

Accelerating Asynchronous Stochastic Gradient Descent for Neural Machine Translation

Nikolay Bogoychev[★] Marcin Junczys-Dowmunt[✈] Kenneth Heafield[★] Alham Fikri Aji[★]

[★] School of Informatics, University of Edinburgh, Edinburgh, United Kingdom

[✈] Microsoft, 1 Microsoft Way, Redmond, WA 98121, USA

{n.bogoych, a.fikri, kheafield}@ed.ac.uk, marcinjd@microsoft.com

Abstract

In order to extract the best possible performance from asynchronous stochastic gradient descent one must increase the mini-batch size and scale the learning rate accordingly. In order to achieve further speedup we introduce a technique that delays gradient updates effectively increasing the mini-batch size. Unfortunately with the increase of mini-batch size we worsen the stale gradient problem in asynchronous stochastic gradient descent (SGD) which makes the model convergence poor. We introduce local optimizers which mitigate the stale gradient problem and together with fine tuning our momentum we are able to train a shallow machine translation system 27% faster than an optimized baseline with negligible penalty in BLEU.

1 Introduction

With training times measured in days, parallelizing stochastic gradient descent (SGD) is valuable for making experimental progress and scaling data sizes. Synchronous SGD sums gradients computed by multiple GPUs into one update, equivalent to a larger batch size. But GPUs sit idle unless workloads are balanced, which is difficult in machine translation and other natural language tasks because sentences have different lengths. Asynchronous SGD avoids waiting, which is faster in terms of words processed per second. However asynchronous SGD suffers from stale gradients (Abadi et al., 2016) that degrade convergence, resulting in an almost no improvement in time to convergence (Hadjis et al., 2016). This paper makes asynchronous SGD even faster and deploys a series of convergence optimizations.

In order to achieve fastest training (and inspired by Goyal et al. (2017) we increase the mini-batch size, making the matrix operations more efficient and reducing the frequency of gradient communica-

tion for the optimizer step. Unlike their task (image classification), text training consumes a lot of GPU Memory (Table 1) for word embedding activations making it impossible to fit mini-batches of similar magnitude as Goyal et al. (2017).

Our main contributions are as follows:

1. We introduce a delayed gradient updates which allow us to work with much larger mini-batches which would otherwise not be possible due to limited GPU memory.
2. We introduce local optimizers which run on each worker to mitigate the extra staleness and convergence issues (Dekel et al., 2010; Keskar et al., 2017) caused by large mini-batches.
3. We highlight the importance of tuning the optimizer momentum and show how it can be used as a cooldown strategy.

VRAM	τ	Words	WPS
3 GB	1	3080	19.5k
7 GB	1	7310	36.6k
10 GB	1	10448	40.2k
20* GB	2	20897	44.2k
30* GB	2	31345	46.0k
40* GB	4	41794	47.6k

Table 1: Relationship between the GPU Memory (VRAM) budget for batches (* means emulated by summing τ smaller batches), number of source words processed in each batch and words-per-second (WPS) measured on a shallow model.

2 Experiments

This section introduces each optimization along with an intrinsic experiment on the WMT 2016 Romanian→English task (Bojar et al., 2016).

The translation system is equivalent to Sennrich et al. (2016), which was the first place constrained system (and tied for first overall in the WMT16

shared task.). The model is a shallow bidirectional GRU (Bahdanau et al., 2014) encoder-decoder trained on 2.6 million parallel sentences. Due to variable-length sentences, machine translation systems commonly fix a memory budget then pack as many sentences as possible into a dynamically-sized batch. The memory allowance for mini-batches in our system is 3 GB (for an average batch size of 2633 words). Adam (Kingma and Ba, 2015) is used to perform asynchronous SGD with learning rate of 0.0001. This is our baseline system. We also compare with a synchronous baseline which uses modified Adam parameters, warmup of 16000 mini-batches and inverse square root cooldown following Vaswani et al. (2017). We used 4 Tesla P 100 GPUs in a single node with the Marian NMT framework for training (Junczys-Dowmunt et al., 2018). Since we apply optimizations over asynchronous SGD we performed a learning rate and mini-batch-size parameter sweep over the baseline system and settled on a learning rate of 0.00045 and 10 GB memory allowance for mini-batches (average batch size of 10449 words). This is the fastest system we could train without sacrificing performance before adding our improvements. In our experiments on Table 2 we refer to this system as "Optimized asynchronous". All systems were trained until 5 consecutive stalls in the cross-entropy metric of the validation set. Note that some systems require more epochs to reach this criteria which indicates poor model convergence.

2.1 Larger Batches and delayed updates

This experiment aims to increase speed, in words-per-second (WPS), by increasing the batch size. Larger batches have two well-known impacts on speed: making more use of GPU parallelism and communicating less often.

After raising the batch size to the maximum that fits on the GPU,¹ we emulate even larger batches by processing multiple mini-batches and summing their gradients locally without sending them to the optimizer. This still increases speed because communication is reduced (Table 1). We introduce parameter τ , which is the number of iterations a GPU performs locally before communicating externally as if it had run one large batch. The Words-

¹The Tesla P100 has 16 GB of GPU memory and we opt to use 10 GBs of mini-batches and the rest is used to store model parameters, shards, optimizers and additional system specific elements such as the cache vectors for gradient dropping (Aji and Heafield, 2017).

per-second (WPS) column on Table 1 shows the effect on corpora processing speed when applying delayed gradients updates for different values of τ . While we reduce the overall training time if we just apply delayed gradient updates we worsen the overall convergence (Table 2).

When increasing the mini-batch size τ times without touching the learning rate, we effectively do τ times less updates per epoch. On the surface, it might seem that these less frequent updates are counterbalanced by the fact that each update is accumulated over a larger batch. But practical optimization heuristics like gradient clipping mean that in effect we end up updating the model less often, resulting in slower convergence. Goyal et al. (2017) recommend scaling the learning rate linearly with the mini-batch size in order to maintain convergence speed.

2.2 Warmup

Goyal et al. (2017) point out that just increasing the learning rate performs poorly for very large batch sizes, because when the model is initialized at a random point, the training error is large. Large error and large learning rate result in bad "jerky" updates to the model and it can't recover from those. Goyal et al. (2017) suggest that initially model updates should be small so that the model will not be pushed in a suboptimal state. Afterwards we no longer need to be so careful with our updates.

2.2.1 Lowering initial learning rate

Goyal et al. (2017) lower the initial learning rate and gradually increase it over a number of mini-batches until it reaches a predefined maximum. This technique is also adopted in the work of Vaswani et al. (2017). This is the canonical way to perform warmup for neural network training.

2.2.2 Local optimizers

We propose an alternative warm up strategy and compare it with the canonical method. Since we emulate large batches by running multiple smaller batches, it makes sense to consider whether to optimize locally between each batch by adapting the concept of local per-worker optimizers from Zhang et al. (2014). In asynchronous SGD setting each GPU has a full copy of the model as well as the master copy of $1/N$ th of the parameters in its capacity as parameter server. We use the local optimizers to update the local model shard in between delayed gradient updates, which helps mitigate staleness.

Unlike prior work, we also update the shard of the global model that happens to be on the same GPU. Local updates are almost free because we avoid remote device communication.

Updating the parameter shard of the global model bears some resemblance to the Hogwild method (Recht et al., 2011) as we don't synchronize the updates to the shard, however, global updates are still synchronised. As before, once every τ iterations we run a global optimizer that updates the sharded parameter set and then distributes the updated model across all devices. Any local model divergences are lost at this point. We found that this strategy improves model convergence in the early epochs but tends to be harmful later on. We hypothesize that initially partial model updates reduce staleness, but when the model starts to converge, local optimizers introduce extra noise in the training, which is harmful. We use local optimizers purely as a warmup strategy, turning them off after the initial phase of the training. Empirically, we found that we can get the best convergence by using them for the first 4000 mini-batches that each device sees. On Table 2 we compare and contrast the two warmup strategies. By itself learning-rate warmup offers slower convergence but to a better point compared to local optimizers. The reader may notice that if we apply delayed gradient updates, the effective batch size that the global optimizer deals with is τ times larger than the mini-batch size on which the local optimizers runs. Therefore we use τ times lower learning rate for the local optimizers compared to the global optimizers.

2.3 Momentum cooldown and tuning

Goyal et al. (2017) and Vaswani et al. (2017) both employ cooldown strategies that lower the learning rate towards the end of training. Inspired by the work of Hadjis et al. (2016) however we decided to pursue a different cooldown strategy by modifying the momentum inside Adam's parameters.

Momentum tuning is not a well explored area in deep learning. Most researchers simply use the default values for momentum for a chosen optimizer (Hadjis et al., 2016) (in the case of NMT, this is usually Adam). Hadjis et al. (2016) argue that this is an oversight especially when it comes to asynchronous SGD, because the asynchronicity adds extra implicit momentum to the training which is not accounted for. Because of this, asynchronous SGD has been deemed ineffective, as without mo-

mentum tuning, the observed increase in training speed is negated by the lower convergence rate, resulting in near-zero net gain (Abadi et al., 2016). However, Hadjis et al. (2016) show that after performing a grid search over momentum values, it is possible to achieve convergence rates typical for synchronous SGD even when working with many asynchronous workers. The downside of momentum tuning is that we can't offer rule-of-thumb values, as they are individually dependent on the optimizer used, the neural model, the number of workers and the batch size. In our experiments, we lowered the overall momentum and in addition performed momentum cooldown where we reduced the momentum of our optimizer (Adam) after the first few thousand batches.

2.4 Results

Table 2 shows the effect of modifying momentum values. When using just delayed gradient updates, training is noticeably faster, but there are significant regressions in BLEU and CE (system 2). In order to mitigate those, when using delayed gradient updates, we tune the momentum and apply momentum cooldown on top of either of our warmup strategies. By doing this we not only further reduce training time, but also recover the loss of accuracy. Compared to the optimized baseline system (1), our best system (4) reduces the training time by 27%. Progression of the training can be seen on figures 1 and 2. Our system starts poorly compared to the baselines in terms of epoch-for-epoch convergence, but catches up in the later epochs. Due to faster training speed however, the desired BLEU score is achieved faster (Figure 2).

Local optimizers as a warmup strategy show faster convergence compared to learning rate warmup at almost no penalty to BLEU or cross-entropy (System 4 vs system 6). Against the system used in WMT 16 (Sennrich et al., 2016), we achieve nearly 4 times faster training time with no discernible penalty in BLEU or CE. In contrast, the other communication reducing method tested, the work of Aji and Heafield (2017), is slower than our work and achieves worse BLEU and CE.

2.5 Using even larger mini-batches

We can achieve even greater processing speed by further increasing τ but we were unable to maintain the same convergence with the Romanian-English shallow model. We found that larger τ values are useful when dealing with the larger deep RNN mod-

System	Time (hours)	Epochs	BLEU	CE	WPS
synchronous	14.3	11	35.3	50.63	15.7k
asynchronous (0)	12.2	13	35.61	50.47	19.5k
(0) + Aji and Heafield (2017)	6.23	12	35.16	50.86	26.9k
optimized asynchronous (1)	4.97	10	35.56	50.90	40.2k
(1) + Aji and Heafield (2017)	4.32	11	35.16	52.02	41.5k
(1) + delayed updates $\tau = 2$ (2)	4.20	11	34.82	51.68	44.2k
(2) + local optimizers (3)	3.66	10	35.45	51.32	44.2k
(3) + momentum tuning (4)	3.66	10	35.48	50.87	44.2k
(2) + warmup (5)	4.87	13	35.29	50.78	44.2k
(5) + momentum tuning (6)	3.98	11	35.76	50.73	44.2k

Table 2: Romanian-English results from our exploration and optimizations. We also compare our methods against the work of Aji and Heafield (2017) which also reduces communication. We use system (1) as our reference baseline upon which we improve. The system that achieved the best training time is bolded.

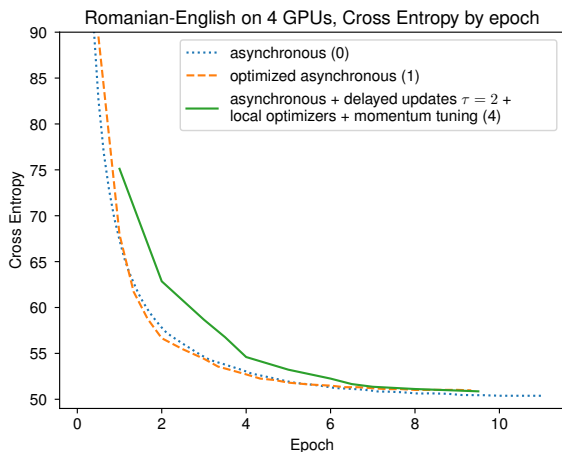


Figure 1: Cross-entropy training progression per epoch for our ro-en systems.

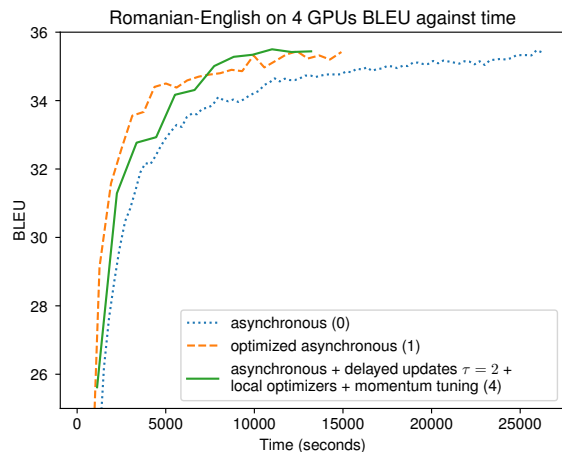


Figure 2: BLEU scores for our ro-en systems.

els. With deep RNN models the parameters take the majority of the available VRAM leaving very little for mini-batches. In this scenario we can apply $\tau = 4$ without negative effect towards convergence. We demonstrate the effectiveness of larger τ on Table 3. The baseline system is equivalent to the winning system for English-German at the WMT 2017 competition (Sennrich et al., 2017). The baseline is trained with synchronous SGD and our system uses asynchronous SGD, delayed gradient updates by a factor of 4, local optimizers and the momentum is tuned and further reduced after the first 16000 mini-batches. We found learning rate of 0.0007 to work the best. We do not report the numbers for asynchronous baseline because we were unable to achieve competitive BLEU scores without using delayed gradient updates. We speculate this is because with this type of deep model, our mini-batch

size is very small leading to very jerky and unstable training updates. Larger mini-batches ensure the gradients produced by different workers are going to be closer to one another. Our training progression can be seen on figures 3 and 4. We show that even though we use 4 times larger mini-batches we actually manage to get lower Cross-Entropy epoch for epoch compared to the baseline (Figure 3). This coupled with our higher training speed makes our method reach the best BLEU score 1.6 times faster than the baseline (Figure 4).

3 Related work

We use larger mini-batches and delay gradient updates in order to increase the speed at which the dataset is processed. The principal reason why this works is because when mini-batch size is increased n (also includes delayed updates) times, communication is reduced by the same amount. This as-

System	Time (h)	BLEU	CE
Baseline	51.3	25.1	47.31
Async (4) + $\tau = 4$	39.7	25.07	46.59

Table 3: Training times for English-German deep RNN system trained on WMT17 data. Our asynchronous system includes the optimizations of system (4) from Table 2.

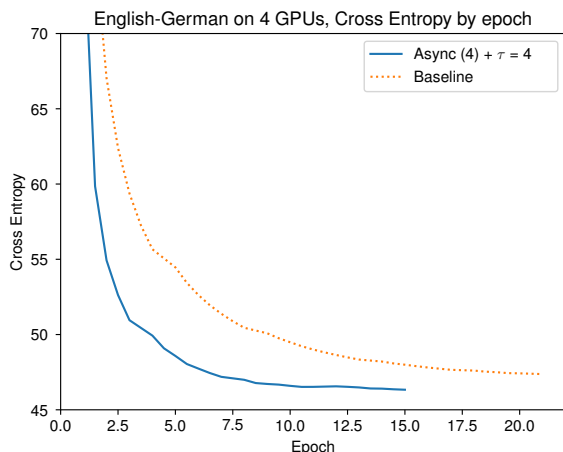


Figure 3: CE scores for our en-de systems.

pect of our work is similar to the work of [Aji and Heafield \(2017\)](#) where they drop the lower 99% of the gradient updates based on absolute value thus reducing the memory traffic. Compared with them we achieve faster dataset processing speed and also better model convergence as shown on Table 2.

Independently from us [Mao et al. \(2018\)](#) extend the work of [Aji and Heafield \(2017\)](#) aiming to reduce gradient communication without suffering any of the negative effects we have noted. In process they independently arrive to some of the methods that we use, notably tuning the momentum and applying warmup to achieve better convergence.

Independently from us [Shazeer and Stern \(2018\)](#) have done further exploratory work on ADAM’s momentum parameters using the Transformer model ([Vaswani et al., 2017](#)) as a case study and have offered a mathematical explanation about why different stages of the training require different momentum values.²

Independently from us [Saunders et al. \(2018\)](#) have employed delayed gradient updates in syntax NMT setting, where the sequences are much longer due to the syntax annotation and delayed updates are necessary because video RAM is limited.³

²This work was published on 11.04.2018.

³This work was published on 01.05.2018.

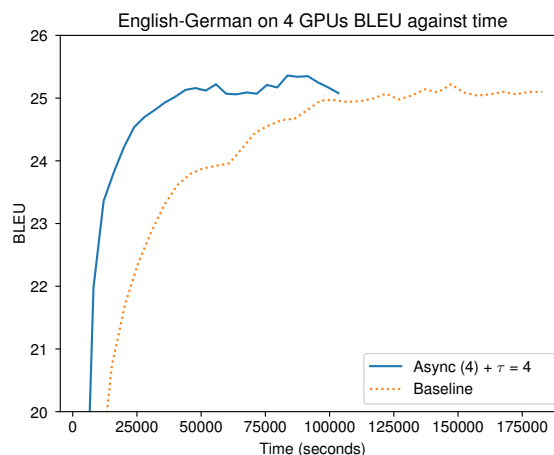


Figure 4: BLEU scores for our en-de systems.

Independently from us [Lin et al. \(2018\)](#) have developed their own local optimizer solution as an alternative to increasing mini-batch sizes.⁴

4 Conclusion and Future work

We show that we can increase speed and maintain convergence rate for very large mini-batch asynchronous SGD by carefully adjusting momentum and applying warmup and cooldown strategies. While we have demonstrated our methods on GPUs, they are hardware agnostic and can be applied to neural network training on any multi-device hardware such as TPUs or Xeon Phis. We were able to achieve end-to-end training on multiple tasks a lot faster than the baseline systems. For our Romanian-English model, we train nearly 3X faster than the commonly used baseline and 1.5X faster over a specifically optimised baseline. When experimenting with English-German we are able to train our model 1.3X faster than the baseline model, achieving practically the same BLEU score and much better model cross-entropy.

In the future we would like to apply local optimizers in distributed setting where the communication latency between local and remote devices varies significantly we could use local optimizers to synchronize remote models less often.

Acknowledgments

We thank Adam Lopez, Sameer Bansal and Naomi Saphra for their help and comments on the paper. We thank our reviewers for their comments and suggestions. Nikolay Bogoychev was funded by an Amazon faculty research award to Adam Lopez.

⁴This work was published on 22.08.2018

References

- Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek Gordon Murray, Benoit Steiner, Paul A. Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zhang. 2016. Tensorflow: A system for large-scale machine learning. *CoRR*, abs/1605.08695.
- Alham Fikri Aji and Kenneth Heafield. 2017. Sparse communication for distributed gradient descent. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, Copenhagen, Denmark.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *CoRR*, abs/1409.0473.
- Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Varvara Logacheva, Christof Monz, Matteo Negri, Aurelie Neveol, Mariana Neves, Martin Popel, Matt Post, Raphael Rubino, Carolina Scarton, Lucia Specia, Marco Turchi, Karin Verspoor, and Marcos Zampieri. 2016. Findings of the 2016 conference on machine translation. In *Proceedings of the First Conference on Machine Translation*, pages 131–198, Berlin, Germany. Association for Computational Linguistics.
- Ofer Dekel, Ran Gilad-Bachrach, Ohad Shamir, and Lin Xiao. 2010. Optimal distributed online prediction using mini-batches. Technical report.
- Priya Goyal, Piotr Dollár, Ross B. Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. 2017. Accurate, large minibatch SGD: training imagenet in 1 hour. *CoRR*, abs/1706.02677.
- S. Hadjis, C. Zhang, I. Mitliagkas, and C. Re. 2016. Omnivore: An optimizer for multi-device deep learning on cpus and gpus.
- Marcin Junczys-Dowmunt, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, André F. T. Martins, and Alexandra Birch. 2018. Marian: Fast neural machine translation in C++. In *Proceedings of ACL 2018, System Demonstrations*, pages 116–121, Melbourne, Australia. Association for Computational Linguistics.
- Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. 2017. On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima. In *Proceedings of ICLR 2017*, ICLR '17.
- Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *ICLR*.
- Tao Lin, Sebastian U. Stich, and Martin Jaggi. 2018. Don't use large mini-batches, use local sgd. *CoRR*, abs/1808.07217.
- Huizi Mao, Song Han, Wang Yu, William Dally, and Yujun Lin. 2018. Deep gradient compression: Reducing the communication bandwidth for distributed training.
- Benjamin Recht, Christopher Re, Stephen Wright, and Feng Niu. 2011. Hogwild: A lock-free approach to parallelizing stochastic gradient descent. In J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett, F. Pereira, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 24*, pages 693–701. Curran Associates, Inc.
- Danielle Saunders, Felix Stahlberg, Adrià de Gispert, and Bill Byrne. 2018. Multi-representation ensembles and delayed sgd updates improve syntax-based nmt. *CoRR*, abs/1805.00456.
- Rico Sennrich, Alexandra Birch, Anna Currey, Ulrich Germann, Barry Haddow, Kenneth Heafield, Antonio Valerio Miceli Barone, and Philip Williams. 2017. The University of Edinburgh's Neural MT Systems for WMT17. In *Proceedings of the Second Conference on Machine Translation, Volume 2: Shared Task Papers*, Copenhagen, Denmark.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Edinburgh Neural Machine Translation Systems for WMT 16. In *Proceedings of the First Conference on Machine Translation*, pages 371–376, Berlin, Germany. Association for Computational Linguistics.
- Noam Shazeer and Mitchell Stern. 2018. Adafactor: Adaptive learning rates with sublinear memory cost. *CoRR*, abs/1804.04235.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need.
- Sixin Zhang, Anna Choromanska, and Yann LeCun. 2014. Deep learning with elastic averaging SGD. *CoRR*, abs/1412.6651.