

Dion: A Communication-Efficient Optimizer for Large Models

Kwangjun Ahn

Byron Xu

Microsoft Research

Technical Report
(April 8, 2025)

Abstract

Training large AI models efficiently requires distributing computation across multiple accelerators, but this often incurs significant communication overhead—especially during gradient synchronization. We introduce Dion, a communication-efficient optimizer that retains the synchronous semantics of standard distributed training (e.g., DDP, FSDP) while substantially reducing I/O costs. Unlike conventional optimizers that synchronize full gradient matrices, Dion leverages orthonormalized updates with device-local momentum buffers, eliminating the need for full gradient exchange. It further supports an efficient sharding strategy that avoids reconstructing large matrices during training.

Optimizer	Data Parallel I/O	Weight Sharding I/O	Optimizer State Memory
Adam	mn	0	$2mn$
Muon	mn	mn	mn
Dion	$(m + n)r$	$(m + 1)r$	$mn + nr$

Table 1: I/O comparison when the parameters shape is $m \times n$ ($m \leq n$). Dion uses a low-rank factor r , which can be much smaller than m or n .

1 Introduction

Efficient training of large AI models requires multiple accelerators (e.g., GPUs, TPUs) to achieve practical runtimes. Distributed training techniques, such as Distributed Data Parallel (DDP) [Li et al., 2020] and Fully Sharded Data Parallel (FSDP) [Zhao et al., 2023] improve scalability by partitioning both the data and model parameters. However, the communication overhead associated with these methods increases significantly with model size.

In practice, gradient synchronization typically depends on high-speed interconnects like InfiniBand, which necessitate costly, tightly coupled infrastructure—often confining accelerators to a single data center. Reducing the volume of synchronized data can help alleviate these constraints, enabling more flexible and cost-effective deployment across different environments.

A main source of this overhead lies in the design of modern optimizers. For instance, using Adam [Kingma and Ba, 2014] to optimize a parameter matrix in $\mathbb{R}^{m \times n}$ requires synchronizing the entire gradient matrix, incurring an $O(mn)$ communication cost. This burden is shared by many practical optimizers that rely on dense, synchronized updates across all devices.

To address communication costs, prior work has primarily focused on two strategies: (i) compressing gradients through quantization or sparsification [Wang et al., 2023], and (ii) reducing synchronization frequency via techniques like federated averaging or local SGD [McMahan et al., 2017]. In contrast, we take a complementary approach. Rather than lowering the precision or frequency of synchronization, we propose a novel **synchronous update rule** that is inherently more communication-efficient. Importantly, our method retains the step-level synchronization behavior of DDP and FSDP, ensuring compatibility with standard distributed training pipelines.

We introduce **Dion** (short for the **d**istributed **o**rthonormalization), a communication-efficient optimizer that reduces I/O costs compared to existing methods (see Table 1), while achieving training speeds on par with state-of-the-art

Algorithm 1 Centralized Dion (see Algorithm 2 for distributed implementation)

Require: Learning rate η , momentum decay μ , rank r .

Parameter $X \in \mathbb{R}^{m \times n}$, momentum $M \in \mathbb{R}^{m \times n}$ ($M_0 = 0$), right factor $Q^{n \times r}$ (randomly initialized).

```
1: for  $t = 1$  to  $T$  do
2:   Compute gradient  $G_t \in \mathbb{R}^{m \times n}$ 
3:    $B_t \leftarrow M_{t-1} + G_t \in \mathbb{R}^{m \times n}$ 
4:    $P_t, R_t \leftarrow \text{PowerIter1}(B_t; Q_{t-1}) \in \mathbb{R}^{m \times r}, \mathbb{R}^{n \times r}$  // rank- $r$  approximate
5:    $M_t \leftarrow B_t - (1 - \mu)P_t R_t^\top \in \mathbb{R}^{m \times n}$  // error feedback
6:    $Q_t \leftarrow \text{ColumnNormalize}(R_t) \in \mathbb{R}^{n \times r}$ 
7:    $X_t \leftarrow X_{t-1} - \eta P_t Q_t^\top \in \mathbb{R}^{m \times n}$  // orthonormal update
8: end for

PowerIter1( $B; Q$ ): // single rank- $r$  power iteration (initialized from  $Q$ )
   $P \leftarrow BQ \in \mathbb{R}^{m \times r}$ 
   $P \leftarrow \text{Orthogonalize}(P) \in \mathbb{R}^{m \times r}$  // using QR decomposition
   $R \leftarrow B^\top P \in \mathbb{R}^{n \times r}$ 
  return  $P, R$ 
```

optimizers. Dion adopts an orthonormalized update rule inspired by Muon, but modifies it to better support distributed training, unlike the original Newton-Schulz approach [Jordan et al., 2024b]. A key feature of Dion is that it allows momentum buffers to evolve independently on each device—an intentional divergence that eliminates the need for full gradient synchronization and enables communication savings. Furthermore, as shown in Algorithm 2, Dion employs an efficient weight sharding strategy that avoids reconstructing the full $m \times n$ matrix required for orthonormalization, unlike Muon.

These I/O savings are especially valuable in hybrid sharding setups for large-scale, geographically distributed training, where higher-level data parallelism (DP) is layered on top of Fully Sharded Data Parallel (FSDP) [Zhao et al., 2023]. In this configuration, each FSDP instance within a DP group holds a full copy of the model and processes a different data shard. As a result, gradient synchronization is required across FSDP instances—an increasingly costly operation when these instances span multiple clusters connected by slower inter-datacenter links.

2 Our Proposed Optimizer

We begin by introducing our main algorithm. For clarity, we describe it in terms of an $m \times n$ parameter matrix, denoted by X . At each step t , let X_t represent the current parameter, G_t the gradient, and M_t the momentum matrix.

We first present the centralized version of Dion in Algorithm 1, along with a brief intuition behind its design. We then introduce the distributed version in Algorithm 2.

2.1 Design Intuition Behind Dion

Our starting point is the Muon optimizer [Jordan et al., 2024b], which introduces a surprisingly effective update rule for matrix parameters using zeroth-power orthonormal updates. Muon maintains a momentum matrix that accumulates gradients over time:

$$M_t \leftarrow \mu M_{t-1} + G_t, \quad \text{for } \mu \in (0, 1),$$

and then applies an update based on the zeroth power of this matrix. Specifically, it computes

$$U_t V_t^\top \quad \text{where} \quad M_t = U_t \Sigma_t V_t^\top$$

is the singular value decomposition (SVD) of the momentum matrix. To approximate this zeroth-power orthonormalization efficiently, Jordan et al. [2024b] employ Newton-Schulz iterations. However, this approach is not well suited to distributed training, as it requires additional communication across the weight-parallel axis (see Table 1).

Algorithm 2 Dion (distributed implementation)

Require: Learning rate η , momentum decay μ , rank r .

- model/states are sharded as: $X^{(i)}, M^{(i)} \in \mathbb{R}^{m \times n_i}, Q^{(i)} \in \mathbb{R}^{n_i \times r}$, with $\sum_i n_i = n$.
- hat variables (e.g., $\hat{M}, \hat{G}, \hat{B}$) are local to each data-parallel worker and not synchronized.
- \mathbb{E}_{DP} denotes data-parallel all-reduce mean; \sum_i denotes weight-parallel all-reduce sum.

```
1: for  $t = 1$  to  $T$  do
2:   Compute local gradient  $\hat{G}_t^{(i)} \in \mathbb{R}^{m \times n_i}$ 
3:    $\hat{B}_t^{(i)} \leftarrow \hat{M}_{t-1}^{(i)} + \hat{G}_t^{(i)} \in \mathbb{R}^{m \times n_i}$ 
4:    $P_t, R_t^{(i)} \leftarrow \text{Distributed-PowerIter1}(\hat{B}_t^{(i)}; Q_{t-1}^{(i)}) \in \mathbb{R}^{m \times r}, \mathbb{R}^{n_i \times r}$ 
5:    $\hat{M}_t^{(i)} \leftarrow \hat{B}_t^{(i)} - (1 - \mu) P_t (R_t^{(i)})^\top \in \mathbb{R}^{m \times n_i}$ 
6:    $Q_t^{(i)} \leftarrow \text{Distributed-ColumnNormalize}(R_t^{(i)})$ 
7:    $X_t^{(i)} \leftarrow X_{t-1}^{(i)} - \eta P_t (Q_t^{(i)})^\top \in \mathbb{R}^{m \times n_i}$ 
8: end for
```

```
Distributed-PowerIter1( $\hat{B}^{(i)}; Q^{(i)}$ ): // distributed single rank- $r$  power iteration
   $\hat{P}^{(i)} \leftarrow \hat{B}^{(i)} Q^{(i)} \in \mathbb{R}^{m \times r}$ 
   $P \leftarrow \text{Orthogonalize}(\sum_i \mathbb{E}_{\text{DP}}[\hat{P}^{(i)}]) \in \mathbb{R}^{m \times r}$ 
   $\hat{R}^{(i)} \leftarrow (\hat{B}^{(i)})^\top P \in \mathbb{R}^{n_i \times r}$ 
   $R^{(i)} \leftarrow \mathbb{E}_{\text{DP}}[\hat{R}^{(i)}] \in \mathbb{R}^{n_i \times r}$ 
  return  $P, R^{(i)}$ 
```

```
Distributed-ColumnNormalize( $R^{(i)}$ ):
   $c^{(i)} \leftarrow \text{ColumnNorm}^2(R^{(i)}) \in \mathbb{R}^r$ 
   $c \leftarrow \sqrt{\sum_i c^{(i)}} \in \mathbb{R}^r$ 
   $Q^{(i)} \leftarrow R^{(i)} / c$  (column-wise)  $\in \mathbb{R}^{n_i \times r}$ 
  return  $Q^{(i)}$ 
```

To address the limitations of Muon, we propose using lower-rank updates, with rank $r \ll m, n$. This is achieved by efficiently computing a rank- r approximation of the momentum matrix using a single iteration of power iteration, followed by column normalization to ensure the update is orthonormal. The idea of using a single power iteration is inspired by [Vogels et al. \[2019\]](#).

However, as we show in [Section 4.1](#), a naive application of this approach fails. To make it effective, we introduce an **error feedback mechanism**, where the momentum buffer tracks the components not captured by the low-rank approximation. Specifically, we update the momentum as

$$M_t \leftarrow B_t - (1 - \mu) P_t R_t^\top,$$

where only the low-rank component is subtracted from the buffer. This mechanism is crucial for faster convergence; see [Figure 3a](#) for an ablation study demonstrating its importance.

Lastly, although omitted from the pseudocode for simplicity, all experiments include a scaling factor of $\sqrt{m/n}$ in the update. This normalization has been found to improve the transferability of hyperparameters across different model scales [\[Bernstein and Newhouse, 2024, Pethick et al., 2025\]](#).

2.2 Distributed Implementation

The key advantage of Dion lies in its efficient support for distributed training. As shown in [Algorithm 2](#), it remarkably allows for a fully synchronous implementation—matching the centralized update rule in [Algorithm 1](#)—without requiring synchronization of the gradient or momentum matrices across workers.

This property enables the reduced communication overhead summarized in [Table 1](#):

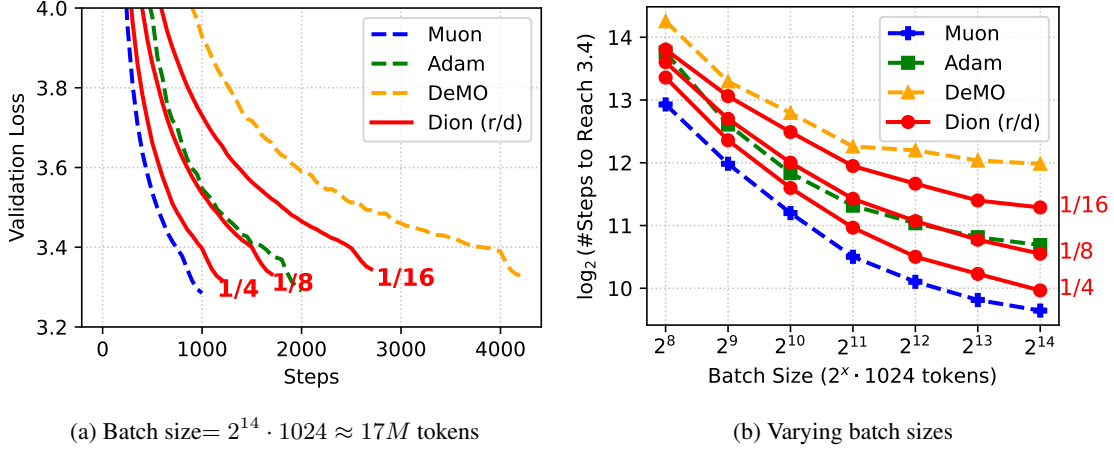


Figure 1: **Optimization performance for lower ranks.** Experimental results for Dion with $r = d/4, d/8, d/16$, where $d = 768$ is the hidden dimension of the GPT model. (a) Validation loss curves for various optimizers with batch size 2^{14} . (b) Number of steps required to reach a validation loss of 3.4 as a function of batch size.

- **Data-parallel axis:** Only the shared factors P and R need to be synchronized via an all-reduce mean, resulting in communication cost of $(m + n)r$ per step.
- **Weight-parallel axis:** Only the aggregated states $P^{(i)}$ and column norms $c^{(i)}$ are synchronized using an all-reduce sum, incurring a communication cost of $(m + 1)r$.
- **Memory usage:** Dion maintains two optimizer states—momentum M_t and the right factor Q_t —leading to total memory complexity of $mn + nr$.

The following result shows that this distributed implementation is functionally equivalent to its centralized counterpart.

Theorem 2.1. *The distributed implementation of Dion (Algorithm 2) is equivalent to the centralized version (Algorithm 1).*

Proof. See Appendix A for the proof. □

3 Experimental Results

In this section, we present preliminary experimental results using our new optimizer, Dion. We train a 120M-parameter GPT-style decoder-only Transformer model based on the nanoGPT codebase [Karpathy, 2023] on the FineWeb dataset [Penedo et al., 2024], following the setup of Jordan et al. [2024a].

Baselines. We compare our optimizer against AdamW [Loshchilov and Hutter, 2019], Muon [Jordan et al., 2024b], and DeMO [Peng et al., 2024]. For AdamW, we use $(\beta_1, \beta_2) = (0.9, 0.95)$ with a weight decay strength of 0.01. For Muon, we choose the default parameter of $\mu = 0.95$. DeMO is a recent optimizer that modifies momentum SGD by decoupling the momentum across distributed accelerators. Similar to Dion, it synchronizes only the fast component of the momentum extracted via the Discrete Cosine Transform (DCT), reducing communication overhead. We use the default hyperparameters from their official implementation and leave a more detailed comparison with DeMO to future work.

Training details. We evaluate the performance of various optimizers when training the Transformer layers. For the token embedding and language modeling head (which share weights), we consistently use Adam with a learning rate of 0.002, without tuning these values across optimizers. The sequence length is chosen as 1024.

For each experiment, we tune the optimizer’s learning rate based on the validation loss after the first 800 training steps. The candidate learning rates are as follows:

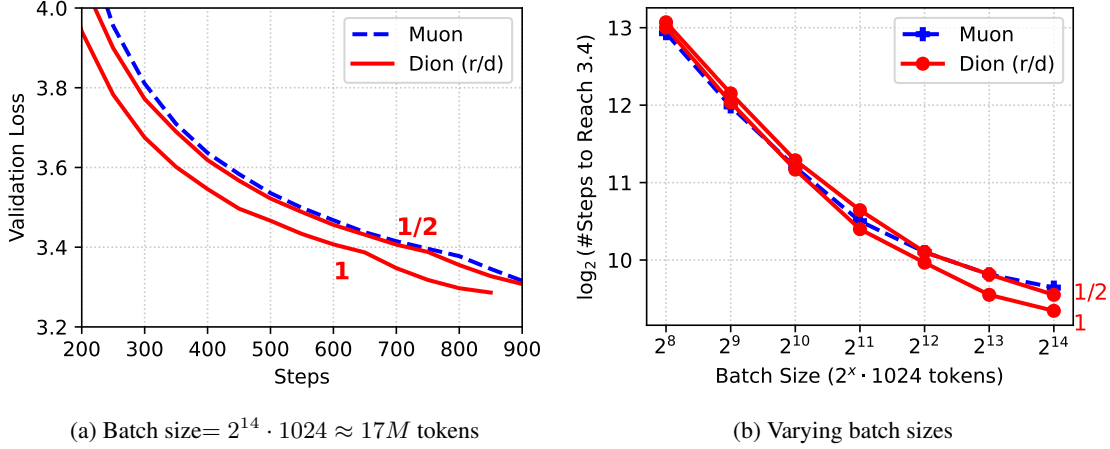


Figure 2: **Optimization performance for higher ranks.** Experimental results for Dion with $r = d, d/2$, where $d = 768$ is the hidden dimension of the GPT model. (a) Validation loss curves for various optimizers with batch size 2^{14} . (b) Number of steps required to reach a validation loss of 3.4 as a function of batch size. As batch size increases, Dion converges faster than Muon. A more detailed investigation is left for future work.

- for Adam and DeMO, we search over $\{0.00025, 0.0005, 0.001, 0.002, 0.004, 0.006, 0.008, 0.01, 0.02\}$,
- and for Muon and Dion, we search over $\{0.01, 0.02, 0.03, 0.04, 0.05\}$.

For Adam and DeMO, we apply learning rate warm-up during the first 100 iterations.

3.1 Optimization Performance

We evaluate the performance of each optimizer in Figure 1 and Figure 2. Figure 1 compares Dion at lower ranks $r = d/4, d/8, d/16$ against other baselines. As expected, Dion’s performance degrades with decreasing rank. However, this degradation is significantly more graceful than that of naive low-rank updates, thanks to the error feedback mechanism, as we shall see in Figure 3a. Remarkably, Dion remains competitive even at low ranks—for instance, at a batch size of 2^{14} , Dion with $r = d/8$ converges faster than Adam.

The curious case of full ranks. In Figure 2, we test Dion at higher ranks, $r = d$ and $r = d/2$. If both Muon and Dion produced perfectly orthonormal updates, their performance would coincide as r approaches d . However, the results in Figure 2 tell a more nuanced story.

We observe that Dion’s performance remains stable as batch size increases, while Muon degrades more noticeably. At batch size 2^{14} , Dion with $r = d/2$ even slightly outperforms Muon. We hypothesize that this gap may be due to Muon’s use of imperfect Newton–Schulz iterations, which could lead to worse updates compared to Dion.

That said, these findings are preliminary, and we leave a deeper investigation of this phenomenon to future work.

4 Ablating Key Design Choices

In this section, we conduct ablation studies to evaluate two core components of Dion: the error feedback mechanism and the use of a single power iteration for low-rank approximation. For the ablation studies, we choose the batch size to be $2048 \cdot 1024 \approx 2M$ tokens.

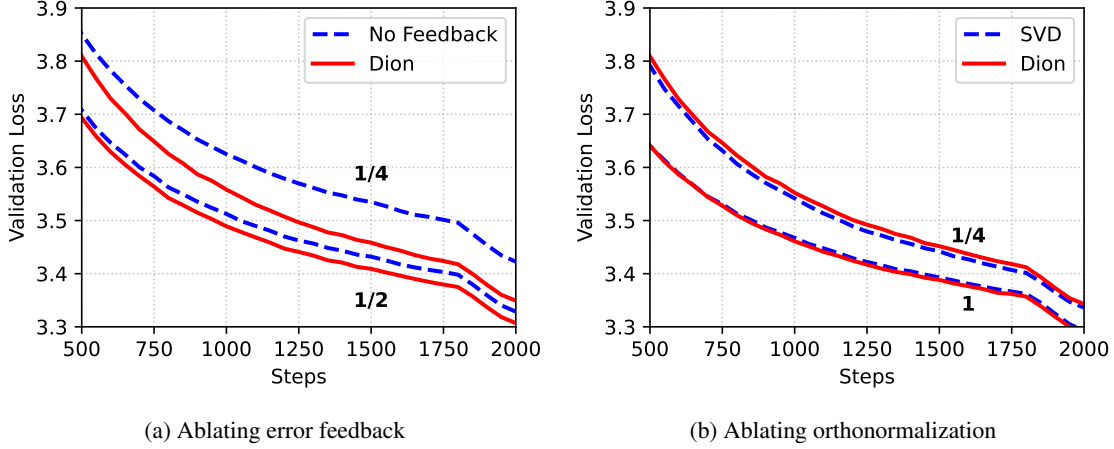


Figure 3: **Testing algorithmic components.** Ablation results on two key design choices in Dion: the error feedback mechanism and the use of a single power iteration. Error feedback proves crucial at lower ranks, while a single power iteration performs on par with full SVD, offering a more efficient alternative without sacrificing performance.

4.1 Importance of Error Feedback

To evaluate the necessity of the error feedback mechanism, we compare Dion against a simplified variant that omits error feedback. In this baseline, the rank- r approximation is applied directly to the momentum buffer:

$$M_t \leftarrow \mu M_{t-1} + G_t, \quad \text{for } \mu \in (0, 1),$$

whereas Dion first compresses an auxiliary buffer B_t and then updates the momentum using error feedback:

$$M_t \leftarrow B_t - (1 - \mu)P_t R_t^\top.$$

As shown in Figure 3a, the version without error feedback suffers from steep performance degradation as the rank decreases—from $r = d/2$ to $r = d/4$ —while Dion maintains stable performance. This highlights the importance of incorporating error feedback to preserve optimization quality under compression.

We note, however, that our current error feedback mechanism is not necessarily optimal. Exploring more effective or adaptive variants is an interesting direction for future work.

4.2 Single Power Iteration vs. Full SVD

We next evaluate the effectiveness of Dion’s single power iteration for computing the rank- r approximation. In this experiment, we compare Dion to a variant that replaces the power iteration with a full singular value decomposition (SVD) at every step—an ideal but computationally expensive alternative.

The results, shown in Figure 3b, reveal negligible differences in convergence behavior between the two approaches. This suggests that Dion’s power iteration—initialized from the previous iteration’s right orthonormal basis—offers a sufficiently accurate approximation at a fraction of the computational cost.

5 Related Work

To reduce communication overhead in distributed training, two main strategies have emerged: gradient quantization/sparsification [Wang et al., 2023] and federated averaging [McMahan et al., 2017, Douillard et al., 2023]. While highly effective, these techniques are orthogonal to our method. As discussed in the introduction, our approach focuses on designing a base optimizer with inherently lower communication costs. It can be used in conjunction with quantization or federated averaging to further improve efficiency.

Closely related to our work is the DeMO optimizer [Peng et al., 2024], which compresses gradients using the discrete cosine transform along with an error feedback mechanism. However, in our preliminary experiments, DeMO underperformed compared to our approach, particularly at large batch sizes. Even at increased sparsity levels, DeMO did not match the performance of state-of-the-art optimizers such as Muon and Adam. In contrast, Dion with full sparsity (i.e., full rank) outperformed Muon in the large-batch regime. A more comprehensive comparison with DeMO is left for future work.

Our approach is also related to recent work on low-rank updates, which have gained attention following the success of GaLore [Zhao et al., 2024] in reducing memory consumption. Since Dion often converges faster than Adam, combining it with GaLore’s core insights may lead to other competitive memory-efficient optimizers. Our findings also offer an interesting perspective on the phenomenon observed by Song et al. [2025], where training fails within the top eigenspace of the Hessian. Roughly speaking, our results suggest that this limitation can, in some sense, be mitigated through an error-feedback mechanism, allowing effective training even with low-rank updates.

References

- Jeremy Bernstein and Laker Newhouse. Old optimizer, new norm: An anthology. *arXiv preprint arXiv:2409.20325*, 2024.
- Arthur Douillard, Qixuan Feng, Andrei A Rusu, Rachita Chhaparia, Yani Donchev, Adhiguna Kuncoro, Marc’Aurelio Ranzato, Arthur Szlam, and Jiajun Shen. Diloco: Distributed low-communication training of language models. *arXiv preprint arXiv:2311.08105*, 2023.
- Keller Jordan, Jeremy Bernstein, Brendan Rappazzo, @fernbear.bsky.social, Boza Vlado, You Jiacheng, Franz Cesista, Braden Koszarsky, and @Grad62304977. modded-nanogpt: Speedrunning the nanogpt baseline, 2024a. URL <https://github.com/KellerJordan/modded-nanogpt>.
- Keller Jordan, Yuchen Jin, Vlado Boza, You Jiacheng, Franz Cecista, Laker Newhouse, and Jeremy Bernstein. Muon: An optimizer for hidden layers in neural networks, 2024b. URL <https://kellerjordan.github.io/posts/muon/>.
- Andrej Karpathy. nanoGPT. <https://github.com/karpathy/nanoGPT>, 2023. Accessed: 2025-01-25.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Shen Li, Yanli Zhao, Rohan Varma, Omkar Salpekar, Pieter Noordhuis, Teng Li, Adam Paszke, Jeff Smith, Brian Vaughan, Pritam Damania, et al. PyTorch distributed: experiences on accelerating data parallel training. *Proceedings of the VLDB Endowment*, 13(12):3005–3018, 2020.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR, 2017.
- Guilherme Penedo, Hynek Kydlíček, Loubna Ben allal, Anton Lozhkov, Margaret Mitchell, Colin Raffel, Leandro Von Werra, and Thomas Wolf. The fineweb datasets: Decanting the web for the finest text data at scale. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024. URL <https://openreview.net/forum?id=n6Sckn2QaG>.
- Bowen Peng, Jeffrey Quesnelle, and Diederik P Kingma. Decoupled momentum optimization. *arXiv preprint arXiv:2411.19870*, 2024.

Thomas Pethick, Wanyun Xie, Kimon Antonakopoulos, Zhenyu Zhu, Antonio Silveti-Falls, and Volkan Cevher. Training deep learning models with norm-constrained lmos. *arXiv preprint arXiv:2502.07529*, 2025.

Minhak Song, Kwangjun Ahn, and Chulhee Yun. Does SGD really happen in tiny subspaces? In *The Fourteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=v6iLQBoIJw>.

Thijs Vogels, Sai Praneeth Karimireddy, and Martin Jaggi. PowerSGD: Practical low-rank gradient compression for distributed optimization. *Advances in Neural Information Processing Systems*, 32, 2019.

Jue Wang, Yucheng Lu, Binhang Yuan, Beidi Chen, Percy Liang, Christopher De Sa, Christopher Re, and Ce Zhang. Cocktailsd: Fine-tuning foundation models over 500mbps networks. In *International Conference on Machine Learning*, pages 36058–36076. PMLR, 2023.

Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuandong Tian. Galore: Memory-efficient llm training by gradient low-rank projection. In *International Conference on Machine Learning*, pages 61121–61143. PMLR, 2024.

Yanli Zhao, Andrew Gu, Rohan Varma, Liang Luo, Chien-Chin Huang, Min Xu, Less Wright, Hamid Shojanazeri, Myle Ott, Sam Shleifer, et al. PyTorch FSDP: experiences on scaling fully sharded data parallel. *Proceedings of the VLDB Endowment*, 16(12):3848–3860, 2023.

A Proof of the Equivalence

Theorem A.1. *The distributed version of Dion (Algorithm 2) is equivalent to the centralized version of Dion (Algorithm 1). In particular, if X_t and \tilde{X}_t denote the iterates produced by Algorithm 1 and Algorithm 2, respectively, then*

$$X_t = \tilde{X}_t, \quad \forall t \geq 0.$$

Proof. We prove the claim by induction on the iteration index t . At $t = 0$, all parameters are initialized identically in both algorithms, so the statement trivially holds. Suppose that for iteration $t - 1$, the following two properties hold:

1. $P_{t-1} = \tilde{P}_{t-1}$, $Q_{t-1} = \tilde{Q}_{t-1}$, $R_{t-1} = \tilde{R}_{t-1}$, and $X_{t-1} = \tilde{X}_{t-1}$.
2. $M_{t-1} = \mathbb{E}_{\text{DP}}[\hat{M}_{t-1}]$, i.e., the centralized momentum equals the averaged momentum (over the data-parallel axis).

We now show these statements also hold at iteration t .

First, observe that

$$\mathbb{E}_{\text{DP}}[\hat{B}_t^{(i)}] = \mathbb{E}_{\text{DP}}[\hat{M}_{t-1}^{(i)}] + \mathbb{E}_{\text{DP}}[\hat{G}_t^{(i)}] = M_{t-1}^{(i)} + G_t^{(i)} = B_t^{(i)}.$$

Hence, each shard’s local $\hat{B}_t^{(i)}$ differs from $B_t^{(i)}$ only by a data-parallel mean.

Next, in the distributed algorithm, each worker computes

$$\hat{P}_t^{(i)} = \hat{B}_t^{(i)} Q_{t-1}^{(i)},$$

which are then aggregated to form

$$\tilde{P}_t = \text{Orthogonalize}\left(\sum_i \mathbb{E}_{\text{DP}}[\hat{P}_t^{(i)}]\right) = \text{Orthogonalize}\left(\sum_i B_t^{(i)} Q_{t-1}^{(i)}\right).$$

Since concatenating local products reproduces the centralized product,

$$\sum_i B_t^{(i)} Q_{t-1}^{(i)} = B_t Q_{t-1},$$

it follows that

$$\tilde{P}_t = \text{Orthogonalize}(B_t Q_{t-1}) = P_t.$$

Similarly, each worker computes

$$\hat{R}_t^{(i)} = (\hat{B}_t^{(i)})^\top P_t,$$

so

$$\tilde{R}_t^{(i)} = \mathbb{E}_{\text{DP}}[\hat{R}_t^{(i)}] = \mathbb{E}_{\text{DP}}[\hat{B}_t^{(i)}]^\top P_t = (B_t^{(i)})^\top P_t = R_t^{(i)}.$$

Hence, the partitioned R_t and the centralized R_t match on all shards.

Next, the distributed algorithm normalizes each shard's $R_t^{(i)}$:

$$\tilde{Q}_t^{(i)} = \text{Distributed-ColumnNormalize}(R_t^{(i)}).$$

Since column norms are computed by the same weight-parallel all-reduce sum in both algorithms, it follows that

$$\tilde{Q}_t^{(i)} = Q_t^{(i)}.$$

Finally, the parameter update in each algorithm is

$$X_t^{(i)} = X_{t-1}^{(i)} - \eta P_t (Q_t^{(i)})^\top,$$

which, upon concatenation over i , matches the centralized update exactly. For the momentum, we have

$$\hat{M}_t^{(i)} = \hat{B}_t^{(i)} - (1 - \mu) P_t (R_t^{(i)})^\top, \quad \text{hence} \quad \mathbb{E}_{\text{DP}}[\hat{M}_t^{(i)}] = B_t^{(i)} - (1 - \mu) P_t (R_t^{(i)})^\top = M_t^{(i)},$$

so the centralized and averaged momentum also coincide at iteration t .

Thus, both induction hypotheses remain true at step t . By induction, the iterates in the two algorithms match for all t , completing the proof. \square