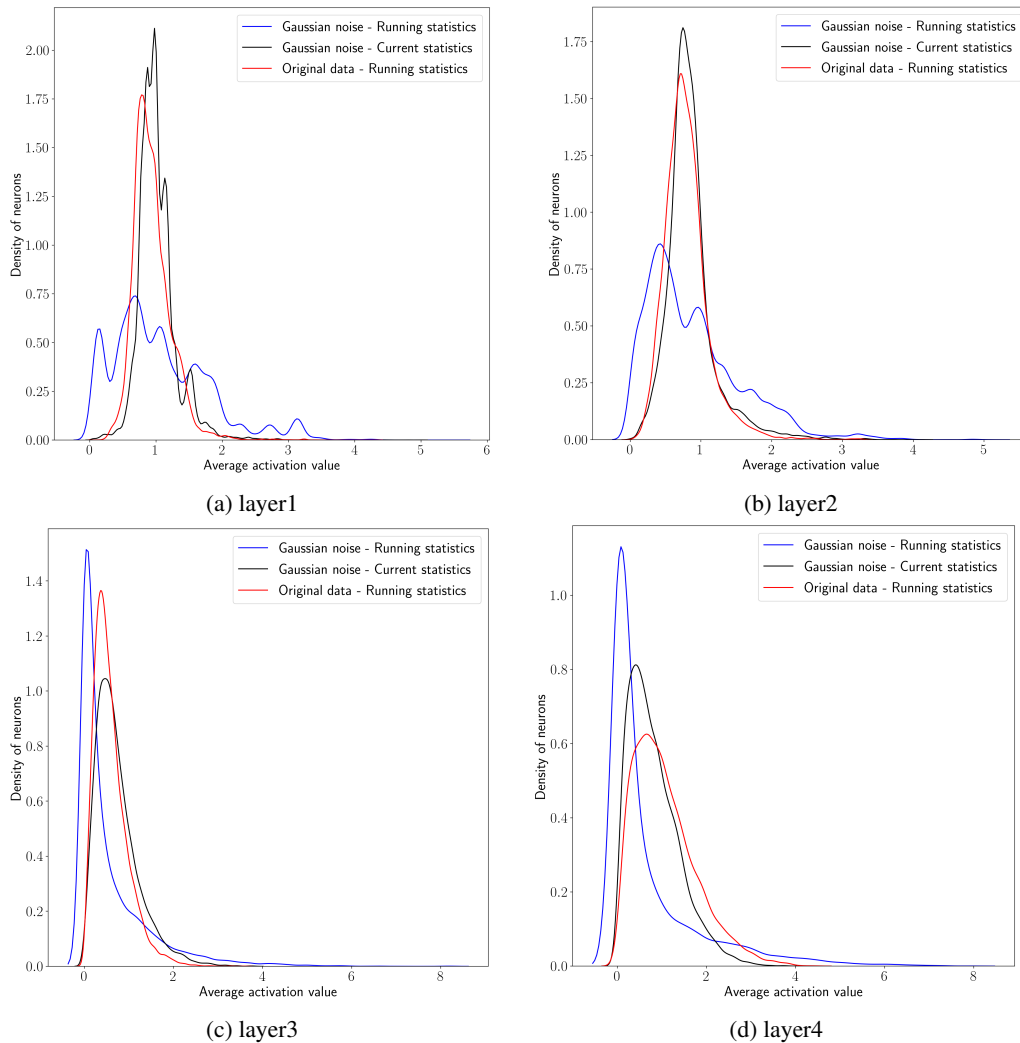


Figure 5: Distributions for average activation value in ResNet-34 teacher network trained on CIFAR-10 for different layers. Note that the layers follow the standard ResNet naming convention.

63 In Figure 5, we observe a similar trend in the activation distribution of all layers. The red curve  
64 denotes the activation distribution of original data with running statistics, which is the ideal case.  
65 However, the trivial use of Gaussian noise with running statistics (blue) results in a significantly  
66 different activation distribution. Our proposed approach makes sure that between Gaussian noise  
67 with running statistics and Gaussian noise with current statistics (black), the activation distribution  
68 of the latter is comparatively similar to the ideal case, thus reducing the shift. We observe a similar  
69 trend throughout all the following plots for SVHN (Figure 6), Food101(Figure 7), and CIFAR100  
70 (Figure 8).

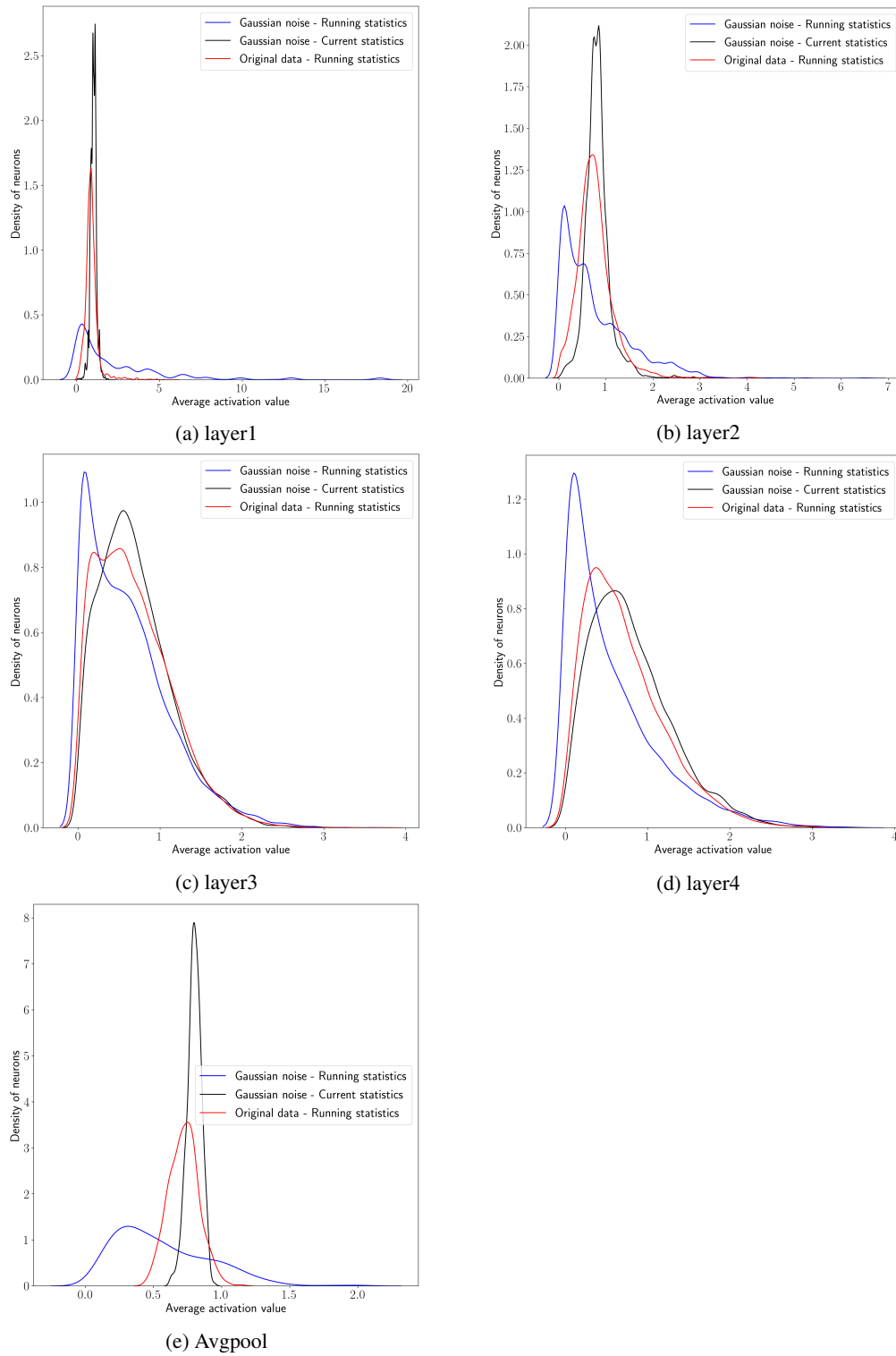


Figure 6: Distributions for average activation value in ResNet-18 teacher network trained on SVHN for different layers. Note that the layers follow the standard ResNet naming convention.

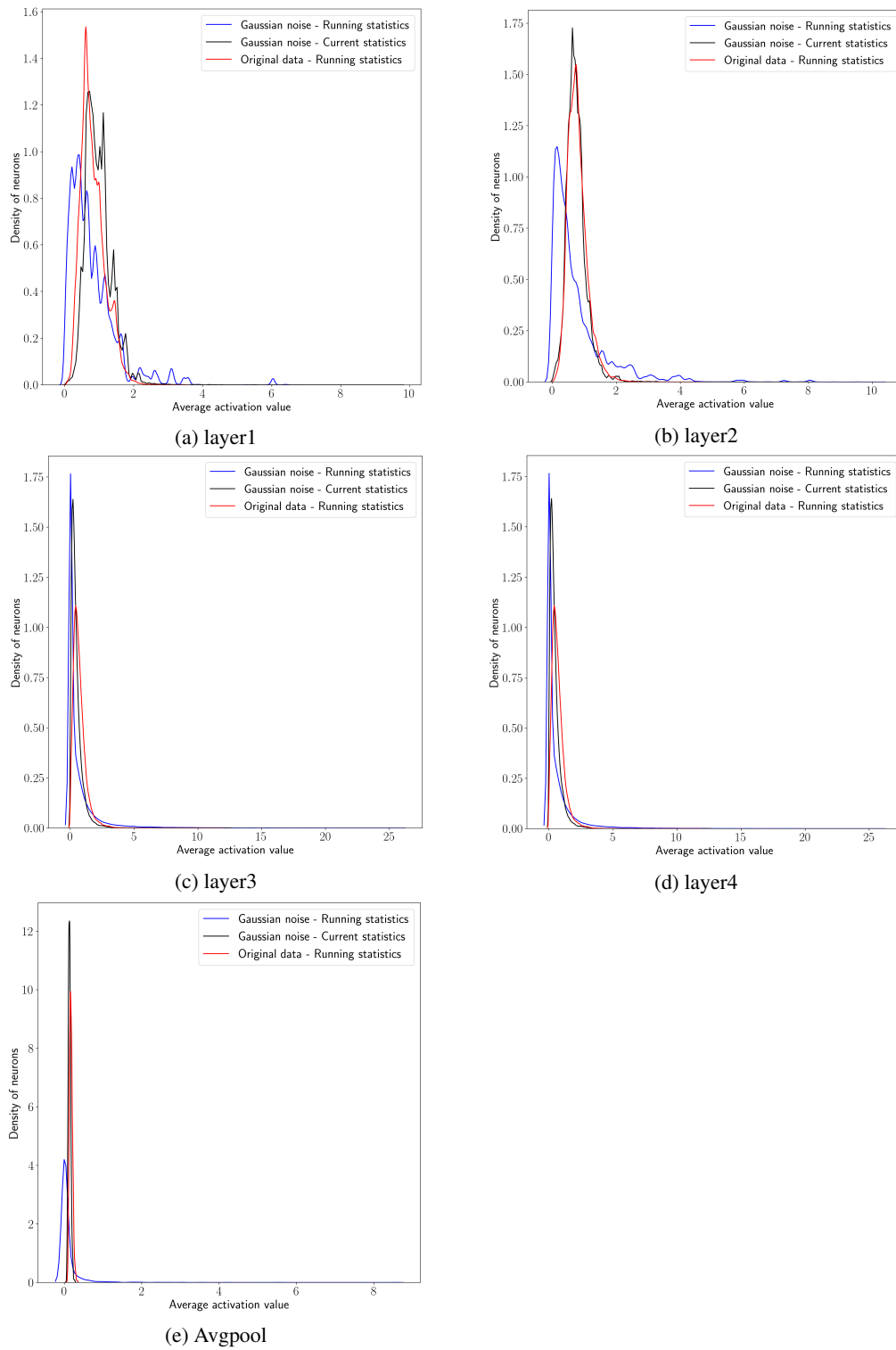


Figure 7: Distributions for average activation value in ResNet-101 teacher network trained on Food101 for different layers. Note that the layers follow the standard ResNet naming convention.

73 **A.6.4 CIFAR100**

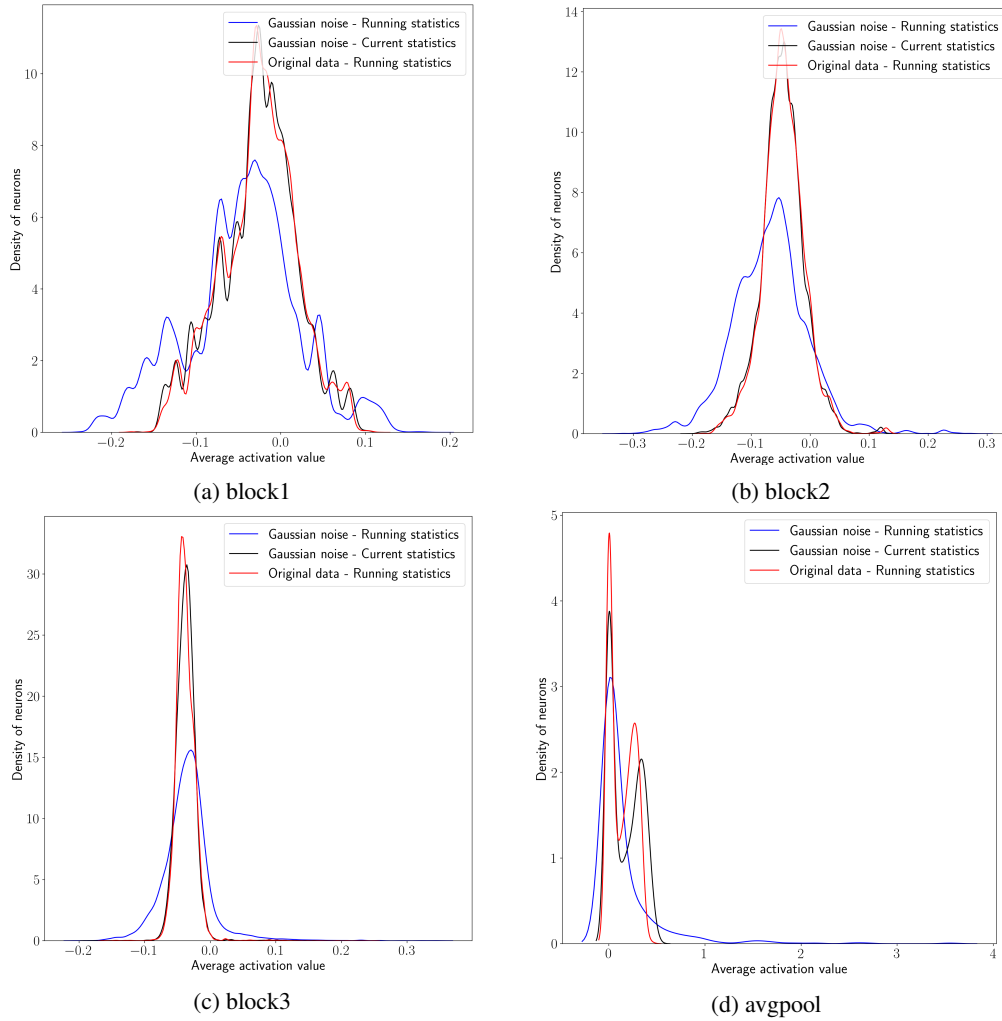


Figure 8: Distributions for average activation value in WideResNet-28-10 teacher network trained on CIFAR100 for different layers. Note that the layers follow the standard Wide-ResNet naming convention.

74 **A.7 Experiment settings**

75 In all of the experiments, we have used standard datasets and model architectures. For Food101, we  
 76 use PyTorch built-in architectures. For SVHN and CIFAR10, we use a custom implementation<sup>2</sup> of  
 77 ResNet architectures, which is tailored towards low-resolution images by making the initial kernel  
 78 shape smaller. For CIFAR100, we use an open-source implementation<sup>3</sup> of Wide-ResNets.

79 In all cases we have used Adam optimizer with a standard learning rate of 0.001, except for the  
 80 training of CIFAR100 teacher, where we use SGD with learning rate of 0.1, momentum of 0.9, and  
 81 weight decay of 0.0005. In all experiments, learning rate scheduler strategy is a standard one to  
 82 reduce learning rate when a plateau hits. We train teacher and student models in each experiment for  
 83 100 epochs and 200 epochs respectively, and wherever not specified we use a batch size of 256.

<sup>2</sup><https://github.com/kuangliu/pytorch-cifar>

<sup>3</sup><https://github.com/xternalz/WideResNet-pytorch>

84 All of the datasets used have predefined train-test splits with labels available for both, hence the same  
85 are used. Note that the train split is used only in case of training the teacher, and test split is used as  
86 evaluation data.

87 **State of BatchNorm.** Handling BatchNorm layers to calculate either running statistics or current  
88 statistics is straightforward that can be done by putting the model in one of the two PyTorch modes  
89 - *train* and *eval*. Traditionally, during distillation using real data, we put the teacher in *eval* mode  
90 where BatchNorm layers use running statistics. For our approach to distill using Gaussian noise, we  
91 need the teacher to use current statistics, which is done by putting the teacher in *train* mode. Note  
92 that, even though the teacher is in *train* mode, it only allows calculation of current statistics for the  
93 BatchNorm layers. No update happens to the teacher’s weights as there is no optimizer bound to the  
94 teacher’s weights.

95 All of the experiments are performed on a system with 16-core Intel x86 CPU with 128 GB RAM  
96 and NVIDIA RTX 3090 GPU.