ASYNCHRONOUS STOCHASTIC GRADIENT DESCENT WITH DECOUPLED BACKPROPAGATION AND LAYER WISE UPDATES

Anonymous authors

Paper under double-blind review

Abstract

The increasing size of deep learning models has created the need for more efficient alternatives to the standard error backpropagation algorithm, that make better use of asynchronous, parallel and distributed computing. One major shortcoming of backpropagation is the interlocking between the forward phase of the algorithm, which computes a global loss, and the backward phase where the loss is backpropagated through all layers to compute the gradients, which are used to update the network parameters. To address this problem, we propose a method that parallelises SGD updates across the layers of a model by asynchronously updating them from multiple threads. Furthermore, since we observe that the forward pass is often much faster than the backward pass, we use separate threads for the forward and backward pass calculations, which allows us to use a higher ratio of forward to backward threads than the usual 1:1 ratio, reducing the overall staleness of the parameters. Thus, our approach performs asynchronous stochastic gradient descent using separate threads for the loss (forward) and gradient (backward) computations and performs layer-wise partial updates to parameters in a distributed way. We show that this approach yields close to state-of-the-art results while running up to $2.97 \times$ faster than Hogwild! scaled on multiple devices (Locally-Partitioned-Asynchronous-Parallel SGD). We theoretically prove the convergence of the algorithm using a novel theoretical framework based on stochastic differential equations and the drift diffusion process, by modeling the asynchronous parameter updates as a stochastic process.

032 033 034

006

008 009 010

011

013

014

015

016

017

018

019

021

023

024

025

026

028

029

031

1 INTRODUCTION

Scaling up modern deep learning models requires massive resources and training time. Asyn chronous parallel and distributed methods for training them using backpropagation play a very
 important role in easing the demanding resource requirements for training these models. Back propagation (BP) (Werbos, 1982) has established itself as the de facto standard method for learning
 in deep neural networks (DNN). Although BP achieves state-of-the-art accuracy on literally all rel evant machine learning tasks, it comes with a number of inconvenient properties that prohibit an
 efficient implementation at scale.

BP is a two-phase synchronous learning strategy in which the first phase (forward pass) computes the
training loss, *L*, given the current network parameters and a batch of data. In the second phase, the
gradients are propagated backwards through the network to determine each parameter's contribution
to the error, using the same weights (transposed) as in the forward pass (see Equation 1). BP suffers
from update locking, where a layer can only be updated after the previous layer has been updated.
Furthermore, the computation of gradients can only be started after the loss has been calculated in
the forward pass.

Moreover, the backward pass usually requires approximately twice as long as the forward pass (Kumar et al., 2021). The bulk of the computational load comes from the number of matrix multiplications required during each phase. If we consider a DNN with M layers, then at any layer $m \le M$ with pre-activations $z_m = \theta_m y_{m-1}$ and post-activations $y_m = f(z_{m-1})$, the computations during the forward pass are dominated by one matrix multiplication $\theta_m y_{m-1}$. During the backward pass, the computations at layer m are dominated by two matrix multiplications:

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}_{m}} = f'(\boldsymbol{\theta}_{m}\boldsymbol{y_{m-1}}) \times \boldsymbol{y_{m-1}}^{\top} \quad \text{and} \quad \frac{\partial \mathcal{L}}{\partial \boldsymbol{y_{m-1}}} = f'(\boldsymbol{\theta}_{m}\boldsymbol{y_{m-1}}) \times \boldsymbol{\theta}_{m}^{\top}, \quad (1)$$

approximately doubling the compute budget required for the backward pass compared to the forward pass. In Eq. 1, θ_m denotes the network weights at layer m and f' the partial derivative with respect to θ_m . This imbalance between forward and backward phase further complicates an efficient parallelization of BP, particularly in heterogeneous settings.

062 In this work, we propose a new approach to parallelize training of deep networks on non-convex 063 objective functions by asynchronously performing the forward and backward passes at a layer-wise 064 granularity in multiple separate threads that make lock-free updates to the parameters in shared 065 memory. The lock-free updates address the locking problem, performing layer-wise updates mit-066 igates the issue of conflicts between parameter updates, and performing the backward pass and 067 updates using more threads than the forward pass mitigates the staleness problem. Specifically, the 068 imbalance in execution time between forward and backward passes is taken care of by having twice as many backward threads than forward threads, breaking the 1:1 ratio of vanilla Backpropagation, 069 therefore significantly speeding-up the training process.

- ⁰⁷¹ In summary, the contributions of this paper are as follows:
 - 1. We introduce a novel asynchronous formulation of Backpropagation which allows the forward pass and the backward pass to be executed separately and in parallel, which allows us to run more backward than forward threads. This approach accounts for the unequal time required by the forward and backward passes.
 - 2. We propose to asynchronously update the model's parameters at a layer-wise granularity without using a locking mechanism which reduces staleness.
 - 3. We give convergence guarantees of the algorithm to a stationary distribution centered around the local optima of conventional BP.
 - 4. We show that the algorithm can reach state-of-the-art performances while being significantly faster than competing asynchronous algorithms.
 - 2 RELATED WORK

056 057

073

074

075

076

077

078 079

080

081

082

085

Asynchronous stochastic gradient descent (SGD). Asynchronous SGD has a long history, starting from Baudet (1978); Bertsekas & Tsitsiklis (2015). Hogwild! (Recht et al., 2011) allows multiple processes to perform SGD without any locking mechanism on shared memory. Kungurtsev et al. (2021) proposed PASSM and PASSM+, where they partition the model parameters across the workers on the same device to perform SGD on the partitions. Chatterjee et al. (2022) decentralizes Hogwild! and PASSM+ to allow parameters or their partitions to be located on multiple devices and perform Local SGD on them. Zheng et al. (2017) compensate the delayed gradients with a gradient approximation at the current parameters. Unlike these methods, we run multiple backward passes in parallel on different devices and don't need any gradient compensation scheme.

Nadiradze et al. (2021) provides a theoretical framework to derive convergence guarantees for a wide 096 variety of distributed methods. Mishchenko et al. (2022) proposes a method of "virtual iterates" to 097 provide convergence guarantees independent of delays. More recently, Even et al. (2024) proposed 098 a unified framework for convergence analysis of distributed algorithms based on the AGRAF frame-099 work. There have been lots of other analysis methods proposed for deriving convergence guarantees for asynchronous distributed SGD (see Assran et al. (2020) for a survey). In our work, we propose an 100 entirely novel framework based on stochastic differential equations, and provide convergence guar-101 antees of the algorithm to a stationary distribution centered around the local optima of conventional 102 BP. 103

Communication-efficient algorithms. One of the bottlenecks when training on multiple devices
 or nodes in parallel is the synchronization step. The bigger or deeper the models get, the more time
 is consumed by synchronization. PowerSGD computes low-rank approximations of the gradients
 using power iteration methods. Poseidon (Zhang et al., 2017) also factorizes gradient matrices
 but interleaves their communication with the backward pass. Wen et al. (2017) and Alistarh et al.

108 (2017) quantize gradients to make them lightweight for communication. Like Zhang et al. (2017), 109 we interleave the backward pass with gradients communication but without gradients averaging. 110

Block local learning. Dividing the network across multiple devices and performing local updates is 111 a widely recognized approach in distributed learning. The backward passes of the different blocks 112 can be done simultaneously. The global loss is used to provide feedback only to the output block 113 while the remaining blocks get learning signals from auxiliary networks which each compute local 114 targets. Jaderberg et al. (2016) models synthetic gradients through auxiliary networks. Ma et al. 115 (2024) uses a shallower version of the network itself as the auxiliary network at each layer. Gomez 116 et al. (2022) allows gradients to flow to k-neighboring blocks. Nøkland & Eidnes (2019) don't 117 allow gradients to flow to neighboring blocks, and instead use an auxiliary matching loss and a local cross-entropy loss to compute the local error. Decoupled Parallel Backpropagation (Huo et al., 118 2018) does full Backpropagation but uses stored stale gradients in the blocks to avoid update locking, 119 therefore needing additional memory buffers. Kappel et al. (2023) take a probabilistic approach by 120 interpreting layers outputs as parameters of a probability distribution. Auxiliary networks provide 121 local targets, which are used to train each block individually. Similar to these distributed paradigms, 122 we mimic the execution of multiple backward passes in parallel by reordering the training sequence 123 but without splitting the network explicitly during forward and backward propagation across devices 124 and needing external buffers or architectural complexities. 125

3 **METHODS**



Figure 1: Illustration of decoupled backpropagation with separate threads for the forward and back-140 ward passes and the layer-wise updates. For each thread, the order of computations for a sample network with three layers denoted L1, L2 and L3, are shown. Arrows denote dependencies across 142 threads. Within each thread, the computations for each layer are performed sequentially, whereas 143 across threads, the dependencies are layer-wise. Interactions for 2 backward threads and a single 144 forward thread are shown. This asynchronous interaction, along with layer-wise updates, reduces 145 the staleness of parameters.

146 147 148

149 150

126 127

128 129 130

131

132 133

134

135

136 137

138

139

141

ASYNCHRONOUS FORMULATION OF BACKPROPAGATION 3.1

We introduce a new asynchronous stochastic gradient descent method where, instead of performing 151 the forward and backward phases sequentially, we execute them in parallel and perform layer-wise 152 parameter updates as soon as the gradients for a given layer are available. The dependencies between 153 forward and backward phases are illustrated in Figure 1. 154

155 Since the gradient computation in the backward pass tends to consume more time than the loss calculation in the forward pass, we decouple these two into separate threads and use one forward 156 thread and two backward threads to counterbalance the disproportionate execution time. 157

158 Figure 1 illustrates the interaction among threads based on one example. Initially, only the first 159 forward pass, $F_0^{(0)}$, is performed. The resulting loss is then used in the first backward pass $B_0^{(1)}$, which starts in parallel to the second forward pass $F_1^{(1)}$. Once $F_1^{(1)}$ ends, its loss is used by $B_1^{(2)}$ running in parallel to the next forward pass and $B_0^{(1)}$. 160 161

162 3.2 LAYER-WISE UPDATES 163

164 Parallelizing the forward and backward passes can speed up training, but it violates several key 165 assumptions of Backpropagation leading to sub-optimal convergence observed in different studies (Keuper & Preundt, 2016; Zheng et al., 2017). This happens because the losses and gradients are 166 often calculated using inconsistent and outdated parameters. 167

168 To alleviate this problem, we update the layers as soon as the corresponding gradients are available from the backward pass. $F_1^{(1)}$ receives partial parameter updates from $B_0^{(1)}$ as soon as they are available. Therefore, the parameters used in $F_1^{(1)}$ will differ from those used in $F_0^{(0)}$ because some 169 170 171 layers of the model would have been already updated by the thread $B_0^{(1)}$. On average, we can expect 172 that the second half of the layers use a new set of parameters. It is important to note that the updates 173 happen without any locking mechanism and asynchronous to the backward pass as done by Zhang 174 et al. (2017). 175

176 177

178

179

181 182

183

185 186

187

188

189

190

191

193

194

195 196

197

3.3 SPEED-UP ANALYSIS

Before discussing experimental results, we study the potential speed-up of the asynchronous with layer-wise updates formulation over standard Backpropagation. To arrive at this result, we make the following assumptions to estimate the performance gain

- We assume that there are no delays between the end of a forward pass, the beginning of its corresponding backward pass and the next forward pass. This implies for example that as soon as $F_0^{(0)}$ ends, $F_1^{(1)}$ and $B_0^{(1)}$ begin immediately. Multiples backward threads are therefore running in parallel.
- Vanilla Backpropagation performs b forward passes. We assume that this number also corresponds to the number of backward passes and the number of batches of data to be trained on.
 - A forward pass lasts T units of time and a backward pass βT units of time, with a scaling factor $\beta > 1$ (expected to be at around 2 as show in appendix A.3).

192 The speed-up factor λ observed can be express as the fraction between the estimated time taken by the standard BP, T_1 , over the Async version of BP, T_2 . Due to the sequential nature of BP, $T_1 = (1 + \beta)bT$. Similarly, $T_2 = (b + \beta)T$ since the backward pass runs parallel to the forward pass. The speed-up factor λ is given by:

 $\lambda = \frac{(1+\beta)b}{(b+\beta)} \,.$

199 Considering a large number of batches, $b \to \infty$, we have

200 201

202

203

204

205 206 Hence, the maximum achievable speedup is expected to be $1 + \beta$, where β is the scaling factor of the backward pass time. In practice, the speed-up factor λ can be influenced by multiples factors like data loading which is sometimes a bottleneck (Leclerc et al., 2023; Isenko et al., 2022), or the system overhead, which reduce the achievable speedup.

 $\lambda = 1 + \beta$.

207 3.4 STALENESS ANALYSIS

208 Here, we demonstrate the advantage of applying layer-wise updates (LU) compared to block updates 209 (BU). BU refers to performing updates only after the entire backward pass is complete, a technique 210 used in various previous asynchronous learning algorithms, e.g. (Recht et al., 2011; Chatterjee et al., 211 2022; Zheng et al., 2017). We use the same notation as in section 3.3. 212

213 To express this formally, we define the relative staleness τ of BU compared to LU as the time delay between when the gradients become available and when they are used to update the model weights. 214 The intuition behind this lies in the fact that the more the updates are postponed, the more likely the 215 gradients will become stale. The staleness will only increase with time and accumulate across the 216 layers. Assuming that the time required to compute the gradients for each layer is uniform and equal 217 to $\frac{\beta T}{M}$, the relative staleness is expressed as $\tau = \frac{\beta T(M-1)}{2}$.

To see this, we use that by definition, the staleness increases as we approach the output layer. At any layer m, the layer-wise staleness $\tau_m = \frac{\beta T}{M}m$. Averaging over the layers, we have

$$\tau = \sum_{m=1}^{M} \tau_m = \frac{\beta T}{M} \sum_{m=1}^{M} m = \beta T \frac{(M-1)}{2}$$

Clearly, τ increases with the network's depth and the time required to perform one backward pass. Thus, the speedup is expected to scale approximately linearly with the network depth, showing the advantage of LU over BU for large M.

3.5 Algorithm

221 222

224

225

226

227 228

229

245

246 247

248

249

250 251

253 254

255

256

257 258

259

260

The Async BP algorithm is illustrated in Figure 1 and described in Algorithm listing 1. The al gorithm consists of two components: a single forward thread and multiple backward threads. All
 threads work independently and asynchronously without any locking mechanism.

The forward thread is solely responsible for computing the loss $\mathcal{L}_i(\theta^u, x_i, y_i)$, given the current mini batch of data $(x_i, y_i) \in \mathcal{D}$ and the latest set of updated weights θ^u . Since the algorithm works asynchronously, the weights θ^u can be updated by any backward thread even while forward pass progresses. Once the forward pass is done, \mathcal{L}_i is sent to one of the backward threads and the forward thread moves to the next batch of data.

In parallel, a backward thread k receives a loss \mathcal{L}_j and performs the backward pass. At each layer m, the gradients $G(\theta_{m,k}^v) = \frac{\partial \mathcal{L}_j}{\partial \theta_{m,k}^v}$ are computed, after which $\theta_{m,k}^v$ is immediately used to update the forward thread parameters. Note that the backward thread here can potentially calculate the gradients for different values of parameters $\theta_{m,k}^v$ than the ones used for the forward pass θ^u . In Section 5 and appendix B we show that this algorithm closely approximates conventional stochastic gradient descent, if asynchronous parameter updates arrive sufficiently frequently.

Algorithm 1 Async BP with decoupled partial updates

Forward thread

Given: Data: $(x_i, y_i) \in \mathcal{D}$, latest up-to-date parameters: θ^u

Compute $\mathcal{L}_i(\theta^u, x_i, y_i)$

send(\mathcal{L}_i) // send loss to a backward thread

Backward Thread *k* (running in parallel to the forward thread)

Given: Loss: \mathcal{L}_j , learning rate: η

for layer $m \in [M,1]$ do

Compute $G(\theta_{m,k}^v)$

$$\theta_{m,k}^{v+1} \leftarrow \theta_{m,k}^v - \eta.G$$

 $\theta_m \leftarrow \theta_{m,k}^{m,n}$ // asynchronously update forward thread

261

262 263 264

265

4 **RESULTS**

end for

We evaluate our method on three vision tasks, CIFAR-10, CIFAR-100 and Imagenet, and on one
sequence modeling task: IMDb sentiment analysis. We use Resnet18 and Resnet50 architectures for
vision tasks and a LSTM network for sequence modelling. These networks are trained on a machine
with 3 NVIDIA A100 80GB PCIe GPUs with two AMD EPYC CPUs sockets of 64 cores each. The
experiment code is based on the C++ frontend of Torch (Paszke et al., 2019) - Libtorch.

The Performance of these tasks is compared to Locally-Asynchronous-Parallel SGD (LAPSGD)
 and Locally-Partitioned-Asynchronous-Parallel SGD (LPPSGD) (Chatterjee et al., 2022). These
 methods extend the well-known Hogwild! algorithm and Partitioned Asynchronous Stochastic Sub gradient (PASSM+) to multiple devices, respectively.

We record the achieved accuracy on the tasks and the wall-clock time to reach a target accuracy (TTA). If not stated otherwise, this accuracy is chosen to be the best accuracy achieved by the worst performing algorithm. We used the code made available by Chatterjee et al. (2022) which uses Pytorch Distributed Data-Parallel API.

- 279
- 279 280 281

282

4.1 ASYNCHNOUS TRAINING OF VISION TASKS

We follow the training protocol of LAPSGD and LPPSGD and chose the number of processes per 283 GPU to be 1 since the GPU utilization was close to 100%. We trained them with a batch size of 284 128 per-rank. We used Stochastic gradient descent (SGD) for both Async BU and LU, with an 285 initial learning rate of 0.005 for 5 epochs and 0.015 after the warm-up phase, a momentum of 0.9 286 and a weight decay of 5×110^{-2} . We use a cosine annealing schedule with a T_{max} of 110. We 287 trained Resnet-50 on Imagenet-1K task (Table 5) for a total of 250 epochs with the same learning 288 rate but with a cosine schedule with T_{max} of 250 and weight decay of 3.5×110^{-2} . Although, 289 the simulations were run with cosine annealing scheduler, implying the ideal number of training 290 epochs, early stopping was applied, i.e. training was stopped if no improvement of the accuracy was 291 achieved for 30 epochs.

292 As shown below, Async LU achieves the highest accuracies while Async BU converges the fastest 293 in terms of time to reach the target accuracy. The CIFAR-10 and CIFAR-100 results are presented 294 in Tables 1, 2 and Tables 3, 4 respectively. In Tables 1 and 3, the time to target accuracy (TTA) is 295 chosen to be the time taken to achieved the best accuracy reached by the worst algorithm. Whereas 296 in Tables 2 and 4, it represents taken by an algorithm achieve its best accuracy. Async BU achieves a 297 speed-up of up to $2.97 \times$ over LPPSGD on CIFAR100 (see Table 3). The poor performance of both 298 LAPSGD and LPPSGD can be explained by the influence of staleness, thus requiring large number of training epochs. 299

We also achieved promising results on the ImageNet-1k dataset (see Table 5). Async LU achieved
 73% accuracy ×3 faster than Backpropagation on single GPU, showing potential of ideal linear
 scaling. An extensive comparison of Async LU with multi-GPU Backpropagation (Data Distributed
 Parallel) is provided in appendix A.2.

Although Async BU converges quicker than Async LU, it reaches lower accuracy. This is particularly visible on CIFAR100, a harder task than CIFAR10 (see Figures 2 and 3). Overall, Async BU showed a good balance between convergence speed and reduction of staleness.

307 308

309	Table 1: Comparison of Async LU, Async BU, LAPSGD and LPPSGD based on time to reach
310	accuracy (TTA): 87% for ResNet18 and 89% for ResNet50, and the number of epochs to reach the
311	target for 3 runs on CIFAR10.

Network architecture	Training method	TTA (in seconds) mean \pm std	Epochs mean \pm std
ResNet-18	Async LU Async BU LAPSGD LPPSGD	$\begin{array}{c} 223.8 \pm 28 \\ 173.4 \pm 2 \\ 706.3 \pm 13 \\ 461.0 \pm 7 \end{array}$	$\begin{array}{c} 65 \pm {}_{10} \\ 64 \pm {}_{2} \\ 104 \pm {}_{2} \\ 86 \pm {}_{1} \end{array}$
ResNet-50	Async LU Async BU LAPSGD LPPSGD	$\begin{array}{c} 737.9 \pm 24 \\ 700.6 \pm 20 \\ 999.5 \pm 38 \\ 863.4 \pm 35 \end{array}$	$\begin{array}{c} 86 \pm 6 \\ 86 \pm 1 \\ 117 \pm 3 \\ 119 \pm 1 \end{array}$

Network architecture	Training method	Best accuracy mean \pm std	TTA (in seconds) mean \pm std	Epochs mean \pm std
ResNet-18	SGD Async LU Async BU LAPSGD LPPSGD	$\begin{array}{c} 93.7 \pm 0.16 \\ 93.7 \pm 0.28 \\ 92.7 \pm 0.16 \\ 88.2 \pm 0.43 \\ 87.8 \pm 0.09 \end{array}$	$\begin{array}{c} \pm \\ 380.7 \pm {}_{33} \\ 308.5 \pm {}_{14} \\ 799.5 \pm {}_{20} \\ 523.4 \pm {}_{15} \end{array}$	$\begin{array}{c} 114 \pm 1 \\ 114 \pm 9 \\ 115 \pm 2 \\ 118 \pm 2 \\ 97 \pm 1 \end{array}$
ResNet-50	SGD Async LU Async BU LAPSGD LPPSGD	$\begin{array}{c} 94.1 \pm 0.20 \\ 93.9 \pm 0.10 \\ 93.2 \pm 0.25 \\ 89.7 \pm 0.27 \\ 89.3 \pm 0.43 \end{array}$	$\begin{array}{c} \pm \\ 1038.8 \pm 61 \\ 953.7 \pm 73 \\ 1098.9 \pm 30 \\ 888.4 \pm 37 \end{array}$	$\begin{array}{c} 99 \pm 5 \\ 121 \pm 6 \\ 116 \pm 8 \\ 117 \pm 3 \\ 119 \pm 1 \end{array}$

Table 2: Comparison of Async LU, Async BU, LAPSGD, and LPPSGD based on best accuracy, time to reach accuracy (TTA), and epoch at the accuracy is achieved for 3 runs on CIFAR10.

Table 3: Comparison of Async LU, Async BU, LAPSGD and LPPSGD based on time to reach accuracy (TTA): 60% for ResNet18 and 63% for ResNet50, and the number of epochs to reach the target for 3 runs on CIFAR100.

Network architecture	Training method	TTA (in seconds) mean \pm std	Epochs mean \pm std
ResNet-18	Async LU Async BU LAPSGD LPPSGD	$\begin{array}{c} 226.4 \pm 10 \\ 155.4 \pm 10 \\ 667.1 \pm 8 \\ 672.9 \pm 11 \end{array}$	69 ± 7 57 ± 4 99 ± 2 100 ± 2
ResNet-50	Async LU Async BU LAPSGD LPPSGD	$\begin{array}{c} 622.7 \pm {}_{32} \\ 641.8 \pm {}_{16} \\ 1006.6 \pm {}_{13} \\ 1016.0 \pm {}_{6} \end{array}$	$70 \pm 5 \\ 81 \pm 3 \\ 107 \pm 1 \\ 107 \pm 1$

353 354 355

356 357

326 327 328

4.2 ASYNCHRONOUS TRAINING OF SEQUENCE MODELLING TASK

For demonstrating Async BP training on sequence modelling, we evaluated an LSTM networks on the IMDb dataset (Maas et al., 2011). Sentiment analysis is the task of classifying the polarity of a given text. We used a 2-Layer LSTM network with 256 hidden dimensions to evaluate this task. We trained the network until convergence using the Adam optimizer with an initial learning rate of 1×10^{-2} . Results are shown in Table 6.

We observe that although both Async LU and Async BU achieve the same accuracy, performing layer-wise updates reduces the training time and number of steps to convergence almost by half (see Figure 4 in Appendix A) highlighting again the importance of this strategy.

365 366 367

368 369

363

364

5 THEORETICAL ANALYSIS OF CONVERGENCE

Here, we theoretically analyse the convergence behavior of the algorithm outlined above. For the theoretical analysis, we consider the general case of multiple threads, acting on the parameter set θ , such that the threads interact asynchronously and can work on outdated versions of the parameters. We model the evolution of the learning algorithm as a continuous-time stochastic process (Bellec et al., 2017) to simplify the analysis. This assumption is justified by the fact that learning rates are typically small, and therefore the evolution of network parameters is nearly continuous.

In the model studied here, the stochastic interaction between threads is modelled as noise induced by random interference of network parameters. To arrive at this model, we use the fact that the dynamics of conventional stochastic gradient descent (SGD) can be modelled as the system of stochastic

Network architecture	Training method	Best accuracy mean \pm std	TTA (in seconds) mean \pm std	Epochs mean \pm std
ResNet-18	SGD	73.7 ± 0.08	±	99 ± 4
	Async LU	73.9 ± 0.38	339.0 ± 10	114 ± 9
	Async BU	71.8 ± 0.08	$277.2\pm$ 17	103 ± 5
	LAPSGD	61.3 ± 0.23	$762.5 \pm {\scriptstyle 15}$	114 ± 3
	LPPSGD	$61.0\pm$ 0.17	$742.0 \pm {\scriptstyle 19}$	$110\pm {\rm 2}$
ResNet-50	SGD	76.8 ± 0.44	±	104 ± 1
	Async LU	76.2 ± 0.57	1071.0 ± 80	122 ± 6
	Async BU	73.7 ± 0.23	$956.2 \pm {\scriptstyle 28}$	123 ± 4
	LAPSGD	$63.7\pm$ 0.33	$1110.1\pm {\scriptstyle 27}$	118 ± 3
	LPPSGD	63.5 ± 0.10	$1122.1 \pm {\scriptstyle 24}$	$119 \pm {\scriptstyle 2}$

Table 4: Comparison of Async LU, Async BU, LAPSGD and LPPSGD based on best accuracy, time
 to reach accuracy (TTA), and epoch at the accuracy is achieved for 3 runs on CIFAR100.

Table 5: Classification accuracy (% correct) for Async LU (3 GPUs) and vanilla backpropagation (single GPU) on the ImageNet task.

Network architecture	training method	test-accuracy	train-accuracy	TTA (in 1000 seconds)
ResNet-50	Async LU	73.42	93.08	134.24
	BP (single GPU)	73.40	91.90	403.77

differential equations that determine the dynamics of the parameter vector θ

$$d\theta_k = -\eta \frac{\partial}{\partial \theta_k} \mathcal{L}(\boldsymbol{\theta}) dt + \frac{\eta \,\sigma_{\text{SGD}}}{\sqrt{2}} \, d\mathcal{W}_k \,, \qquad (2)$$

with learning rate η and where dW_k are stochastic changes of the Wiener processes.

407 Eq. 2 describes the dynamics of a single parameter θ_k . The dynamics is determined by the gradient 408 of the loss function \mathcal{L} , and the noise induced by using small mini-batches modelled here as idealized 409 Wiener process with amplitude σ_{SGD} . Because of this noise, SGD does not strictly converge to a 410 local optimum but maintains a stationary distribution $p^*(\theta_k) \propto e^{-\frac{1}{\eta}\mathcal{L}(\theta_k)}$, that assigns most of the 411 probability mass to parameter vectors that reside close to local optima (Bellec et al., 2017).

In the concurrent variant of SGD studied here, however, the dynamics is determined by perturbed gradients for different stale parameters. When updating the network using the described asynchronous approach without locking, we potentially introduce noise in the form of partially stale parameters or from one thread overwriting the updates of another. This noise will introduce a deviation from the ideal parameter vector θ . We model this deviation as additive Gaussian noise $\boldsymbol{\xi} \sim \mathcal{N}(0, \sigma_{\text{STALE}})$ to the current parameter vector with variance σ_{STALE} . To approximate the noisy loss function, we use a first-order Taylor expansion around the noise-free parameters:

419 420 421

424

380

382

384

394

395

397

399 400 401

 $\mathcal{L}(\boldsymbol{\theta} + \boldsymbol{\xi}) = \mathcal{L}(\boldsymbol{\theta}) + \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})^{\top} \boldsymbol{\xi} + \mathcal{O}(\sigma^2) \\ \approx \mathcal{L}(\boldsymbol{\theta}) + \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})^{\top} \boldsymbol{\xi},$ (3)

and thus the gradient can be approximated as

$$\nabla_{\theta} \mathcal{L}(\boldsymbol{\theta} + \boldsymbol{\xi}, \boldsymbol{X}, \boldsymbol{Y}) \approx \nabla_{\theta} \mathcal{L}(\boldsymbol{\theta}) + \nabla_{\theta}^{2} \mathcal{L}(\boldsymbol{\theta})^{\top} \boldsymbol{\xi}.$$
(4)

Based on this, we can express the update rule as a Stochastic Differential Equation (SDE) and model the various noise terms using a Wiener Process W. The noise sources in the learning dynamics come from two main sources, (1) noise caused by stochastic gradient descent, and (2) noise caused by learning with outdated parameters. We model the former as additive noise with amplitude σ_{STALE} and the latter using the Taylor approximation Eq. (4). Using this, we can write the approximate dynamics of the parameter vector θ as the stochastic differential equation

 $d\theta_k = \mu_k(\boldsymbol{\theta}, t) + \sqrt{D_k(\boldsymbol{\theta})} \, d\mathcal{W}_k \,, \tag{5}$



Table 6: Comparison of Async LU, Async BU based on best accuracy, time to reach accuracy (TTA), and epoch at the accuracy is achieved for 3 runs on IMDb

Figure 2: Learning curves of Asynchronous SGD with layer-wise updates (Async LU) and Block updates (Async BU) on the CIFAR100 dataset. 3 independent runs are shown for each class.

with

$$\mu_{k}(\boldsymbol{\theta}) = -\eta \frac{\partial}{\partial \theta_{k}} \mathcal{L}(\boldsymbol{\theta})$$

$$D_{k}(\boldsymbol{\theta}) = \frac{\eta^{2} \sigma_{\text{SGD}}^{2}}{2} + \frac{\eta^{2} \sigma_{\text{STALE}}^{2}}{2} \sum_{l} \frac{\partial^{2}}{\partial \theta_{k} \partial \theta_{l}} \mathcal{L}(\boldsymbol{\theta}) , \qquad (6)$$

where μ_k is the drift and D_k the diffusion of the SDE.

In Appendix B we study the stationary distribution of this parameter dynamics. We show that the stationary distribution is a close approximation to p^* of SGD, which is perfectly recovered if σ_{STALE} is small compared to σ_{SGD} , i.e. if the effect of staleness is small compared to the noise induced by minibatch sampling.

6 DISCUSSION

 In this work, we introduced a novel asynchronous approach to train deep neural networks that decouples the forward and backward passes and performs layer-wise parameter updates. Our method addresses key limitations of standard backpropagation by allowing parallel execution of forward and backward passes and mitigating update locking through asynchronous layer-wise updates.

485 The experimental results demonstrate that our approach can achieve comparable or better accuracy than synchronous backpropagation and other asynchronous methods across multiple vision and lan-

guage tasks, while providing significant speedups in training time. On CIFAR-10 and CIFAR-100, we observed speedups of up to 2.97× compared to asynchronous SGD covering a broad range of paradigms. The method also showed promising results on a sentiment analysis task and the ImageNet classification task where it reached close to ideal scaling.

Our theoretical analysis, based on modeling the learning dynamics as a continuous-time stochastic
 process, provides convergence guarantees and shows that the algorithm converges to a stationary
 distribution closely approximating that of standard SGD under certain conditions. This offers a
 solid foundation for understanding the behavior of our asynchronous approach.

While our implementation using C++ and LibTorch demonstrated the potential of this method, we also identified some limitations related to GPU resource allocation in SIMT architectures. Future work could explore optimizing the implementation for more efficient GPU utilization, or investigating hybrid CPU-GPU approaches to fully leverage the benefits of asynchronous execution.

Overall, this work presents a promising direction for scaling up deep learning through asynchronous,
 decoupled updates. The approach has the potential to enable more efficient training of large-scale
 models, particularly in distributed and heterogeneous computing environments. Further research
 could explore extensions to even larger models, additional tasks, and more diverse hardware setups
 to fully realize the potential of this asynchronous training paradigm.

505 REPRODUCIBILITY

506

504

We ensure that the results presented in this paper are easily reproducible using just the information provided in the main text as well as the supplement. Details of the models used in our simulations are presented in the main paper and further elaborated in the supplement. We provide additional details and statistics over multiple runs in the supplement section A.4. We use publicly available libraries and datasets in our simulations. We will further provide the source code to the reviewers and ACs in an anonymous repository once the discussion forums are opened. This included code will also contain "readme" texts to facilitate easy reproducibility. The theoretical analysis provided in section 5 is derived in the supplement.

514 515

520

521

522

523

526

527

528

529

530

516 REFERENCES

- 517 Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. Qsgd:
 518 Communication-efficient sgd via gradient quantization and encoding, 2017. URL https:
 519 //arxiv.org/abs/1610.02132.
 - Mahmoud Assran, Arda Aytekin, Hamid Reza Feyzmahdavian, Mikael Johansson, and Michael G Rabbat. Advances in asynchronous parallel and distributed optimization. *Proceedings of the IEEE*, 108(11):2013–2031, 2020.
- 524 Gerard M Baudet. Asynchronous iterative methods for multiprocessors. *Journal of the ACM* 525 (*JACM*), 25(2):226–244, 1978.
 - Guillaume Bellec, David Kappel, Wolfgang Maass, and Robert Legenstein. Deep rewiring: Training very sparse deep networks. *arXiv preprint arXiv:1711.05136*, 2017.
 - Dimitri Bertsekas and John Tsitsiklis. *Parallel and distributed computation: numerical methods*. Athena Scientific, 2015.
- Bapi Chatterjee, Vyacheslav Kungurtsev, and Dan Alistarh. Scaling the wild: Decentralizing hogwild!-style shared-memory sgd, 2022. URL https://arxiv.org/abs/2203.06638.
- Mathieu Even, Anastasia Koloskova, and Laurent Massoulie. Asynchronous SGD on Graphs: A
 Unified Framework for Asynchronous Decentralized and Federated Optimization. In *Proceedings* of *The 27th International Conference on Artificial Intelligence and Statistics*, pp. 64–72. PMLR, April 2024. URL https://proceedings.mlr.press/v238/even24a.html.
- Aidan N. Gomez, Oscar Key, Kuba Perlin, Stephen Gou, Nick Frosst, Jeff Dean, and Yarin Gal.
 Interlocking backpropagation: improving depthwise model-parallelism. J. Mach. Learn. Res., 23 (1), jan 2022. ISSN 1532-4435.

546

547

548

579

588

589

- Zhouyuan Huo, Bin Gu, qian Yang, and Heng Huang. Decoupled parallel backpropagation with convergence guarantee. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 2098–2106. PMLR, 10–15 Jul 2018. URL https://proceedings.mlr.
 press/v80/huo18a.html.
 - Alexander Isenko, Ruben Mayer, Jeffrey Jedele, and Hans-Arno Jacobsen. Where is my training bottleneck? hidden trade-offs in deep learning preprocessing pipelines. In *Proceedings of the* 2022 International Conference on Management of Data, pp. 1825–1839, 2022.
- Max Jaderberg, Wojciech Marian Czarnecki, Simon Osindero, Oriol Vinyals, Alex Graves, David
 Silver, and Koray Kavukcuoglu. Decoupled Neural Interfaces using Synthetic Gradients.
 arXiv:1608.05343 [cs], August 2016. URL http://arxiv.org/abs/1608.05343.
- David Kappel, Khaleelulla Khan Nazeer, Cabrel Teguemne Fokam, Christian Mayr, and Anand Subramoney. Block-local learning with probabilistic latent representations, 2023.
- Janis Keuper and Franz-Josef Preundt. Distributed training of deep neural networks: Theoretical and practical limits of parallel scalability. In 2016 2nd workshop on machine learning in HPC environments (MLHPC), pp. 19–26. IEEE, 2016.
- Adarsh Kumar, Kausik Subramanian, Shivaram Venkataraman, and Aditya Akella. Doing more by doing less: how structured partial backpropagation improves deep learning clusters. In *Proceedings of the 2nd ACM International Workshop on Distributed Machine Learning*, pp. 15–21, 2021.
- Vyacheslav Kungurtsev, Malcolm Egan, Bapi Chatterjee, and Dan Alistarh. Asynchronous optimization methods for efficient training of deep neural networks with guarantees. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(9):8209–8216, May 2021. doi: 10. 1609/aaai.v35i9.16999. URL https://ojs.aaai.org/index.php/AAAI/article/ view/16999.
- Guillaume Leclerc, Andrew Ilyas, Logan Engstrom, Sung Min Park, Hadi Salman, and Alek sander Madry. Ffcv: Accelerating training by removing data bottlenecks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12011–12020, 2023.
- 571 Chenxiang Ma, Jibin Wu, Chenyang Si, and Kay Chen Tan. Scaling supervised local learning with
 572 augmented auxiliary networks. *arXiv preprint arXiv:2402.17318*, 2024.
- Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In Dekang Lin, Yuji Matsumoto, and Rada Mihalcea (eds.), *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pp. 142–150, Portland, Oregon, USA, June 2011. Association for Computational Linguistics. URL https://aclanthology.org/ P11–1015.
- Konstantin Mishchenko, Francis Bach, Mathieu Even, and Blake Woodworth. Asynchronous SGD
 Beats Minibatch SGD Under Arbitrary Delays, June 2022. URL http://arxiv.org/abs/
 2206.07638.
- Giorgi Nadiradze, Ilia Markov, Bapi Chatterjee, Vyacheslav Kungurtsev, and Dan Alistarh. Elastic
 Consistency: A Practical Consistency Model for Distributed Stochastic Gradient Descent. Proceedings of the AAAI Conference on Artificial Intelligence, 35(10):9037–9045, May 2021. ISSN 2374-3468. doi: 10.1609/aaai.v35i10.17092. URL https://ojs.aaai.org/index.
 php/AAAI/article/view/17092.
 - Arild Nøkland and Lars Hiller Eidnes. Training neural networks with local error signals. In *Inter*national conference on machine learning, pp. 4839–4850. PMLR, 2019.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward
 Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep

594 595 596	<pre>learning library. In Advances in Neural Information Processing Systems 32, pp. 8024-8035. Cur- ran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/ paper/2019/file/bdbca288fee7f92f2bfa9f7012727740-Paper.pdf.</pre>
597 598 599 600 601 602	Benjamin Recht, Christopher Re, Stephen Wright, and Feng Niu. Hogwild!: A lock-free approach to parallelizing stochastic gradient descent. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K.Q. Weinberger (eds.), <i>Advances in Neural Information Processing Systems</i> , volume 24. Curran Associates, Inc., 2011. URL https://proceedings.neurips.cc/paper_files/paper/2011/file/218a0aefdldla4be65601cc6ddc1520e-Paper.pdf.
603 604 605	Wei Wen, Cong Xu, Feng Yan, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Terngrad: Ternary gradients to reduce communication in distributed deep learning, 2017. URL https: //arxiv.org/abs/1705.07878.
606 607 608	Paul Werbos. Applications of advances in nonlinear sensitivity analysis. <i>System Modeling and Optimization</i> , pp. 762–770, 1982.
609 610 611 612 613 614	Hao Zhang, Zeyu Zheng, Shizhen Xu, Wei Dai, Qirong Ho, Xiaodan Liang, Zhiting Hu, Jin- liang Wei, Pengtao Xie, and Eric P. Xing. Poseidon: An efficient communication archi- tecture for distributed deep learning on GPU clusters. In 2017 USENIX Annual Technical Conference (USENIX ATC 17), pp. 181–193, Santa Clara, CA, July 2017. USENIX Asso- ciation. ISBN 978-1-931971-38-6. URL https://www.usenix.org/conference/ atc17/technical-sessions/presentation/zhang.
615 616 617 618	Shuxin Zheng, Qi Meng, Taifeng Wang, Wei Chen, Nenghai Yu, Zhi-Ming Ma, and Tie-Yan Liu. Asynchronous stochastic gradient descent with delay compensation. In <i>International conference on machine learning</i> , pp. 4120–4129. PMLR, 2017.
619 620 621 622	
623 624	
625 626 627	
628 629 630	
631 632 633	
634 635 636	
637 638 639	
640 641 642	
643 644	
645 646 647	

А ADDITIONAL RESULTS

A.1 LEARNING CURVES

Here, we provide additional details to the results provided in the main text. Figures 3 and 4 show the learning dynamics of Asynchrounous Backpropagation with blocks updates (Async BU) and with layer-wise updates (Async LU) on CIFAR10 and IMDb respectively. The difference in convergence speed and accuracy observed with CIFAR100 2 is less noticeable on CIFAR10, probably because it is a simpler task. However, we clearly see the advantage of Async LU on the IMDb, where it not only converges faster but also to similar accuracy.

Figures 5 and 6 compare the training curses of sequential SGD with Async LU respectively. We observed that Async LU needs more epochs to converge, ~ 15 epochs more. When trained to a larger dataset (Figures 7 and 8), we can see that both Async LU and SGD seem to converge within the same number of epochs while Async LU scales almost linearly with the number of GPUs. We should take these results carefully given that accuracies are plots against the number of epochs and not the time since both are trained on different numbers of GPUs. Appendix A.2 gives a comparison with equal number of GPUs.



Figure 3: testing curves of Asynchronous SGD with layer-wise updates (Async LU) and Block updates (Async BU) for ResNet18 (top plots) and ResNet50 (bottom plots) on the CIFAR10 dataset.

A.2 SPEED-UP COMPARISON WITH MULTI-GPU BACKPROPAGATION

Here we do a comparison of Asynchronous Backpropagation with layer-wise updates (Async LU) and multi-GPU Data Distributed Parallel(DDP) both trained on 3 GPUs residing on the same ma-chine (described in section 4) to achieve their accuracies. Since Async LU uses only one forward pass, we set its batch size to be 128 and that of BP to $3 \times$ higher (384). Async LU was imple-mented on the c++ library of Pytorch, Libtorch, while multi-GPU SGD is implemented using Pytorch DataDistributedParallel (DDP) API. The hyperparameters used are the same as described in section 4.

To make the comparison fair across implementation platforms, the relative speed-up is calculated with respect to their single GPU implementation respectively.



Figure 4: testing curves of Asynchronous SGD with layer-wise updates (Async LU) and Block updates (Async BU) on the IMDb dataset.



Figure 5: training curves of Asynchronous SGD with layer-wise updates (Async LU) on 3 GPUs and Single GPU SGD on CIFAR-10 dataset on ResNet18 (left) and ResNet50 (right).

In this Settings, DDP should clearly be at advantage since it uses a bigger batch size, all the GPUs are on the same machine and we optimized DDP loading process by persisting the data on the GPUs, reducing considerably the communication bottleneck, hence making the synchronization step faster. What we however observe is that Async LU achieves comparable relative speed-up over a single GPU compared to DDP on both CIFAR10 and CIFAR100. This shows the effectiveness of our asynchronous formulation (figure 1). We can expect Async LU to have greater advantage in a multinode or heterogeneous setting because the synchronization barrier becomes problem.

Table 7: Comparison of Async LU and DDP both trained on 3GPUs based on their relative speed-up to single GPU implemention for 3 runs on CIFAR10

Network architecture	Training method	Relative speed-up mean \pm std
ResNet-18	Async LU DDP	$\begin{array}{c} 1.86 \pm 0.21 \\ 1.90 \pm 0.08 \end{array}$
ResNet-50	Async LU DDP	$\begin{array}{c} 2.02 \pm 0.10 \\ 2.38 \pm 0.19 \end{array}$

A.3 TIME MEASUREMENTS

Here we provide the results of a small-scale experiment on the timing measurement of forward and backward passes for CIFAR-100 with batch size 128 in table 9. As expected, a single backward pass requires $\sim 2 \times$ than that of a single forward pass. Extensive experiments on this is provided by Kumar et al. (2021)

Figure 6: training curves of Asynchronous SGD with layer-wise updates (Async LU) on 3 GPUs and Single GPU SGD on the CIFAR-100 dataset on ResNet18 (left) and ResNet50 (right).

Figure 7: training curves of Asynchronous SGD with layer-wise updates (Async LU) on 3 GPUs and Single GPU SGD on the ImageNet dataset.

A.4 HYPERPARAMETERS FOR THE EXPERIMENTS

Hyperparameters used in training experiments presented in section 4 are documented in table 10

B CONVERGENCE PROOF

Here we provide the proof that the stochastic parameter dynamics, Eq. (5) of the main text, converges to a stationary distribution $p^*(\theta)$ given by

$$p^*(\boldsymbol{\theta}) = \frac{1}{\mathcal{Z}} \exp\left(\sum_k h_k(\boldsymbol{\theta})\right), \text{ with } h_k(\boldsymbol{\theta}) = \int \frac{\mu_k(\boldsymbol{\theta})}{D_k(\boldsymbol{\theta})} d\boldsymbol{\theta} - \ln|D_k(\boldsymbol{\theta})| + C.$$
 (7)

798 799

796 797

768

769

782

783

784 785 786

787

788 789 790

791 792

The proof is analogous to the derivation given in Bellec et al. (2017), and relies on stochastic calculus to determine the parameter dynamics in the infinite time limit. Since the dynamics include a noise term, the exact value of the parameters $\theta(t)$ at a particular point in time t > 0 cannot be determined, but we can describe the distribution of parameters using the Fokker-Planck formalism, i.e. we describe the parameter distribution at time t by a time-varying function $p_{\rm FP}(\theta, t)$.

To arrive at an analytical solution for the stationary distribution, $p^*(\theta)$ we make the adiabatic assumption that noise in the parameters only has local effects, such that the diffusion due to noise in any parameter θ_j has negligible influence on dynamics in θ_k , i.e. $\frac{\partial}{\partial \theta_j} D_k(\theta) = 0, \forall j \neq k$. Using this assumption, it can be shown that, for the dynamics (6), $p_{\rm FP}(\theta, t)$ converges to a unique stationary distribution in the limit of large t and small noise σ_{STALE} . To prove the convergence to the stationary distribution, we show that it is kept invariant by the set of SDEs Eq. (6) and that it can be

Figure 8: testing curves of Asynchronous SGD with layer-wise updates (Async LU) on 3 GPUs and Single GPU SGD on the ImageNet dataset.

Table 8: Comparison of Async LU and DDP both trained on 3GPUs based on their relative speed-up to single GPU implemention for 3 runs on CIFAR100

Network architecture	Training method	Relative speed-up mean \pm std
ResNet-18	Async LU DDP	$\frac{1.83 \pm 0.06}{1.86 \pm 0.21}$
ResNet-50	Async LU DDP	$\begin{array}{c} 2.02 \pm 0.13 \\ 2.30 \pm 0.13 \end{array}$

reached from any initial condition. Eq. 6 implies a Fokker-Planck equation given by

$$\frac{\partial}{\partial t} p_{\rm FP}(\boldsymbol{\theta}, t) = -\sum_{k} \frac{\partial}{\partial \theta_{k}} [\mu_{k}(\boldsymbol{\theta}, t) p_{\rm FP}(\boldsymbol{\theta}, t)] + \frac{\partial^{2}}{\partial \theta_{k}^{2}} [D_{k}(\boldsymbol{\theta}, t) p_{\rm FP}(\boldsymbol{\theta}, t)]$$
(8)

We show that, under the assumptions outlined above, the stochastic parameter dynamics Eq. (6) of the main text, converges to the stationary distribution $p^*(\theta)$ (Eq. (7)).

To arrive at this result, we plug in the assumed stationary distribution into Eq. (8) and show the equilibrium $\frac{\partial}{\partial t} p_{\text{FP}}(\boldsymbol{\theta}, t) = 0$, i.e.

 $\frac{\partial}{\partial t} p_{\rm FP}(\boldsymbol{\theta}, t) = -\sum_{k} \frac{\partial}{\partial \theta_{k}} [\mu_{k}(\boldsymbol{\theta}) p_{\rm FP}(\boldsymbol{\theta}, t)]$

$$\begin{split} &+ \frac{\partial^2}{\partial \theta_k^2} [D_k(\boldsymbol{\theta}) p_{\rm FP}(\boldsymbol{\theta}, t)] = 0 \\ \leftrightarrow &- \sum_k \frac{\partial}{\partial \theta_k} [\mu_k(\boldsymbol{\theta}) p_{\rm FP}(\boldsymbol{\theta}, t)] \end{split}$$

(9)

$$+ \frac{\partial}{\partial \theta_k} \left[\left(\frac{\partial}{\partial \theta_k} D_k(\boldsymbol{\theta}) \right) p_{\rm FP}(\boldsymbol{\theta}, t) \right]$$

$$+ \frac{\partial}{\partial \theta_k} \left[D_k(\boldsymbol{\theta}) \left(\frac{\partial}{\partial \theta_k} h_k(\boldsymbol{\theta}) \right) p_{\text{FP}}(\boldsymbol{\theta}, t) \right]$$

Network architecture

ResNet-18

ResNet-50

		14				
	CIFA	R-100	CIFA	R-10	Imagenet	IMDb
hyperparameter	Resnet-18	Resnet-50	Resnet-18	Resnet-50	Resnet-50	LSTM
batch_size	128	128	128	128	250	75
lr	0.015	0.015	0.015	0.015	0.015	0.001
momentum	0.9	0.9	0.9	0.9	0.9	0.9
T₋max	100	120	105	120	250	150
warm_up_epochs	5	5	3	5	5	0
warm_up lr	0.005	0.005	0.005	0.005	0.0035	-
weight_decay	0.005	0.005	0.005	0.005	0.0035	0

864	Table 9: Timing measurement of forward and backward passes for CIFAR-100 with batch size 128.
865	Averaged over all batches for 15 epochs.

Forward pass (s)

mean \pm std

 $0.0049 \pm 1E-04$

 $0.0166 \pm \text{5e-05}$

Backward pass (s)

mean \pm std

 $0.0102 \pm 1E-04$

 $0.0299 \pm \textbf{4E-05}$

where we used the simplifying assumption, $\frac{\partial}{\partial \theta_j} D_k(\boldsymbol{\theta}) = 0, \forall j \neq k$, as outlined above. Next, using $\frac{\partial}{\partial \theta_j} h_k(\boldsymbol{\theta}) = \frac{1}{D_k(\boldsymbol{\theta},t)} \left(\mu_k(\boldsymbol{\theta},t) - \frac{\partial}{\partial \theta_j} D_k(\boldsymbol{\theta},t) \right)$, we get

$$\frac{\partial}{\partial t} p_{\rm FP}(\boldsymbol{\theta}, t) = 0 \quad \leftrightarrow \quad -\sum_{k} \frac{\partial}{\partial \theta_{k}} [\mu_{k}(\boldsymbol{\theta}, t) p_{\rm FP}(\boldsymbol{\theta}, t)] \\
+ \frac{\partial}{\partial \theta_{k}} \left[\left(\frac{\partial}{\partial \theta_{k}} D_{k}(\boldsymbol{\theta}, t) \right) p_{\rm FP}(\boldsymbol{\theta}, t) \right] \\
+ \frac{\partial}{\partial \theta_{k}} \left[\left(\mu_{k}(\boldsymbol{\theta}, t) - \frac{\partial}{\partial \theta_{k}} D_{k}(\boldsymbol{\theta}, t) \right) p_{\rm FP}(\boldsymbol{\theta}, t) \right] = 0$$
(10)

This shows that the simplified dynamics, Eq. 6, leave the stationary distribution (7) unchanged. This stationary distribution $p^*(\theta)$ is a close approximation to SGD. To see this, we study the maxima of the distribution, by taking the derivative

$$\frac{\partial}{\partial \theta_k} h_k(\boldsymbol{\theta}) = \frac{\mu_k(\boldsymbol{\theta})}{D_k(\boldsymbol{\theta})} - \frac{\partial}{\partial \theta_k} \ln |D_k(\boldsymbol{\theta})| , \qquad (11)$$

which by inserting (6) can be written as

$$\frac{\partial}{\partial \theta_k} h_k(\boldsymbol{\theta}) = -\frac{1}{\eta} \frac{\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) + \sigma_{\text{STALE}}^2 \nabla_{\boldsymbol{\theta}}^3 \mathcal{L}(\boldsymbol{\theta})}{\sigma_{\text{SGD}}^2 + \sigma_{\text{STALE}}^2 \nabla_{\boldsymbol{\theta}}^2 \mathcal{L}(\boldsymbol{\theta})}$$
(12)

If σ_{STALE} is small compared to σ_{SGD} we recover the cannonical results for SGD $\frac{\partial}{\partial \theta_k} h_k(\theta) \approx -\frac{1}{\eta} \frac{\nabla_{\theta} \mathcal{L}(\theta)}{\sigma_{\text{SCD}}^2}$, where smaller learning rates η make the probability of reaching local optima more peaked. Distortion of local optima, which manifests in the second term in the nominator, only depend on third derivatives, which can be expected to be small for most neural network architectures with well-behaved non-linearities.