

---

# Appendix

---

## [KAKURENBO: Adaptively Hiding Samples in Deep Neural Network Training]

---

Anonymous Author(s)

Affiliation

Address

email

### 1 Appendix A. Proof of Lemma 1

2 **Lemma 1.** Let  $F(\mathbf{w}) = \mathbb{E}[f_i(\mathbf{w})]$  be non-convex. Set  $\sigma^2 = \mathbb{E}[\|\nabla f_i(\mathbf{w}_M)\|^2]$  with  $\mathbf{w}^* :=$   
 3  $\operatorname{argmin} F(\mathbf{w})$ . Suppose  $\eta \leq \frac{1}{\sup_i L_i}$ . Let  $\Delta_t = \mathbf{w}_t - \mathbf{w}$ . After  $T$  iterations, SGD satisfies:

$$\mathbb{E} [\|\Delta_T\|^2] \leq (1 - 2\eta\hat{C})^T \|\Delta_0\|^2 + \eta R_\sigma$$

4 where  $\hat{C} = \lambda(1 - \eta \sup_i L_i)$  and  $R_\sigma = \frac{\sigma^2}{\hat{C}}$ .

5 *Proof.*  $\|\nabla f_i(\mathbf{w})\| = 0$  in the noiseless setting, and so  $\sigma := 0$ . For  $\mathbf{x}_k$  being the input at  $i$  random  
 6 index for iteration  $k$ , there exists a parameter  $\lambda_{\mathbf{w}_t}$  for  $\lambda_{max}$  (Eq. 7), and  $w = w_\lambda$ , we have for step  
 7 size  $\gamma$

$$\begin{aligned} \mathbb{E} [\|\Delta_T\|^2] &= \|\mathbf{x}_k - \mathbf{x}_* - \gamma \nabla f_i(\mathbf{x}_k)\|^2 \\ &= \|(\mathbf{x}_k - \mathbf{x}_*) - \gamma(\nabla f_i(\mathbf{x}_k) - \nabla f_i(\mathbf{x}_*)) - \gamma \nabla f_i(\mathbf{x}_*)\|^2 \\ &= \|\mathbf{x}_k - \mathbf{x}_*\|^2 - 2\gamma \mathbf{x}_k - \mathbf{x}_* * \nabla f_i(\mathbf{x}_k) + \gamma^2 \|\nabla f_i(\mathbf{x}_k) - \nabla f_i(\mathbf{x}_*) + \nabla f_i(\mathbf{x}_*)\|^2 \\ &\leq \|\mathbf{x}_k - \mathbf{x}_*\|^2 - 2\gamma \mathbf{x}_k - \mathbf{x}_* * \nabla f_i(\mathbf{x}_k) + 2\gamma^2 \|\nabla f_i(\mathbf{x}_k) - \nabla f_i(\mathbf{x}_*)\|^2 + 2\gamma^2 \|\nabla f_i(\mathbf{x}_*)\|^2 \\ &\leq \|\mathbf{x}_k - \mathbf{x}_*\|^2 - 2\gamma \mathbf{x}_k - \mathbf{x}_* * \nabla f_i(\mathbf{x}_k) \\ &\quad + 2\gamma^2 L_i \|\mathbf{x}_k - \mathbf{x}_* + \nabla f_i(\mathbf{x}_k) - \nabla f_i(\mathbf{x}_*)\|^2 + 2\gamma^2 \|\nabla f_i(\mathbf{x}_*)\|^2 \end{aligned}$$

8 where we employ Jensen's inequality in the first inequality for  $\sigma^2 = \mathbb{E}[\|\nabla f_i(\mathbf{w}_M)\|^2]$ . Then  
 9  $\nabla f_i(\mathbf{x}) = F(\mathbf{x})$ , and we obtain

$$\begin{aligned} \mathbb{E} [\|\Delta_T\|^2] &\leq \|\mathbf{x}_k - \mathbf{x}_*\|^2 - 2\gamma \mathbf{x}_k - \mathbf{x}_* * F(\mathbf{x}_k) + 2\gamma^2 [L_i \|\mathbf{x}_k - \mathbf{x}_* + \nabla f_i(\mathbf{x}_k) - \nabla f_i(\mathbf{x}_*)\|^2 \\ &\quad + 2\gamma^2 \|\nabla f_i(\mathbf{x}_*)\|^2] \\ &\leq \|\mathbf{x}_k - \mathbf{x}_*\|^2 - 2\gamma \mathbf{x}_k - \mathbf{x}_* * F(\mathbf{x}_k) + 2\gamma^2 \sup_i L_i \|\mathbf{x}_k - \mathbf{x}_* + \nabla f_i(\mathbf{x}_k) - \nabla f_i(\mathbf{x}_*)\|^2 \\ &\quad + 2\gamma^2 \|\nabla f_i(\mathbf{x}_*)\|^2 \\ &= \|\mathbf{x}_k - \mathbf{x}_*\|^2 - 2\gamma \mathbf{x}_k - \mathbf{x}_* * F(\mathbf{x}_k) + 2\gamma^2 \sup L \|\mathbf{x}_k - \mathbf{x}_* + \nabla f_i(\mathbf{x}_k) - \nabla f_i(\mathbf{x}_*)\|^2 + 2\gamma^2 \sigma^2 \end{aligned}$$

Table 1: Datasets and Models Used in Experiments (\* Down-stream training using the pre-trained model).

Model	Dataset	#Samples	#Epoch	#GPUs	minibatch (per GPU)	Task
Resnet50 <a href="#">He u. a. (2016)</a>	ImageNet-1K <a href="#">Deng u. a. (2009)</a>	1.2M	100	32	64	Image Classification
EfficientNet-b3 <a href="#">Tan und Le (2019)</a>					32	
WideResNet-28-10 <a href="#">Zagoruyko und Komodakis (2016)</a>	CIFAR-100 <a href="#">Krizhevsky und Hinton (2009)</a>	50K	200	32	32	Image Classification
DeepCAM <a href="#">Kurth u. a. (2018)</a>	DeepCAM <a href="#">Kurth u. a. (2018)</a>	~ 122K	35	1024	1	Image Segmentation
DeiT-Tiny-224 <a href="#">Touvron u. a. (2021)</a>	Fractal-3K <a href="#">Kataoka u. a. (2022)</a>	3M	80	32	16	Image Classification
	(*) CIFAR-10 <a href="#">Krizhevsky und Hinton (2009)</a>	50K	1000	8	96	
	(*) CIFAR-100 <a href="#">Krizhevsky und Hinton (2009)</a>	50K	1000	8	96	

10 when  $\gamma \leq \frac{1}{\sup L}$ . Recursively applying this bound over the first  $k$  iterations yields the desired result

$$\begin{aligned} \mathbb{E} [\|\Delta_T\|^2] &\leq \left(1 - 2\gamma\mu(1 - \gamma)\right)^k \|\mathbf{x}_0 - \mathbf{x}_*\|^2 + 2 \sum_{j=0}^{k-1} \left(1 - 2\gamma\mu(1 - \gamma)\right)^j \gamma^2 \sigma^2 \\ &\leq \left(1 - 2\gamma\mu(1 - \gamma)\right)^k \|\mathbf{x}_0 - \mathbf{x}_*\|^2 + \frac{\gamma\sigma^2}{\mu(1 - \gamma)}. \end{aligned}$$

11

□

## 12 Appendix B. Experiments Details

### 13 B.1. System detail

14 We run our experiments on a supercomputer with 1000s of compute nodes, each equipped with 2  
15 Intel Xeon Gold 6148 CPUs, 384 GiB of RAM, 4 NVidia V100 GPUs, and Infiniband EDR NICs  
16 (100Gbps×2). We run 4 MPI ranks per compute node so that each rank has a dedicated access to a  
17 GPU.

### 18 B.2. Model training method details and dataset information:

19 Table 1 summarizes the models and datasets used in this work. In details, we evaluate KAKURENBO  
20 using several models on various datasets as the following:

- 21 • **ImageNet-1K** [Deng u. a. \(2009\)](#): We use the subset of the ImageNet dataset containing  
22 1000 classes each containing around 1300 images (1,282,048 images in total). We also  
23 test the trained model on the validation set of 50,000 samples. We train ResNET-50 and  
24 EfficientNet-b3 provided by ‘torchvision v0.12.0’ on ImageNet-1K dataset.
- 25 • **CIFAR-10/CIFAR-100** [Krizhevsky und Hinton \(2009\)](#): The CIFAR-10/CIFAR-100 dataset  
26 dataset consists of 60,000 colour images. It has 100 categories each containing 600 images.  
27 The dataset provides 50,000 training images and 10,000 test images with a size of  $32 \times 32$   
28 pixels. CIFAR-100 dataset is available at <https://www.cs.toronto.edu/~kriz/cifar.html>.
- 29 • **DeepCAM** [Kurth u. a. \(2018\)](#): DeepCAM dataset for image segmentation, which consists of  
30 approximately 122K samples and requires 8.8TBs of storage. We use the settings in [Kurth  
31 u. a. \(2018\)](#) to train DeepCAM with the top learning rate of 0.0055.
- 32 • **Fractal-3K** [Kataoka u. a. \(2022\)](#) A rendered dataset from the Visual Atom method [Kataoka  
33 u. a. \(2022\)](#). Fractal-3K dataset comprise of 3 million images of visual atoms, where the  
34 number of classes is  $C = 3000$  and the number of images per class is  $N = 1000$ . We train the  
35 DeiT-Tiny-224 model on Fractal-3K dataset and fine tune it with CIFAR-10 and CIFAR-100  
36 datasets.

Table 2: Hyper-parameters used for different training in the paper and the baseline top-1 testing accuracy. We also considers different hyper-parameters for ResNet-50 model on ImageNet-1K dataset.

	ImageNet-1K				CIFAR-100	Fractal-3K	CIFAR-10	CIFAR-100
	ResNet-50	ResNet-50 (A)	ResNet-50 (B)	EfficientNet-b3	WideResNet-28-10	DeiT-Tiny-224		
Train Res	224	224	224	224	32	224	224	224
Test Res	232	224	232	224	32	-	32	32
Epochs	600	100	600	100	200	80	1000	1000
Number of workers	32	32	32	32	32	32	8	8
Batch size	2048	1024	1024	1024	1024	512	768	768
Optimizer	SGD	SGD	SGD	SRMSProp	SGD	adamw	SGD	SGD
Momentum	0.9	0.9	0.9	0.9	0.9	-	0.9	0.9
LR	0.11	0.0125	0.125	0.01	0.025	0.001	0.01	0.01
Weight decay	1e-5	5e-5	2e-5	5e-5	5e-4	0.05	1e-4	1e-4
LR decay	cosineLR	step	cosineAnnealing	step	step	Cosine_iter	Cosine_iter	Cosine_iter
Decay rate	-	0.1	-	0.9	0.2	-	-	-
Decay epochs	-	[30, 60, 80]	-	2	[60, 120, 160]	-	-	-
Warmup epochs	5	5	5	5	1	5	5	5
Warmup method	linear	linear	linear	linear	linear	linear	linear	linear
Label Smoothing	0.1	-	-	-	-	-	0.1	0.1
H.flip	YES	YES	YES	YES	YES	YES	YES	YES
Erasing prob.	0.1	-	0.1	-	-	0.5	0.5	0.5
Auto augment	ta_wide	-	ta_wide	-	-	rand-m9-mstd0.5-inc1		
Interpolation	bilinear	-	bilinear	-	-	bicubic	bicubic	bicubic
Train crop	176	-	176	-	-	224	224	224
Test crop	224	-	224	-	-	-	-	-
EMA	YES	-	-	-	-	-	-	-
EMA steps	32	-	-	-	-	-	-	-
EMA decay	0.99998	-	-	-	-	-	-	-
Loss	Cross Entropy	Cross Entropy	Cross Entropy	Cross Entropy	Cross Entropy	Cross Entropy	Soft Target Cross Entropy	
<b>Baseline acc.</b>	74.89	73.68	76.58	76.63	77.49	-	95.03	79.69
Max fraction	0.3	0.3	0.3	0.3	0.3	0.3	-	-
Max fraction decay	[1, 0.8, 0.6]	[1, 0.8, 0.6]	[1, 0.8, 0.6]	[1, 0.8, 0.6]	[1, 0.8, 0.6]	[1, 0.8, 0.6]	-	-
Fraction decay epoch	[200, 400, 600]	[30, 60, 80]	[200, 400, 600]	[30, 60, 80]	[60, 120, 160]	[30, 60, 80]	-	-
<b>KAKURENBO acc.</b>	75.15	73.52	76.62	76.23	77.21	-	95.28	79.35

### 37 B.3. Hyper-parameters

38 It is worth noting that we follow the hyper-parameters reported in Vryniotis for training ResNet-  
39 50, Zagoruyko und Komodakis (2016) for training WideResNet-28-10, Tan und Le (2019) for training  
40 EfficientNet-b3, and Kurth u. a. (2018) for DeepCAM. We also use the setting in Kataoka u. a. (2022)  
41 for both pretrain and finetune tasks in Fractal-3K. Table 2 shows the detail of our hyper-parameters.  
42 Specifically, We follow the guideline of ‘TorchVision’ to train the ResNet-50 that uses the CosineLR  
43 learning rate scheduler<sup>1</sup>, auto augments, and random erasing, etc Vryniotis. We also set the weight  
44 decay to  $1e - 05$  and crop the input image to  $176 \times 176$  pixels and train for a long number of epochs,  
45 i.e., 600 (The ResNet-50 setting). We train the WideResNet-28-4 on the CIFAR-100 dataset in  
46 200 epochs following the setting in Zagoruyko und Komodakis (2016). Specifically, we use the  
47 base learning rate of  $0.025 \times k$ , momentum 0.9, and weight decay 0.0005. For EfficientNet-b3, we  
48 use RMSProp optimizer with momentum 0.9; batch norm momentum 0.99 weight decay  $1e - 5$   
49 (following Tan und Le (2019)). We use an initial learning rate of 0.016 that decays by 0.9 every  
50 2 epochs. We set the minibatch size per worker (GPU) to  $b$ , e.g., the global batch size of  $b \times p$  in  
51 the case of  $p$  GPUs. The minibatch size per GPU and the number of GPUs in each experiments are  
52 shown in Table 1.

53 To show the robustness of KAKURENBO, we also train ResNet-50 with different settings, e.g.,  
54 marked as (A) and (B) in the Table 2 and discuss the result in Appendix C.3. For example, in  
55 ResNet-50 (A) setting, we follow the hyper-parameters reported in Goyal u. a. (2017). Specifically,  
56 we use using the Stochastic Gradient Descent (SGD) optimizer with a Nesterov momentum of 0.9  
57 and weight decay of 0.00005. We trained all the models for 100 epochs and apply the linear scaling  
58 rule with the base learning rate of  $0.0125 \times k$  where  $k$  is the number of workers. We reduce the  
59 learning rate by 0.1 at the 30th, 60th, and 80th epoch. We gradual warmup which starts with 0 and is

<sup>1</sup>implemented by timm <https://github.com/huggingface/pytorch-image-models/tree/main/timm>

60 linearly increased to the base learning rate over 5 epochs. We also use scale and aspect ratio data  
 61 augmentation. The input image is a  $224 \times 224$  pixel random crop from an augmented image or its  
 62 horizontal flip.

#### 63 B.4. Implementation

64 It is worth noting that KAKURENBO merely hides samples before the input pipeline. As a result,  
 65 KAKURENBO can be easily implemented with simple extensions to PyTorch and TensorFlow  
 66 implementations<sup>2</sup>. Using KAKURENBO with new models and datasets can be added to any training  
 67 code by indicating so in the model launch parameters.

### 68 Appendix C. Ablation Studies

#### 69 C.1. Analysis of the Factors Affecting KAKURENBO’s Performance

70 In this section, we present an analysis of the factors affecting KAKURENBO’s performance, e.g., the  
 71 lagging loss and the prediction confidence.

72 Figure 1 shows the histogram of the loss as the number of epochs increases when training ResNet-50  
 73 (A) on the ImageNet-1K dataset. At the first few epochs, the histogram of the loss follows a Gaussian  
 74 distribution. As the number of epochs increases, the number of samples with small loss increases  
 75 significantly. For example, starting from epoch 30, more than 50% of the samples have a loss which  
 76 is lower than 5% of the highest loss. As a result, there is an increase in the number of samples that  
 77 provide about the same absolute contribution to the update, e.g., in the latter epochs. Hiding a fraction  
 78 of (fixed)  $F$  samples during training in this case may lead to a relatively higher negative impact on  
 79 the accuracy than that at some early epochs. Thus, we reduce the maximum hidden fraction at the  
 80 epoch number increases (as mentioned in Section 3.4 in the main manuscript).

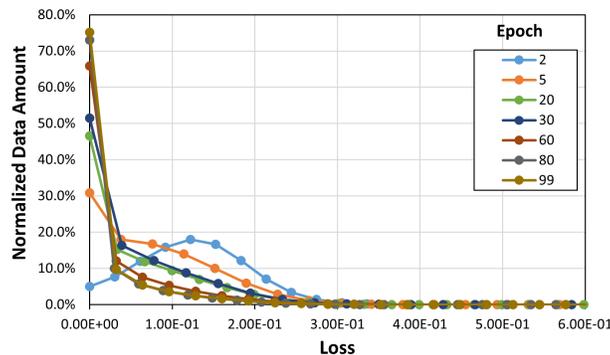


Figure 1: Histogram of the lagging-loss as the number of epoch increases during training (ResNet-50 w/ ImageNet-1K).

81 In addition, as the number of samples with the same absolute loss increases, there is a high probability  
 82 that samples classified as important are in fact unimportant. To this end, we propose moving back  
 83 samples from the hidden set based on their prediction confidence score (as per Section 3.2). Unlike  
 84 the ahead-of-time method proposed by the authors in [Toneva u. a. \(2019\)](#), instead of computing the  
 85 loss of all the samples before training and selecting samples to be removed from the training process,  
 86 we compute the loss of the samples on the fly. With this method, at each epoch, a dynamic hiding  
 87 fraction  $F^*$  is applied. Figure 2 shows the number of hidden samples of each class in KAKURENBO  
 88 (ResNet-50, ImageNet-1K). The figure shows the result of the first 50 classes. The number on top of  
 89 each column shows the rank over 1000 classes (a lower rank indicates a higher number of hidden  
 90 samples). The result shows that our method could dynamically hide the samples at each epoch. For  
 91 example, fewer samples in the class 25 are hidden while more and more samples in class 13 are  
 92 selected to hide as epochs increase.

<sup>2</sup> Our PyTorch implementation is available at <https://anonymous.4open.science/r/kakurenbo-8F10/>

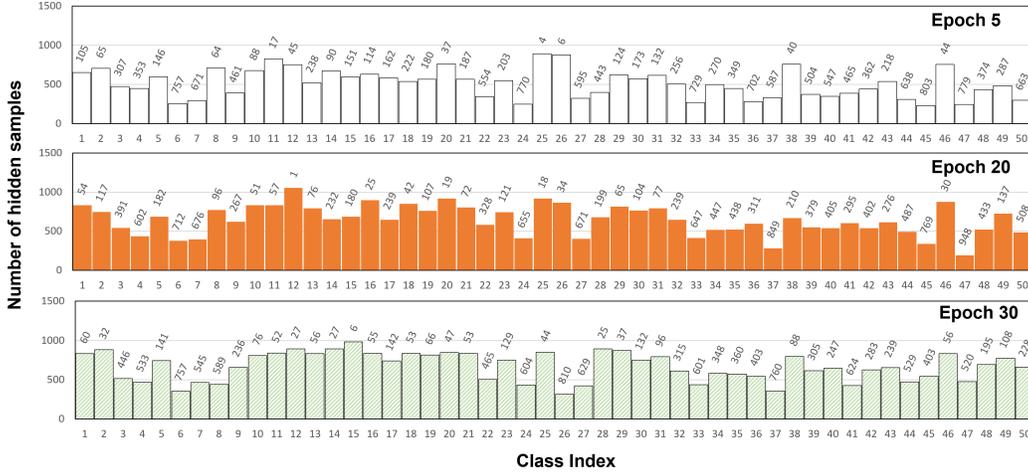


Figure 2: Number of hidden samples of each class in KAKURENBO (ResNet-50, ImageNet-1K). The figure shows the result of the first 50 classes. The number on top of each column shows the rank over 1000 classes (a lower rank indicates a higher number of hidden samples).

93 **C.2. Evolution of the Hiding Fraction**

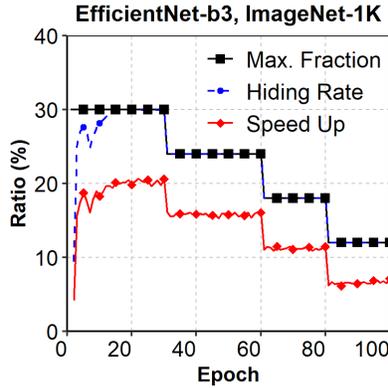


Figure 3: Reduction of hiding fraction, per epoch, and the resulting speedup.

94 Figure 3 shows how KAKURENBO adapts the size of the hidden set during the training of  
 95 EfficientNet-b3. At the beginning of the training, the maximum hiding fraction is set to 30 %.  
 96 This fraction is progressively reduced after a few epochs followed by our fraction adjustment rule.  
 97 The figure also reports the effective proportion of samples that are hidden at each epoch (*Hiding*  
 98 *rate* in the Figure). As described in Section 3 in the main manuscript, KAKURENBO first cuts a  
 99 part of the dataset before moving back samples that are mispredicted or correctly predicted but with  
 100 low confidence. Figure 3 shows that the moving back strategy mostly impacts the beginning of the  
 101 training when the model is still inaccurate.

102 Figure 3 also reports the measured speedup per epoch as compared to the baseline epoch duration.  
 103 The speedup follows the same trend as the hiding rate. This is because reducing the number of  
 104 samples in the training set impacts the speed of the training. The measured speedup does not reach  
 105 the maximum hiding rate because of additional hidden sample selection and due to the need for  
 106 computing the forward pass on samples in the hidden list.

107 **C.3. Robustness of Our Method**

Table 3: Test accuracy (Top-1) in percentage and total training time in seconds of KAKURENBO in the comparison with those of the baseline.

Setting	ResNet-50(A) + ImageNet-1K		ResNet-50(B) + ImageNet-1K	
	Accuracy	Time (sec)	Accuracy	Time (sec)
Baseline	73.68	16118	76.58	64060
KAKURENBO-0.2	-	-	76.11	61723
KAKURENBO-0.3	73.52	12984	76.17	59063
KAKURENBO-0.4	-	-	75.62	57582

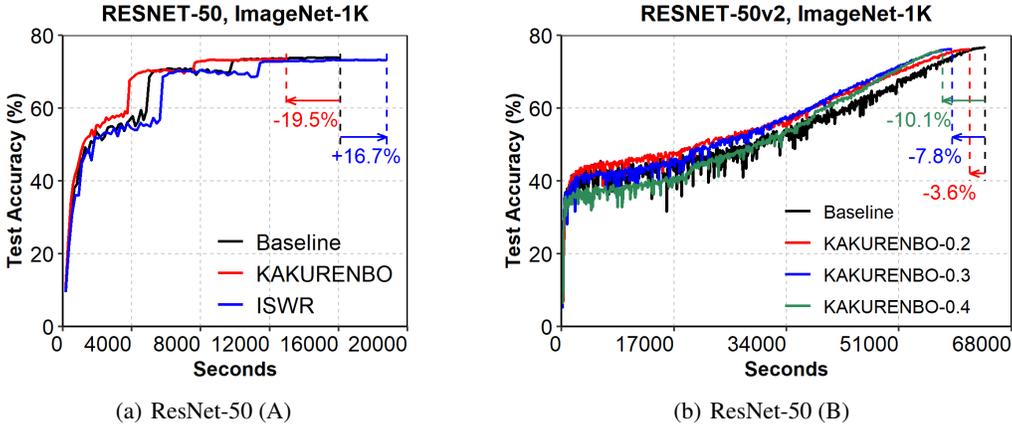


Figure 4: Convergence and speedup of KAKURENBO with different setting of ResNet-50.

108 In this section, we demonstrate the robustness of KAKURENBO with different settings during  
 109 training, e.g. (1) when using different techniques to improve accuracy and (2) the batch size is  
 110 changed.

111 We first measure the robustness of KAKURENBO when using SoTA techniques in training, the  
 112 ResNet-50 (A) and (B) described in Table 2. The result in Figure 4 and Table 3 show that our proposed  
 113 method is also stable with different learning techniques. For example, KAKURENBO could reduce  
 114 the total training time to 19.5% (7.8%) with only 0.2% (0.41) percent of accuracy reduction when  
 115 the maximum hidden fraction is set to 30% for RESNET-50 (A) and (B), respectively.

Table 4: Test accuracy (Top-1) in percentage of KAKURENBO in comparison with those of the baseline when the batch size changes.

Setting	ResNet-50 (A) + ImageNet-1K			
#GPUs	32	64	128	256
Batch size	1024	2048	4096	8192
Baseline	73.68	73.98	73.59	73.81
KAKURENBO-0.4	73.60	73.21	73.03	72.84

116 We now fix the mini-batch size per worker to 32 and then increase the number of workers (GPUs),  
 117 i.e., we increase the global batch size in the case of ResNet-50 (A). Table 4 shows the top-1 testing  
 118 accuracy of ResNet-50 (A) on the ImageNet-1K dataset when the batch size changes from 1024 to  
 119 8192. The result shows that KAKURENBO can maintain the accuracy (or with a trivial reduction of  
 120 accuracy) even with large batch sizes. KAKURENBO could help with large-scale training which has  
 121 become common when training DL models on a large supercomputer or cluster.

Table 5: The impact of different components of KAKURENBO on testing accuracy (Resnet-50 (A), ImageNet-1K,  $F = 0.4$ ) including **HE**: Hiding  $F\%$  lowest-loss examples, **MB**: Moving Back, **RF**: Reducing the Fraction by epoch, **LR**: Adjusting Learning Rate. Numbers inside the (.) indicate the gap in percentage compared to the full version of KAKURENBO.

	Component				Accuracy
	HE	MB	RF	LR	
Baseline	×	×	×	×	73.68
v1000	✓	×	×	×	72.25 (-1.8%)
v1001	✓	×	×	✓	73.08 (-0.7%)
v1010	✓	×	✓	×	72.81 (-1.1%)
v1011	✓	×	✓	✓	73.27 (-0.4%)
v1100	✓	✓	×	×	72.37 (-1.7%)
v1101	✓	✓	×	✓	73.09 (-0.7%)
v1110	✓	✓	✓	×	72.96 (-0.9%)
KAKUR. (v1111)	✓	✓	✓	✓	73.6

#### 122 C.4. Impact of different components of KAKURENBO

123 We evaluate how KAKURENBO’s individual internal strategies, and their combination, affect the  
 124 testing accuracy of a neural network. Table 5 reports the results we obtained when training ResNet-  
 125 50 on ImageNet-1K with a maximum hiding fraction of 40% . The results show that when only  
 126 HE (Hiding Examples) of the 40% lowest loss samples is performed, accuracy slightly degrades.  
 127 Combining HE with other strategies, namely MB (Move-Back), RF (Reducing Fraction), and LR  
 128 (Learning Rate adjustment) gradually improves testing accuracy. In particular, all combinations  
 129 with RF achieve higher accuracy than the ones without it. For example, the accuracy of v1110 is  
 130 higher than that of v1100 by about 0.59%. We also observe that using LR helps to improve the  
 131 training accuracy by a significant amount, i.e., from 0.46% to 0.83%. The MB strategy also improves  
 132 accuracy. For example, the accuracy of v1010 is 72.81%, compared to v1110 which is 72.96%. This  
 133 small impact of MB on the accuracy is due to moving back samples at the beginning of the training,  
 134 as seen in Appendix C.3. By using all the strategies, KAKURENBO achieves the best accuracy of  
 135 73.6%, which is very close to the baseline of 73.68%.

#### 136 Appendix D. Discussion on DeepCAM

137 We have shown how KAKURENBO’s internal strategies, and their combination, affect the testing  
 138 accuracy of a neural network in the case of ResNet-50 and the ImageNet-1K dataset. Figure 5  
 139 presents the same result on the DeepCAM dataset. In this experiment, we evaluate two combinations:  
 140 **v1000** and **v1001**. For **v1000** we hide  $F\%$  lowest-loss samples only (Hiding Example or HE for  
 141 short). For **v1001** we combine HE and learning rate adjustment techniques. It is worth noting that our  
 142 proposed method, **KAKURENBO**, is the combination of HE, LR, MB (Moving Back sample), and  
 143 FR (Reducing the Fraction by epoch). The result with different maximum hidden fractions, e.g.  $F$   
 144 from 0.2 to 0.4, shows that using LR helps to improve the training accuracy by a significant amount,  
 145 and KAKURENBO achieves the best accuracy which is very close to the baseline. This result is  
 146 similar to what we observed with ResNet-50 and the ImageNet-1K dataset.

147 For DeepCAM, we also observed that the loss of the samples with the highest loss does not decrease  
 148 significantly during the last few epochs of training and remain substantially above the rest of other  
 149 samples. Those samples may be hard to learn or represent noise in the data. Figure 6 demonstrates  
 150 this phenomena showing the loss distributions of the full, bottom 98% and top 2% of the dataset  
 151 according to the loss values, respectively. As seen, the top 2%’s loss distribution remains high until  
 152 the very last epoch.

153 This observation motivated us to consider a version in which we cut 2% of the highest-loss samples  
 154 at each epoch (DropTop). Interestingly, it helps to improve the testing accuracy of DeepCAM, e.g.,  
 155 from 77.16% in KAKURENBO to 77.37% with a maximum fraction of 0.3. For version v1001,  
 156 Droptop increases the accuracy by 0.82%.

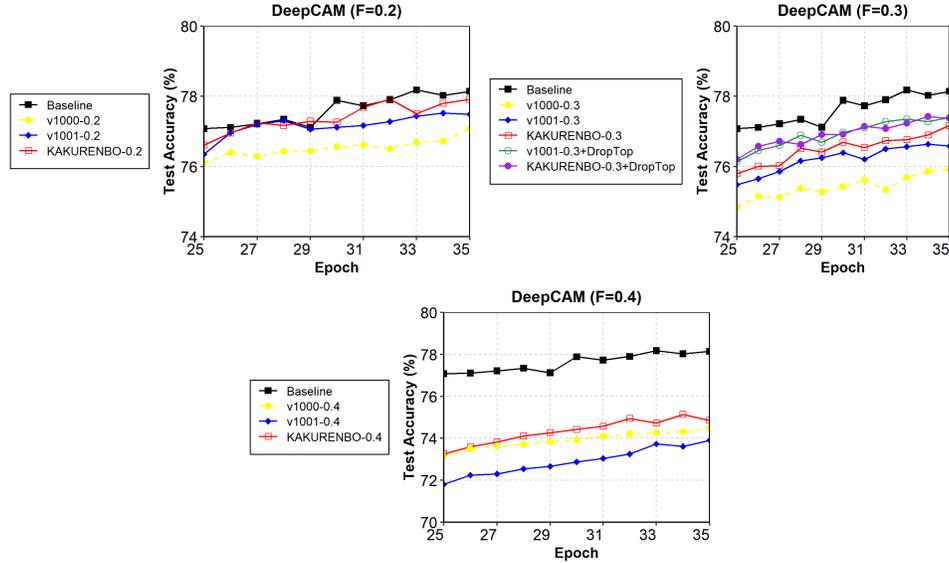


Figure 5: The impact of different components of KAKURENBO on testing accuracy (DeepCAM). **v1000**: Hiding  $F\%$  lowest-loss samples only (HE). **v1001**: HE + LR (Adjusting Learning Rate). **KAKURENBO**: our proposed method with HE + LR + MB (Moving Back) + FR (Reducing the Fraction by epoch). We also consider the version in which we cut 2% of the highest-loss samples at each epoch (DropTop).

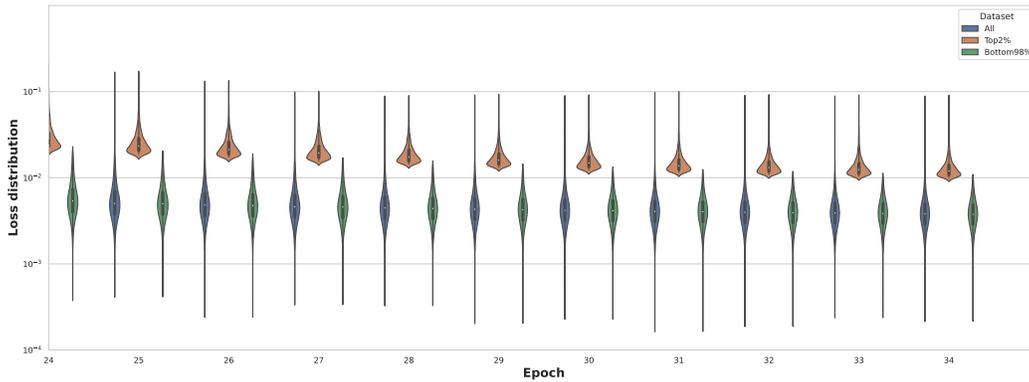


Figure 6: Loss distributions of DeepCAM training samples (full dataset, bottom 98% and top 2%) in the last 10 epoch of training.

## 157 Appendix E. Related work

158 As the size of training datasets and the complexity of deep-learning models increase, the cost of  
 159 training neural networks becomes prohibitive. Several approaches have been proposed to reduce this  
 160 training cost without degrading accuracy significantly.

161 **Biased with-Replacement Sampling** has been proposed as a method to improve the convergence rate  
 162 in SGD training [Katharopoulos und Fleuret \(2018\)](#); [Mindermann u. a. \(2022\)](#). Importance sampling  
 163 is based on the observation that not all samples are of equal *importance* when it comes to training,  
 164 and accordingly replaces the regular uniform sampling used to draw samples from datasets with  
 165 a biased sampling function that assigns a likelihood to a sample being drawn proportional to its  
 166 importance; the more important the sample is, the higher the likelihood it would be selected. The  
 167 with-replacement strategy of importance sampling maintains the total number of samples the network  
 168 trains on.

169 Several improvements over importance sampling have been proposed. Reducible Holdout Loss  
170 Selection (RHO-LOSS) [Mindermann u. a. \(2022\)](#) is a selection function that quantifies by how  
171 much each sample would reduce the loss on unseen data had it been trained on. Mercury uses an  
172 importance-aware data sharding technique in order to speed up distributed training [Zeng u. a. \(2021\)](#).  
173 It distributes important samples across workers between iterations. This allows important samples to  
174 be uniformly distributed between workers, and it reduces the number of samples to communicate for  
175 each epoch since non-important samples are kept local.

176 The importance of a sample can be estimated with several methods. In [Wu u. a. \(2017\)](#), authors  
177 use distance weighted sampling to determine the importance of samples. [Zhao und Zhang \(2015\)](#)  
178 uses stochastic optimization to reduce the stochastic variance. [Allen-Zhu u. a. \(2016\)](#) selects each  
179 coordinate with a probability proportional to the square root of its smoothness parameter (applied  
180 to accelerated coordinate descent). RAIS [Johnson und Guestrin \(2018\)](#) proposes approximating the  
181 ideal sampling distribution, which introduces little computational overhead.

182 Overall, biased with-replacement sampling aims at increasing the convergence speed of SGD by  
183 focusing on samples that induce a measurable change in the model parameters, which would al-  
184 low a reduction in the number of epochs. While these techniques promise to converge in fewer  
185 epochs on the whole dataset, each epoch requires computing the importance of samples which is  
186 time-consuming; and the actual speedup in terms of time-to-solution remains unclear. Moreover,  
187 existing studies [Katharopoulos und Fleuret \(2018\)](#); [Mindermann u. a. \(2022\)](#); [Zeng u. a. \(2021\)](#) only  
188 evaluate small datasets. Our experiments show that the biased with-replacement, importance sam-  
189 pling [Katharopoulos und Fleuret \(2018\)](#), the algorithm does not speedup the training when applied to  
190 large-scale datasets (demonstrated in the evaluation section in the paper).

191 **Data Pruning techniques** are used to reduce the size of the dataset by removing less important  
192 samples. Pruning the dataset requires training on the full dataset and adds significant overheads for  
193 quantifying individual differences between data points [Sorscher u. a. \(2022\)](#). However, the assumption  
194 is that the advantage would be a reduced dataset that replaces the original datasets when used by  
195 others to train. Several studies investigate the selection of the samples to discard from a dataset.  
196 In [Toneva u. a. \(2019\)](#), authors detect unforgettable samples that are correctly classified during the  
197 course of training. EL2N [Paul u. a. \(2021\)](#) uses the loss gradient norm of samples to identify the  
198 important ones and prune the unimportant samples from the dataset after a few epochs. While this  
199 work does not require fully training the model before pruning, it remains unclear if EL2N reduces the  
200 total training time. Another work uses memorization to identify outliers or mislabeled samples in a  
201 given dataset [Feldman und Zhang \(2020\)](#). Removing these atypical samples accelerates the training  
202 without altering the trained model accuracy. Ensemble Active Learning [Chitta u. a. \(2021\)](#) trains an  
203 ensemble of networks and uses ensemble uncertainty to identify which samples are hard to learn.  
204 They manage to reduce the ImageNet dataset by 20% without degrading the accuracy of the trained  
205 model, but again, their method is prohibitive for models and datasets that require excessive resources  
206 for training.

207 Pruning the dataset does reduce the training time without significantly degrading the accuracy [Toneva](#)  
208 [u. a. \(2019\)](#); [Feldman und Zhang \(2020\)](#). However, these techniques require fully training the model  
209 on the whole dataset to identify the samples to be removed, which is compute intensive. While  
210 most of the proposed solutions perform well on small datasets such as CIFAR, many fail to maintain  
211 accuracy on larger datasets like ImageNet [Sorscher u. a. \(2022\)](#).

212 **Selective-Backprop** [Jiang u. a. \(2019\)](#) combines importance sampling and online data pruning. It  
213 reduces the number of samples to train on by using the output of each sample’s forward pass to  
214 estimate the sample’s importance and cuts a fixed fraction of the dataset at each epoch. While this  
215 method shows notable speedups, it has been evaluated only on tiny datasets without providing any  
216 measurements on how accuracy is impacted. In addition, the authors allow up to 10% reduction in  
217 test error in their experiments. EIF [Wu u. a. \(2020\)](#) is similar to Selective-Backprop: it reduces the  
218 computation cost of training by filtering out the samples with the lowest loss. E<sup>2</sup>-Train [Wang u. a.](#)  
219 [\(2019\)](#) shows that the combination of randomly dropping samples during training with selective layer  
220 update in CNNs can significantly reduce the training time, while slightly degrading the accuracy.  
221 However, E<sup>2</sup>-Train targets edge environments and is evaluated only on very small datasets.

222 **GRAD-MATCH** [Killamsetty u. a. \(2021\)](#) is an online method that selects a subset of the samples  
223 that would minimize the gradient matching error, where the error of the gradients of a matched subset  
224 samples (and their weights) becomes minimum. To avoid the impractical storing and computation

225 of the optimization of the gradients of all instances, the authors approximate the gradients by only  
 226 using the gradients of the last layer, use a per-class approximation, and run data selection every  $R$   
 227 epochs, in which case, the same subsets and weights will be used between epochs. The infrequent  
 228 selection, however, means the model is limited in its capacity to learn in intermediate epochs - where  
 229 selection occurs - since it trains on the same limited subset of samples. This often leads to a longer  
 230 numbers of epochs needed to converge to the same validation accuracy that can be achieved by the  
 231 baseline or the baseline reaching much higher accuracy Pooladzandi u. a. (2022). Another important  
 232 point worth mentioning is that GRAD-MATCH is impractical in distributed training, which is a de  
 233 facto requirement in large dataset and models (e.g., the DeepCAM model/dataset). That is since the  
 234 approximation of the classes would require very expensive high-volume collective communication  
 235 operations to gather the gradients scattered across different samples belonging to the same class. The  
 236 communication cost would be  $O(N.R.G)$  where  $N$  is the number of samples,  $R$  is the frequency of  
 237 selection, and  $G$  is the gradients (of the last layer, if gradient approximation is to be used). Distributed  
 238 GRAD-MATCH would require a scatter communication to collect the class approximations and a  
 239 collective all-reduce of the gradients to then do the matching optimization. This is practically a very  
 240 high cost for communication per epoch that could even exceed the average time per epoch. Finally,  
 241 the mini-batch variant of GRAD-MATCH can only be effective for small mini-batches. However,  
 242 since in distributed training the mini-batch grows with the scale (i.e., the mini-batch aggregates the  
 243 local mini-batch of all workers), the cost of communication amplifies by  $B$  (where  $B$  is mini-batch  
 244 size).

## 245 References

- 246 [Allen-Zhu u. a. 2016] ALLEN-ZHU, Zeyuan ; QU, Zheng ; RICHTÁRIK, Peter ; YUAN, Yang: Even  
 247 faster accelerated coordinate descent using non-uniform sampling. In: *International Conference*  
 248 *on Machine Learning* PMLR (Veranst.), 2016, S. 1110–1119 9
- 249 [Chitta u. a. 2021] CHITTA, Kashyap ; ÁLVAREZ, José M ; HAUSSMANN, Elmar ; FARABET,  
 250 Clément: Training data subset search with ensemble active learning. In: *IEEE Transactions on*  
 251 *Intelligent Transportation Systems* (2021) 9
- 252 [Deng u. a. 2009] DENG, Jia ; DONG, Wei ; SOCHER, Richard ; LI, Li-Jia ; LI, Kai ; FEI-FEI, Li:  
 253 Imagenet: A large-scale hierarchical image database. In: *Conf. on comp. vision and pattern recogn.*  
 254 *IEEE* (Veranst.), 2009, S. 248–255 2
- 255 [Feldman und Zhang 2020] FELDMAN, Vitaly ; ZHANG, Chiyuan: What Neural Networks  
 256 Memorize and Why: Discovering the Long Tail via Influence Estimation. In: *Proceedings of*  
 257 *the 34th International Conference on Neural Information Processing Systems*, 2020 (NIPS'20). –  
 258 ISBN 9781713829546 9
- 259 [Goyal u. a. 2017] GOYAL, Priya ; DOLLÁR, Piotr ; GIRSHICK, Ross ; NOORDHUIS, Pieter ;  
 260 WESOLOWSKI, Lukasz ; KYROLA, Aapo ; TULLOCH, Andrew ; JIA, Yangqing ; HE, Kaiming:  
 261 Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour. In: *CoRR* abs/1706.02677 (2017)  
 262 3
- 263 [He u. a. 2016] HE, Kaiming ; ZHANG, Xiangyu ; REN, Shaoqing ; SUN, Jian: Deep residual  
 264 learning for image recognition. In: *CVPR*, 2016, S. 770–778 2
- 265 [Jiang u. a. 2019] JIANG, Angela H. ; WONG, Daniel L. K. ; ZHOU, Giulio ; ANDERSEN, David G. ;  
 266 DEAN, Jeffrey ; GANGER, Gregory R. ; JOSHI, Gauri ; KAMINKSY, Michael ; KOZUCH, Michael ;  
 267 LIPTON, Zachary C. ; PILLAI, Padmanabhan: *Accelerating Deep Learning by Focusing on the*  
 268 *Biggest Losers*. 2019. – URL <https://arxiv.org/abs/1910.00762> 9
- 269 [Johnson und Guestrin 2018] JOHNSON, Tyler B. ; GUESTRIN, Carlos: Training deep models faster  
 270 with robust, approximate importance sampling. In: *Advances in Neural Information Processing*  
 271 *Systems* 31 (2018) 9
- 272 [Kataoka u. a. 2022] KATAOKA, Hirokatsu ; HAYAMIZU, Ryo ; YAMADA, Ryosuke ; NAKASHIMA,  
 273 Kodai ; TAKASHIMA, Sora ; ZHANG, Xinyu ; MARTINEZ-NORIEGA, Edgar J. ; INOUE, Nakamasa ;  
 274 YOKOTA, Rio: Replacing Labeled Real-Image Datasets With Auto-Generated Contours. In:  
 275 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022,  
 276 S. 21232–21241 2, 3

- 277 [Katharopoulos und Fleuret 2018] KATHAROPOULOS, Angelos ; FLEURET, François: Not all  
 278 samples are created equal: Deep learning with importance sampling. In: *International conference*  
 279 *on machine learning* PMLR (Veranst.), 2018, S. 2525–2534 8, 9
- 280 [Killamsetty u. a. 2021] KILLAMSETTY, Krishnateja ; DURGA, S ; RAMAKRISHNAN, Ganesh ;  
 281 DE, Abir ; IYER, Rishabh: Grad-match: Gradient matching based data subset selection for efficient  
 282 deep model training. In: *International Conference on Machine Learning* PMLR (Veranst.), 2021,  
 283 S. 5464–5474 9
- 284 [Krizhevsky und Hinton 2009] KRIZHEVSKY, A. ; HINTON, G.: Learning multiple layers of  
 285 features from tiny images. In: *Master’s thesis, Dep. of Comp. Sci. Univ. of Toronto* (2009) 2
- 286 [Kurth u. a. 2018] KURTH, Thorsten ; TREICHLER, Sean ; ROMERO, Joshua ; MUDIGONDA,  
 287 Mayur ; LUEHR, Nathan ; PHILLIPS, Everett ; MAHESH, Ankur ; MATHESON, Michael ;  
 288 DESLIPPE, Jack ; FATICA, Massimiliano u. a.: Exascale deep learning for climate analytics.  
 289 In: *SC18: International Conference for High Performance Computing, Networking, Storage and*  
 290 *Analysis* IEEE (Veranst.), 2018, S. 649–660 2, 3
- 291 [Mindermann u. a. 2022] MINDERMANN, Sören ; BRAUNER, Jan M. ; RAZZAK, Muhammed T. ;  
 292 SHARMA, Mrinank ; KIRSCH, Andreas ; XU, Winnie ; HÖLTGEN, Benedikt ; GOMEZ, Aidan N. ;  
 293 MORISOT, Adrien ; FARQUHAR, Sebastian ; GAL, Yarin: Prioritized Training on Points that  
 294 are Learnable, Worth Learning, and not yet Learnt. In: *Proceedings of the 39th International*  
 295 *Conference on Machine Learning*, PMLR, Jun 2022, S. 15630–15649. – ISSN 2640-3498 8, 9
- 296 [Paul u. a. 2021] PAUL, Mansheej ; GANGULI, Surya ; DZIUGAITE, Gintare K.: Deep  
 297 Learning on a Data Diet: Finding Important Examples Early in Training. In: RAN-  
 298 ZATO, M. (Hrsg.) ; BEYGELZIMER, A. (Hrsg.) ; DAUPHIN, Y. (Hrsg.) ; LIANG,  
 299 P.S. (Hrsg.) ; VAUGHAN, J. W. (Hrsg.): *Advances in Neural Information Processing*  
 300 *Systems* Bd. 34, URL [https://proceedings.neurips.cc/paper/2021/file/  
 301 ac56f8fe9eea3e4a365f29f0f1957c55-Paper.pdf](https://proceedings.neurips.cc/paper/2021/file/ac56f8fe9eea3e4a365f29f0f1957c55-Paper.pdf), 2021, S. 20596–20607 9
- 302 [Pooladzandi u. a. 2022] POOLADZANDI, Omead ; DAVINI, David ; MIRZASOLEIMAN, Baharan:  
 303 Adaptive second order coresets for data-efficient machine learning. In: *International Conference*  
 304 *on Machine Learning* PMLR (Veranst.), 2022, S. 17848–17869 10
- 305 [Sorscher u. a. 2022] SORSCHER, Ben ; GEIRHOS, Robert ; SHEKHAR, Shashank ; GANGULI,  
 306 Surya ; MORCOS, Ari S.: *Beyond neural scaling laws: beating power law scaling via data pruning.*  
 307 2022. – URL <https://arxiv.org/abs/2206.14486> 9
- 308 [Tan und Le 2019] TAN, Mingxing ; LE, Quoc: EfficientNet: Rethinking Model Scaling for  
 309 Convolutional Neural Networks. In: CHAUDHURI, Kamalika (Hrsg.) ; SALAKHUTDINOV, Ruslan  
 310 (Hrsg.): *Proceedings of the 36th International Conference on Machine Learning* Bd. 97, PMLR, 09–  
 311 15 Jun 2019, S. 6105–6114. – URL [https://proceedings.mlr.press/v97/tan19a.  
 312 html](https://proceedings.mlr.press/v97/tan19a.html) 2, 3
- 313 [Toneva u. a. 2019] TONEVA, Mariya ; SORDONI, Alessandro ; COMBES, Remi T. des ;  
 314 TRISCHLER, Adam ; BENGIO, Yoshua ; GORDON, Geoffrey J.: An Empirical Study of Example  
 315 Forgetting during Deep Neural Network Learning. In: *7th International Conference on Learning*  
 316 *Representations (ICLR)*, URL <https://openreview.net/forum?id=BJlxm30cKm>,  
 317 2019 4, 9
- 318 [Touvron u. a. 2021] TOUVRON, Hugo ; CORD, Matthieu ; DOUZE, Matthijs ; MASSA, Francisco ;  
 319 SABLAYROLLES, Alexandre ; JEGOU, Herve: Training data-efficient image transformers &  
 320 distillation through attention. In: MEILA, Marina (Hrsg.) ; ZHANG, Tong (Hrsg.): *Proceedings of*  
 321 *the 38th International Conference on Machine Learning* Bd. 139, PMLR, 18–24 Jul 2021, S. 10347–  
 322 10357. – URL <https://proceedings.mlr.press/v139/touvron21a.html> 2
- 323 [Vryniotis ] VRYNIOTIS, Vasilis: *How to Train State-Of-The-Art*  
 324 *Models Using TorchVision’s Latest Primitives.* – URL [https://  
 325 pytorch.org/blog/how-to-train-state-of-the-art-models/  
 326 -using-torchvision-latest-primitives/](https://pytorch.org/blog/how-to-train-state-of-the-art-models/-using-torchvision-latest-primitives/) 3

- 327 [Wang u. a. 2019] WANG, Yue ; JIANG, Ziyu ; CHEN, Xiaohan ; XU, Pengfei ; ZHAO, Yang ;  
328 LIN, Yingyan ; WANG, Zhangyang: E2-train: Training state-of-the-art cnns with over 80% energy  
329 savings. In: *Advances in Neural Information Processing Systems* 32 (2019) 9
- 330 [Wu u. a. 2017] WU, Chao-Yuan ; MANMATHA, R ; SMOLA, Alexander J. ; KRAHENBUHL,  
331 Philipp: Sampling matters in deep embedding learning. In: *Proceedings of the IEEE International*  
332 *Conference on Computer Vision*, 2017, S. 2840–2848 9
- 333 [Wu u. a. 2020] WU, Yawen ; WANG, Zhepeng ; SHI, Yiyu ; HU, Jingtong: Enabling on-device  
334 cnn training by self-supervised instance filtering and error map pruning. In: *IEEE Transactions on*  
335 *Computer-Aided Design of Integrated Circuits and Systems* 39 (2020), Nr. 11, S. 3445–3457 9
- 336 [Zagoruyko und Komodakis 2016] ZAGORUYKO, Sergey ; KOMODAKIS, Nikos: Wide Residual  
337 Networks. In: *BMVC*, September 2016, S. 87.1–87.12 2, 3
- 338 [Zeng u. a. 2021] ZENG, Xiao ; YAN, Ming ; ZHANG, Mi: Mercury: Efficient On-Device  
339 Distributed DNN Training via Stochastic Importance Sampling. In: *Proceedings of the 19th ACM*  
340 *Conference on Embedded Networked Sensor Systems*, 2021, S. 29–41 9
- 341 [Zhao und Zhang 2015] ZHAO, Peilin ; ZHANG, Tong: Stochastic optimization with importance  
342 sampling for regularized loss minimization. In: *international conference on machine learning*  
343 PMLR (Veranst.), 2015, S. 1–9 9