

ASYNCMESH: FULLY ASYNCHRONOUS OPTIMIZATION FOR DATA AND PIPELINE PARALLELISM

Anonymous authors

Paper under double-blind review

ABSTRACT

Data and pipeline parallelism are key strategies for scaling neural network training across distributed devices, but their high communication cost necessitates co-located computing clusters with fast interconnects, limiting their scalability. We address this communication bottleneck by introducing *asynchronous updates across both parallelism axes*, relaxing the co-location requirement at the expense of introducing *staleness* between pipeline stages and data parallel replicas. To mitigate staleness, for pipeline parallelism, we adopt a weight look-ahead approach, and for data parallelism, we introduce an *asynchronous sparse averaging* method equipped with an exponential moving average based correction mechanism. We provide convergence guarantees for both sparse averaging and asynchronous updates. Experiments on large-scale language models (up to *1B parameters*) demonstrate that our approach matches the performance of the fully synchronous baseline, while significantly reducing communication overhead.

1 INTRODUCTION

Distributed optimization approaches enable large-scale model training by partitioning computation across multiple interconnected devices, primarily through Data Parallelism (DP) (Goyal, 2017; Li et al., 2020; Zhao et al., 2023) and Model Parallelism (MP) (Huang et al., 2019; Krizhevsky et al., 2017; Shoeybi et al., 2019). While DP replicates the model across devices with partitioned data, MP partitions the model itself across devices. Combining these approaches allows training foundation models at the frontier scale (Dubey et al., 2024; Liu et al., 2024a). However, both DP and MP rely on high-bandwidth interconnects due to high communication costs, limiting distributed training to co-located computing clusters. Therefore, scaling beyond a centralized infrastructure requires addressing communication bottlenecks inherent to both DP and MP setups.

Existing approaches mitigate communication costs via compression (Bernstein et al., 2018; Wang et al., 2023; Wangni et al., 2018) and/or overlapping computation with communication using scheduling (Narayanan et al., 2019; Qi et al., 2023) or asynchronous methods (Agarwal & Duchi, 2011; Stich & Karimireddy, 2019). We focus on asynchronous methods, which by design, offer full utilization of the distributed infrastructure and support heterogeneous hardware, by eliminating synchronization barriers. In this work, we study asynchronous training in a *2D mesh* combining DP and Pipeline Parallelism (PP) (Huang et al., 2019) – a special case of MP that partitions the model into sequential stages.

In a synchronous mesh, each stage within a pipe¹ is optimized in a lock-step manner, ensuring that the weights and gradients are synchronized at each step, and the model replicas are explicitly averaged by pausing optimization. In contrast, **ASynCMesh** eliminates both synchronization points allowing decoupled optimization in each stage as well as decoupling optimization in a pipe and averaging across replicas. This enables *continuous data processing without communication bottlenecks*. Fig. 1 for an illustrates this. Asynchronism, however, introduces *staleness* in model weights, necessitating correction mechanisms to ensure consensus among model replicas and convergence (Agarwal & Duchi, 2011; Ajanthan et al., 2025a; Stich & Karimireddy, 2019; Zheng et al., 2017).

For PP, the staleness is stage-dependent (Narayanan et al., 2019), and extrapolating the previous update direction using the Nesterov method is shown to effectively compensate for this delay (Ajanthan

¹A pipe is a sequence of stages that form a full model.

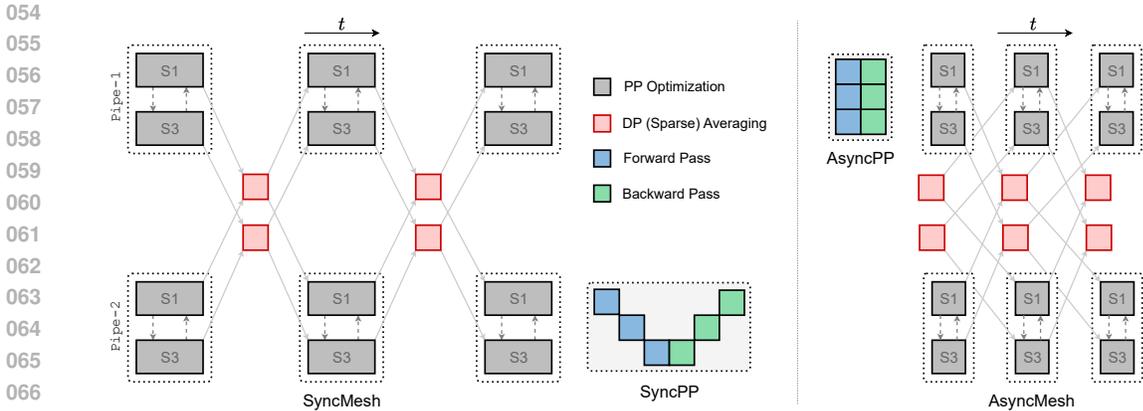


Figure 1: *SyncMesh vs AsyncMesh, for a 3-stage, 2-DP replica setup (only 2 stages: S1 and S3, are shown for clarity). Notably, AsyncMesh eliminates idle time due to communication for both PP and DP. In synchronous DP, devices are idle while parameters are averaged, whereas asynchronous DP eliminates this idle time by using the “old” average. Moreover, in AsyncPP, each stage alternates between forward and backward passes without any communication delay.*

et al., 2025a). Due to its simplicity and empirical effectiveness, we adopt this strategy (Ajanthan et al., 2025a) to optimize each pipe asynchronously. For DP, the staleness depends on the interconnect bandwidth and the data transfer volume. Therefore, to minimize staleness, we *exchange only 5% of weights* across replicas, similar to sparse averaging (Beton et al., 2025; Fournier et al., 2024), and communicate them asynchronously to *completely mask the DP communication*. To address resulting weight discrepancies, we design an *Exponential Moving Average (EMA)* based correction mechanism that approximates the average staleness. Theoretical analysis under the homogeneous setting (identical hyperparameters and i.i.d. data splits) shows that our method ensures consensus among replicas on expectation, despite asynchronous sparse averaging.

We validate our approach on large-scale language modeling tasks with decoder-only transformer architectures (Karpathy, 2022; Vaswani et al., 2017). Experiments demonstrate that our method *matches the performance of the fully synchronous baseline*, while eliminating the communication overhead via sparse asynchronous updates. Notably, for the first time, we train a *1B parameter model* to convergence in AsyncMesh matching the performance of the synchronous alternative. Our results show the feasibility of distributed training over bandwidth constrained interconnects using asynchronous optimization.

Our contributions can be summarized as follows:

- We introduce a new **AsyncMesh** setup where *both DP and PP are asynchronous*, and present a fully asynchronous method that can match the performance of the fully synchronous method.
- We provide theoretical justification of convergence in the presence of staleness in a homogeneous setup where only a small subset of weights is communicated between DP replicas.
- We empirically demonstrate the robustness and scalability of our approach across varying subset sizes, staleness levels, DP communication intervals, and $PP \times DP$ mesh sizes.

2 PRELIMINARIES

We first define the 2D mesh configuration with DP (Goyal, 2017; Li et al., 2020) and PP (Huang et al., 2019), and then briefly review the Asynchronous Pipeline Parallel (AsyncPP) method (Ajanthan et al., 2025a) and Sparse Parameter Averaging (SPARTA) (Beton et al., 2025; Fournier et al., 2024) upon which we build our work. We refer the reader to the respective papers for more details.

2.1 PROBLEM SETUP: 2D MESH

Let P be the number of pipeline stages and m be the number of replicas. We consider a symmetric 2D mesh, where each stage is replicated m times, constituting Pm devices (or workers) in total. A *pipe* is a sequence of stages that form the full model, which can be defined using its forward and

backward functions. Consider a pipe (or replica) $i \in \{1, \dots, m\}$, and let the forward function for stage j be $f_{ij} := f_j(\mathbf{w}_{ij}; \mathbf{x}_{j-1})$ with weights \mathbf{w}_{ij} , and input \mathbf{x}_{j-1} . Then, the forward and backward functions for pipe i can be defined as:

$$\begin{aligned} F(\mathbf{W}_i; \mathbf{x}_0) &:= f_{iP} \circ f_{iP-1} \circ \dots \circ f_{i1}(\mathbf{x}_0), & \text{forward,} \\ G(\mathbf{W}_i; \mathbf{e}_P) &:= g_{i1} \circ g_{i2} \circ \dots \circ g_{iP}(\mathbf{e}_P), & \text{backward,} \end{aligned} \quad (1)$$

where $\mathbf{W}_i = \{\mathbf{w}_{iP}, \dots, \mathbf{w}_{i1}\}$ and \mathbf{x}_0 is an input data point. Here, $g_{ij} := g_j(\mathbf{w}_{ij}; \mathbf{e}_j)$ is the backward function for stage j corresponding to f_{ij} and \mathbf{e}_P is the error signal corresponding to \mathbf{x}_0 . In our 2D mesh, there are m such pipes, and the goal is to optimize the following consensus objective (Boyd et al., 2011):

$$\begin{aligned} \min_{\mathbf{W} \in \mathbb{R}^d} F(\mathbf{W}; \mathcal{D}) &:= \min_{\mathbf{W}_i \in \mathbb{R}^d} \sum_{i=1}^m F(\mathbf{W}_i; \mathcal{D}_i), \\ \text{s.t. } \mathbf{W}_i &= \mathbf{W}, \quad \forall i \in \{1, \dots, m\}, \end{aligned} \quad (2)$$

where \mathcal{D}_i is an i.i.d. subset of the dataset \mathcal{D} and d is the number of learnable parameters of the full model. Each pipe is optimized independently on its own data split and using its own optimizer, and its weights are synchronized periodically with other pipes – typically after every optimization step.

2.2 PIPELINE PARALLELISM

Pipeline Parallelism (PP) methods (Guan et al., 2024) provide communication efficient ways to optimize the model parameters in a pipe. Specifically, pipeline scheduling strategies such as GPipe (Huang et al., 2019), 1F1B (Narayanan et al., 2021b), and ZeroBubble (Qi et al., 2023) design the order of processing forward and backward passes of microbatches in a pipe to reduce communication overhead between stages and improve device utilization. In contrast, asynchronous PP methods (Ajanthan et al., 2025a; Narayanan et al., 2019) eliminate the requirement to synchronize the weights and gradients across stages at each update step, offering full pipeline utilization.²

Asynchronous Pipeline Parallel. The main challenge in the Asynchronous Pipeline Parallel (AsyncPP) method (Ajanthan et al., 2025a) is the discrepancy between the gradients at a particular stage and the corresponding weights due to asynchronous updates. Specifically, since the weights of a particular stage are updated multiple times between the forward and backward passes of a microbatch, outdated gradients are used for weight updates. Formally, the weight update for pipe i and stage j can be written as (omitting optimizer specifics):

$$\mathbf{w}_{ij}^{t+1} = \mathbf{w}_{ij}^t - \eta_t \nabla f_j(\mathbf{w}_{ij}^{t-\delta_j}; \mathcal{B}_i^{t-\delta_j}), \quad (3)$$

where $\eta_t > 0$ is the learning rate, $\mathcal{B}_i^{t-\delta_j}$ is the minibatch, and $\delta_j \geq 0$ is the stage-dependent constant delay due to asynchronous PP updates.

To compensate for this delay, AsyncPP (Ajanthan et al., 2025a) extrapolates the last update direction ($\mathbf{w}_{ij}^t - \mathbf{w}_{ij}^{t-1}$) using the Nesterov Accelerated Gradient (NAG) framework (Nesterov, 1983), and shows that it acts as gradient delay correction in the weight space. This was shown to surpass the synchronous GPipe method on some language modeling tasks, and we employ this approach to optimize each pipe in our experiments.

2.3 DATA PARALLELISM

In a typical Data Parallelism (DP) setup (Goyal, 2017; Li et al., 2020), each worker computes the gradient of the full model on a minibatch and communicates it to the central parameter server. The server averages the gradients from all workers, performs an optimization step of the server model, and distributes the updated parameters to the workers for the next iteration. This is equivalent to performing gradient based optimization using a larger minibatch, and the consensus constraint in Eq. (2) is maintained.

²Asynchronous PP methods (Ajanthan et al., 2025a; Narayanan et al., 2021a) assume comparable compute and communication times per stage, but if the inter-stage bandwidth is low, compression (Ramasinghe et al., 2025) is needed to fully eliminate PP communication overhead.

We consider a serverless scenario, which better reflects a decentralized training setup (Ryabinin et al., 2023). In this, each worker performs a local update using its optimizer and averages the updated weights with others (McMahan et al., 2017), or equivalently averages gradients before applying local updates (Ryabinin et al., 2021). Synchronizing weights (instead of gradients) is becoming popular as they can be communicated infrequently to reduce data transfer (Douillard et al., 2025; 2023; Reddi et al., 2020). To this end, we consider the setup where only a small subset is averaged periodically (Beton et al., 2025; Fournier et al., 2024; Lee et al., 2023) which is shown to perform similarly to traditional DP.

Sparse Parameter Averaging. In Sparse Parameter Averaging (SPARTA) (Beton et al., 2025; Fournier et al., 2024), after each local update, a small subset of parameters (*e.g.*, 5%) is averaged across workers, reducing data transfer. Adopting this to our 2D mesh is straightforward, as the local update is performed on each pipe, and the sparse averaging is done between stage replicas with a randomly sampled subset. Formally, sparse averaging for stage j with subset $\mathcal{S}_j^t \subset \{1, \dots, d_j\}$ can be written as (omitting optimizer specifics):

$$\begin{aligned} \hat{\mathbf{w}}_{ij}^t &= \mathbf{w}_{ij}^{t-1} - \eta_{t-1} \nabla f_j(\mathbf{w}_{ij}^{t-1}; \mathcal{B}_i^{t-1}), & \forall i, & \text{local update,} \\ w_{ij;\mu}^t &= \begin{cases} \frac{1}{m} \sum_i \hat{w}_{ij;\mu}^t, & \text{if } \mu \in \mathcal{S}_j^t, \\ \hat{w}_{ij;\mu}^t, & \text{if } \mu \notin \mathcal{S}_j^t, \end{cases} & \forall i, \mu, & \text{sparse averaging.} \end{aligned} \quad (4)$$

Here, $\eta_t > 0$ is the learning rate, \mathcal{B}_i^{t-1} is the minibatch, and $w_{ij;\mu}^t$ is the μ -th element of vector \mathbf{w}_{ij}^t .

3 OUR METHOD: OPTIMIZATION IN ASYNCMESH

We consider AsyncMesh, where both PP and DP axes in the consensus optimization problem (Eq. (2)) are optimized asynchronously. By making the mesh fully asynchronous, our setup ensures *full pipeline utilization* throughout training without interruption, while encouraging consensus among model replicas. Our setup is clearly illustrated in Fig. 1. For optimizing each pipe, we employ the AsyncPP method (refer to Sec. 2.2) that uses a variant of NAG for delay correction within a pipe. For DP, we introduce an asynchronous version of sparse parameter averaging that eliminates the communication overhead due to DP.

3.1 ASYNCHRONOUS SPARSE PARAMETER AVERAGING

Let us consider a particular stage, and drop the stage index for simplified notation. In asynchronous DP, the local updates in each worker (*i.e.*, PP optimization) do not wait for the averaging operation (*i.e.*, DP communication) to complete. Therefore, the weights at each stage are updated multiple times while the averaging operation is performed between the stage replicas. Thus, the averaged weights are older than the weights at each worker, which leads to *staleness*. Let τ be the corresponding delay, then, the delayed sparse averaging can be written as:

$$w_{i;\mu}^t = \begin{cases} \bar{w}_{i;\mu}^{t-\tau}, & \text{if } \mu \in \mathcal{S}^{t-\tau}, \\ \hat{w}_{i;\mu}^t, & \text{if } \mu \notin \mathcal{S}^{t-\tau}, \end{cases} \quad \forall i, \mu, \quad \text{where } \bar{\mathbf{w}}^t = \frac{1}{m} \sum_i \hat{\mathbf{w}}_i^t. \quad (5)$$

Compared to Eq. (4), the only difference is that elements in the subset $\mathcal{S}^{t-\tau}$ are set to the *old average* $\bar{w}_{i;\mu}^{t-\tau}$ instead of the new one at time t . This delay is detrimental to training as 1) it ignores the local updates between $t-\tau$ and t for the subset, and 2) the weight vector has a discrepancy as some weights correspond to time $t-\tau$ and the rest at time t . Therefore, it is essential to compensate for this delay.

3.2 DELAY CORRECTION VIA ESTIMATING AVERAGE STALENESS

Our idea is to approximate the new average $\bar{\mathbf{w}}^t$ using the old average $\bar{\mathbf{w}}^{t-\tau}$ and the estimated *average staleness*. Specifically, we estimate the average staleness in each stage independently using Exponential Moving Average (EMA) of staleness (*i.e.*, weight differences) throughout training. Precisely, we estimate the new average as,

$$\mathbf{d}_i^t \llbracket \mathcal{S}^{t-\tau} \rrbracket = (1 - \lambda_t) \mathbf{d}_i^{t-1} \llbracket \mathcal{S}^{t-\tau} \rrbracket + \lambda_t (\hat{\mathbf{w}}_i^t \llbracket \mathcal{S}^{t-\tau} \rrbracket - \hat{\mathbf{w}}_i^{t-\tau} \llbracket \mathcal{S}^{t-\tau} \rrbracket), \quad \text{EMA of staleness,} \quad (6)$$

$$\hat{\mathbf{w}}_i^t \llbracket \mathcal{S}^{t-\tau} \rrbracket = \bar{\mathbf{w}}^{t-\tau} \llbracket \mathcal{S}^{t-\tau} \rrbracket + \mathbf{d}_i^t \llbracket \mathcal{S}^{t-\tau} \rrbracket, \quad \text{estimate of average at } t,$$

where $\mathbb{I}[\cdot]$ denotes the indicator and $\lambda_t \in (0, 1)$ is the EMA coefficient. Then, the delay corrected asynchronous sparse averaging takes the following form:

$$w_{i:\mu}^t = \begin{cases} \tilde{w}_{i:\mu}^t, & \text{if } \mu \in \mathcal{S}^{t-\tau}, \\ \hat{w}_{i:\mu}^t, & \text{if } \mu \notin \mathcal{S}^{t-\tau}, \end{cases} \quad \forall i, \mu. \quad (7)$$

Intuitively, the idea here is that \mathbf{d}_i^t being the EMA of staleness, robustly estimates the average staleness. Therefore, $\tilde{\mathbf{w}}_i^t$ approximates the average at time t . Concretely,

$$\tilde{\mathbf{w}}_i^t = \bar{\mathbf{w}}^{t-\tau} + \mathbf{d}_i^t \approx^{\mathbf{a}} \bar{\mathbf{w}}^{t-\tau} + \mathbb{E}[\hat{\mathbf{w}}_i^t - \hat{\mathbf{w}}_i^{t-\tau}] \approx^{\mathbf{b}} \bar{\mathbf{w}}^{t-\tau} + \bar{\mathbf{w}}^t - \bar{\mathbf{w}}^{t-\tau} = \bar{\mathbf{w}}^t. \quad (8)$$

Here, the approximation \mathbf{a} is due to EMA being a stochastic approximation of the expected value (Robbins & Monro, 1951), and \mathbf{b} is due to the empirical average of DP replicas. Note the accuracy of the EMA approximation depends on the delay τ (lower the better) and the smoothness of the optimization trajectory.³ Whereas the approximation error of \mathbf{b} reduces with increasing number of DP replicas.

Note that sparse averaging is performed after every local step in our description so far. However, this requirement can be relaxed and averaging can be performed every K steps, similar to (Douillard et al., 2023; McMahan et al., 2017). In this setup, our approach is a strict generalization of the concurrent idea of eager DiLoCo (Kale et al., 2025), in which, all the parameters are communicated, the delay $\tau = K$, and $\mathbf{d}_i^t = \frac{1}{m} (\hat{\mathbf{w}}_i^t - \hat{\mathbf{w}}_i^{t-\tau})$.

3.3 THEORETICAL INSIGHTS

Recall that our aim is to optimize the consensus objective Eq. (2) with asynchronous updates for both PP and DP. AsyncPP (Ajanthan et al., 2025a) proved convergence with fixed delay, providing a theoretical justification for the single pipeline setup. However, it is unclear if sparse averaging ensures consensus, and since we also make it asynchronous, the effect of staleness in achieving consensus needs to be studied.

To this end, we take a first step in theoretically understanding the conditions required for achieving consensus for asynchronous sparse averaging. First, we show that sparse averaging ensures consensus on expectation if the learning rate is chosen proportional to the subset size. Then, we provide a theoretical insight showing that under standard assumptions of stochastic approximation (Robbins & Monro, 1951) with an identical setup for each replica, the EMA approximates the average staleness and consequently ensures consensus on expectation. Both of these provide a theoretical justification that our approach ensures consensus, and when coupled with the standard convergence proof for Stochastic Gradient Descent (SGD) (Bottou et al., 2018), show that our approach converges to a fixed point of Eq. (2).

Suppose the consensus error is defined as:

$$\|\Delta^t\|^2 := \sum_{i=1}^m \|\mathbf{w}_i^t - \mathbf{w}^t\|^2, \quad \text{where } \mathbf{w}^t = \frac{1}{m} \mathbf{w}_i^t. \quad (9)$$

We first show that, on expectation, the consensus error vanishes, *i.e.*, all model replicas converge to their average. Now, by invoking the standard convergence proof for SGD (Bottou et al., 2018) on the averaged model, one can show convergence for sparse parameter averaging.

Theorem 1 (Sparse averaging ensures consensus). *Let f be a L -smooth function, the stochastic gradient be an unbiased estimate of ∇f and have bounded variance, and $p > 0$ be the averaging probability for an element $\mu \in \{1, \dots, d\}$, then, for an appropriate choice of learning rate $\eta_t > 0$, the consensus error for updates in Eq. (4) diminishes on expectation, *i.e.*, $\lim_{t \rightarrow \infty} \mathbb{E}[\|\Delta^t\|^2] = 0$.*

Proof. We first show that sparse averaging shrinks the consensus error by a factor of $(1-p)$ on expectation. Then, we derive a recursion on $\mathbb{E}[\|\Delta^t\|^2]$ and choose a learning rate proportional to p to ensure that it is a contraction. Detailed proof is provided in the appendix. \square

³Since EMA is estimated across different time steps, a smoother trajectory (*i.e.*, slow drift in average staleness) leads to a better approximation, conditions for this are stated in the next section.

270 Consensus of partial averaging has been previously studied in federated learning for pre-determined
 271 subsets (Lee et al., 2023), where a similar relationship between the learning rate and subset size
 272 is observed. This suggests convergence may be slow for small subset sizes (equivalently, small p)
 273 due to small learning rates. However, the tightness of this result is unclear, and we leave any such
 274 analysis for future work. In practice, we use the standard learning rate value and do *not* adjust it
 275 based on the subset size.

276 For the theoretical analysis of asynchronous sparse averaging, we consider a *homogeneous setup*
 277 with the same initialization, optimizer hyperparameters, and i.i.d. data subsets, in which we can show
 278 that the expected staleness can be independently estimated in each replica. Suppose the expected
 279 staleness and its drift be:

$$280 \mathbf{D}^t := \mathbb{E}[\hat{\mathbf{w}}_i^t - \hat{\mathbf{w}}_i^{t-\tau}], \quad \alpha_t := \mathbf{D}^t - \mathbf{D}^{t-1}. \quad (10)$$

282 With standard assumptions from stochastic approximation theory (Robbins & Monro, 1951; Robbins
 283 & Siegmund, 1971) such as diminishing EMA coefficient, and diminishing drift in expected stale-
 284 ness, we show that asynchronous (*i.e.*, delayed) averaging ensures consensus on expectation. This,
 285 together with Theorem 1, provides a theoretical justification for asynchronous sparse averaging.

286 **Theorem 2 (Delayed averaging with EMA ensures consensus).** *Consider a homogeneous setup*
 287 *where the average staleness \mathbf{D}^t is bounded and its drift is diminishing, i.e., $\lim_{t \rightarrow \infty} \|\alpha_t\| = 0$.*
 288 *Then, the EMA based delay correction as in Eq. (6) with λ_t satisfying $\sum_t \lambda_t = \infty$, $\sum_t \lambda_t^2 < \infty$,*
 289 *and $\sum_t \frac{\|\alpha_t\|^2}{\lambda_t} < \infty$ ensures consensus on expectation, i.e., $\lim_{t \rightarrow \infty} \mathbb{E}[\|\Delta^t\|^2] = 0$.*

291 *Proof.* Intuitively, in a homogeneous setup, a particular weight trajectory is equally probable for all
 292 replicas, and hence, the expected staleness can be independently estimated in each replica. More-
 293 over, classical stochastic approximation theory (Robbins & Monro, 1951) shows that EMA approx-
 294 imates the expected value. Here, since \mathbf{D}^t is time varying, we use the additional assumption on the
 295 interplay between the drift and the EMA coefficient to bound the consensus error. Detailed proof is
 296 provided in the appendix. \square

297 The assumption of diminishing $\|\alpha_t\|$ and $\sum_t \frac{\|\alpha_t\|^2}{\lambda_t} < \infty$ can be interpreted as the optimization
 298 trajectory being smoother such that the average staleness changes slowly with time and its difference
 299 diminishes towards convergence. This also puts a restriction on the allowed delay τ . To this end,
 300 sparse averaging with asynchronous updates is appealing as we can reduce the delay τ by only
 301 communicating a small subset at each iteration over a bandwidth limited interconnect.

302 Note the theory recommends diminishing EMA coefficient λ_t . In our implementation, we initialize
 303 λ_t to 0.5 and after 1k iterations, we gradually decay it to 0.01 using a cosine schedule. This slightly
 304 improves over a constant λ_t , and we keep this schedule fixed for all the experiments.

306 4 RELATED WORK

307 **Communication Efficient DP Methods.** DP is a traditional distributed training setup (Goyal,
 308 2017; Li et al., 2020), where each device computes the gradients of the full model in its data split,
 309 aggregates the gradients on a central server, performs a central optimization step, and distributes the
 310 updated model parameters back to the devices for the next iteration. To reduce the communication
 311 overhead, the gradients can be compressed using low-rank approximation (Vogels et al., 2019; Zhao
 312 et al., 2024), sparsification (Wang et al., 2017; Wangni et al., 2018), and quantization (Bernstein
 313 et al., 2018; Wang et al., 2023). Alternatively, in the serverless setup, each device individually up-
 314 dates the model and periodically synchronizes the model parameters, where partial averaging (Be-
 315 ton et al., 2025; Fournier et al., 2024; Lee et al., 2023) and/or infrequent communication as in
 316 DiLoCo (Douillard et al., 2025; 2023; Reddi et al., 2020) reduce communication overhead.

317 **Asynchronous DP Methods.** Unlike synchronous methods, asynchronous DP methods can fully
 318 eliminate the communication overhead. These are well-studied for the parameter server setup that
 319 communicates the gradients, and many gradient delay correction mechanisms have been devel-
 320 oped (Agarwal & Duchi, 2011; Assran et al., 2020; Stich & Karimireddy, 2019). Notable meth-
 321 ods that improve over the simple asynchronous SGD (Recht et al., 2011) include delay dependent
 322 learning rate (Barkai et al., 2019; Mishchenko et al., 2022), gradient forecasting with second-order
 323 information (Zheng et al., 2017), and look-ahead in the weight space (Hakimi et al., 2019). Apart
 from this, training dynamics of such asynchronous DP methods have also been analyzed (Mitliagkas

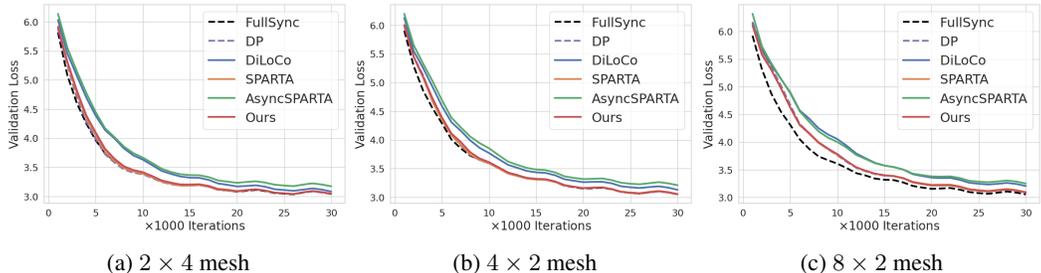


Figure 2: Results on WikiText with varying mesh configurations with AsyncPP for all methods except FullSync. In all scenarios, our method matches the performance of the fully synchronous method, while outperforming the fully asynchronous baseline AsyncSPARTA.

et al., 2016). These methods are not directly applicable in our setup, as we communicate weights instead of gradients between replicas. Asynchronous methods are underexplored in this context, although there are a few recent empirical methods designed for DiLoCo (Ajanthan et al., 2025b; Kale et al., 2025; Liu et al., 2024b), we show some of them (Kale et al., 2025) can be viewed as special cases of our method.

Asynchronous PP Methods. Asynchronous PP methods eliminate the synchronization bottleneck in PP (Guan et al., 2024) to achieve 100% pipeline utilization at the cost of gradient staleness. In PP, weight stashing is used to ensure correct backpropagation is performed (Narayanan et al., 2019; 2021a), however, the discrepancy (*i.e.*, staleness) between weights and gradients still needs to be compensated. For this, many correction mechanisms such as learning rate discounting (Yang et al., 2021), direct weight prediction (Chen et al., 2018; Guan et al., 2019), and extrapolation in the weight space (Ajanthan et al., 2025a; Zuo et al., 2025) have been developed.

In this work, for the first time, we consider AsyncMesh, where both DP and PP are asynchronous. For PP, we adopt the recent weight extrapolation method (Ajanthan et al., 2025a), and for DP, we combine sparse averaging (Beton et al., 2025) with asynchronous updates to fully eliminate the DP communication overhead.

5 EXPERIMENTS

5.1 EXPERIMENTAL SETUP

We evaluate on four large-scale language modelling datasets, namely, WikiText (WT) (Merity et al., 2016), BookCorpus (BC) (Zhu et al., 2015), OpenWebText (OWT) (Gokaslan et al., 2019), and FineWeb (FW) (Penedo et al., 2024), using decoder-only architectures with varying mesh configurations, denoted by PP-stages \times DP-replicas. Our architecture is based on NanoGPT (Karpathy, 2022) with no dropout. The base model has a context length of 1024, an embedding dimension of 768, 12 attention heads, and 12 layers, with approximately 163M parameters. We use the GPT2 tokenizer (Radford et al., 2019) and train the model from scratch. For AsyncPP, NAdamW optimizer (Dozat, 2016) with momentum 0.99 is used as per (Ajanthan et al., 2025a), and all other hyperparameters are set to the standard values for this task (refer to appendix), and kept constant for all the experiments.

Our aim is to show that neither asynchronous updates for both PP and DP, nor the sparse averaging for DP, deteriorate the validation performance, in various configurations. For fair comparison, we compute the validation loss (and perplexity) on the model averaged across all DP replicas (*i.e.*, consensus model), for all the methods. Primarily, we compare our method against **FullSync**, which is the ideal synchronized setup *without any communication efficient scheduling* for PP and only utilizes $\frac{1}{P}$ % of the pipeline (Huang et al., 2019; Yang et al., 2021). In addition, we evaluate the following synchronous DP methods: DP that averages all the parameters at every step, SPARTA (Beton et al., 2025) that averages a small subset, DiLoCo (Douillard et al., 2023) that performs infrequent communication, and AsyncSPARTA that performs asynchronous sparse averaging without delay correction. For sparse averaging methods, the subset size is 5%, the averaging interval is 1, and the asynchronous methods incur a 10-step delay, unless specified otherwise.

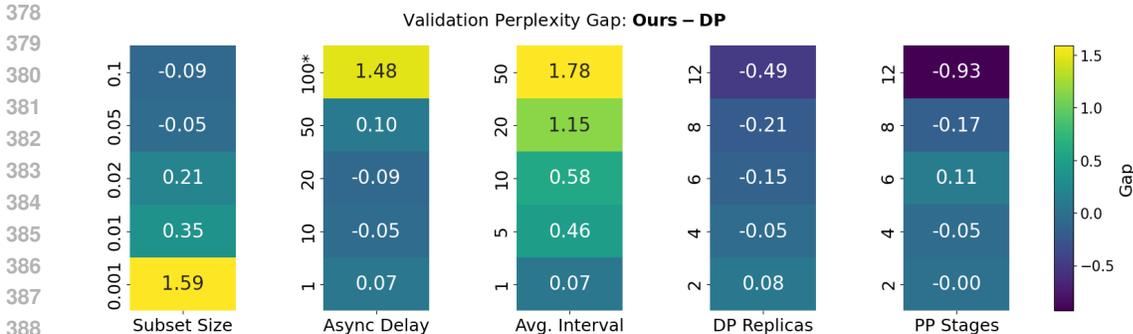


Figure 3: Perplexity gap (lower the better) between our method and DP on WikiText for different configurations. Our method is robust to a range of configurations and scales favourably with the number of DP replicas and PP stages. Here, 100* indicates, `async-delay = 10` with `avg-interval = 10`, simulating 100-step effective delay, as `async-delay = 100` did not converge.

Our method is implemented in PyTorch (Paszke et al., 2019) using the publicly available codes of SPARTA⁴, and AsyncPP⁵. Unless otherwise specified, all experiments use the base architecture described above and are performed on the WikiText dataset. All our experiments are performed on a system equipped with 8 A100 GPUs, and where needed, multiple such instances were used.

5.2 MAIN RESULTS

We analyze the validation loss trajectories for different mesh configurations with AsyncPP in Fig. 2. In all scenarios, our method outperforms AsyncSPARTA and matches the performance of FullSync. As shown in the appendix, the behaviour is similar for synchronous PP updates as well. Except AsyncSPARTA and DiLoCo, all methods show near identical performance in most cases. Note DiLoCo is a prominent method in the DP setup with full model in each replica, however, it seems to be inferior in the mesh setup with AsyncPP. To our knowledge, this is the first time DiLoCo is tested with AsyncPP.

In Table 1, we report the validation perplexities on multiple datasets after 30k iterations for the 4×2 mesh. Our method outperforms AsyncSPARTA by **3 – 6** points, and yields similar perplexities as FullSync, even surpassing it in 2 out of 4 datasets. This is remarkable as our method is fully asynchronous and only averages 5% of the parameters at each iteration. Note these results are with 2 DP replicas, and as shown in Fig. 3, the more replicas, the better for our method. For completeness, we trained to the compute-optimal point (Hoffmann et al., 2022) on FineWeb, where the perplexities are 19.92 for FullSync, and 20.10 for ours, confirming the merits of our method beyond doubt. These results demonstrate that our delay correction method effectively compensates for staleness.

Method	WT	BC	OWT	FW
FullSync	21.23	35.99	35.71	36.77
DP	21.26	35.24	35.97	36.81
SPARTA	21.30	35.15	35.73	37.10
AsyncSPARTA	24.80	37.83	41.41	43.20
Ours	21.14	35.09	36.13	37.31

Table 1: Validation perplexity scores at 30k iterations for the 4×2 mesh. Our method outperforms AsyncSPARTA, and matches FullSync while eliminating the DP communication overhead. Except FullSync, all other methods use AsyncPP.

5.3 VARYING CONFIGURATIONS

To test the robustness of our method in a variety of setups, we consider the base model in the following default setup: `subset-size = 5%`, `async-delay = 10`, `avg-interval = 1`, `DP-replicas = 4`, `PP-stages = 4`, and change one criterion at a time. To isolate the effect of our asynchronous sparse averaging, we measure the gap in validation perplexity between our method and DP, which synchronously averages all parameters, while pipeline is optimized with AsyncPP for both methods.

⁴<https://github.com/matttreed/diloco-sim>

⁵<https://github.com/PluralisResearch/AsyncPP>

As reported in Fig. 3, our method is robust for a wide range of configurations, and it is consistently better than DP in many scenarios. Notably, our method is robust for subset size as low as 1% (only a fractional change in perplexity), tolerates delay up to 50 steps, and works reasonably even if we average every 10 steps (*i.e.*, 10× less communication). More encouraging results are that *our method becomes better with a larger mesh* – with a larger number of DP replicas or a larger number of PP stages, the Gap: Ours - DP becomes more negative. This shows the superior scalability of our method compared to synchronously averaging all the parameters as in DP.

Note that our theory predicts (refer to Sec. 3.2) that a larger number of DP replicas would yield a better approximation of average staleness, leading to better delay correction. These results seem to corroborate that. Meanwhile, AsyncPP is shown to degrade with more PP stages (Ajanthan et al., 2025a), however, the gains from our sparse averaging may help offset this.

Practical Implications. Compared to SPARTA, our method achieves 1.5× – 3.7× speed-up according to our communication time measurements (refer to appendix for detailed results), with larger meshes yielding better speed-ups. Furthermore, the above results show that our method tolerates a 50-step delay with only a 5% subset size, reducing communication by 10× per iteration (subset and indices need to be exchanged). For simplicity, assuming 1s per forward-backward pass in a stage, this gives 50s for DP communication – enough to support up to 1.5B parameters (FP32) per stage over a 100 Mbps connection, highlighting the viability of decentralized training over the internet.

5.4 INCREASING THE MODEL SIZE

To demonstrate the scalability of our approach, we train a 1B parameter model in the asynchronous 2D mesh. We maintain the number of stages at 4, but increase the embedding dimension to 2304, with 24 attention heads. As illustrated in Fig. 4, the results align with those of the base model. Specifically, our approach matches FullSync throughout training, yielding validation perplexity of 19.53 compared to 19.62 for FullSync. This large-scale experiment demonstrates the merits of our method and the feasibility of asynchronous PP and DP optimization for language model training.

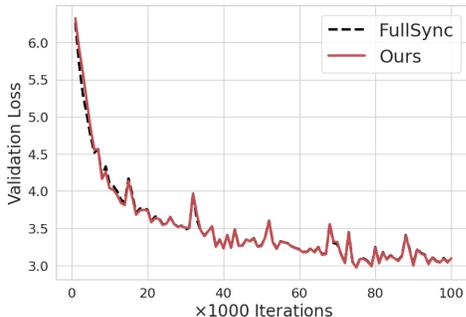


Figure 4: Results on the 1B parameter model on FineWeb for 4 × 2 mesh. Our method matches FullSync, similar to the base model.

5.5 ABLATION STUDY

To understand the effect of different components of our method, we perform an ablation study in 8 × 2 mesh on WikiText. Simply using weight differences within each replica, *i.e.*, $\mathbf{d}_i^t = \hat{\mathbf{w}}_i^t - \hat{\mathbf{w}}_i^{t-\tau}$, substantially improve AsyncSPARTA, where EMA and adaptive momentum coefficient as predicted by our theory further improves. Note that the effect of EMA is more prominent in the early phase of training, where the step sizes are larger (*i.e.*, noisier). This aligns with our intuition that EMA robustly approximates the average staleness.

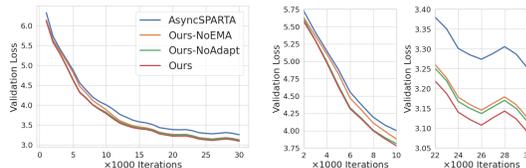


Figure 5: Ablation of our method on WikiText for the 8 × 2 configuration. EMA improves the early phase of training, and the adaptive momentum coefficient yields further marginal improvement.

6 CONCLUSION

In this paper, we studied a new AsyncMesh setup where both DP and PP are asynchronous, and introduced a fully asynchronous approach that can match the performance of the fully synchronized method. We theoretically showed that our method converges to a fixed point of the consensus objective on expectation, despite sparse averaging and asynchronous communication. Our experiments on a wide range of configurations demonstrate the merits of our approach, and show the feasibility of asynchronous optimization for large scale language model training. By alleviating communication overhead without any performance penalty, our approach takes a step towards realizing large-scale collaborative training over the internet.

486 REPRODUCIBILITY STATEMENT

487

488 All assumptions related to the theoretical claims are clearly mentioned in the theorem statements.
 489 Furthermore, necessary details to understand the experiments and results are provided in the main
 490 paper, and all hyperparameter settings are provided in the appendix for reproducibility. Our
 491 anonymized code is submitted as supplementary material and will be made publicly available upon
 492 acceptance.

493

494 REFERENCES

495

496 Alekh Agarwal and John C Duchi. Distributed delayed stochastic optimization. *Advances in neural*
 497 *information processing systems*, 24, 2011.

498 Thalaiyasingam Ajanthan, Sameera Ramasinghe, Gil Avraham, Yan Zuo, and Alexander Long. Nes-
 499 terov method for asynchronous pipeline parallel optimization. *ICML, 2025a*.

500

501 Thalaiyasingam Ajanthan, Sameera Ramasinghe, Gil Avraham, Yan Zuo, and Alexander Long. Mo-
 502 mentum look-ahead for asynchronous distributed low-communication training. In *ICLR Workshop*
 503 *on Modularity for Collaborative, Decentralized, and Continual Deep Learning, 2025b*.

504 Mahmoud Assran, Arda Aytakin, Hamid Reza Feyzmahdavian, Mikael Johansson, and Michael G
 505 Rabbat. Advances in asynchronous parallel and distributed optimization. *Proceedings of the*
 506 *IEEE*, 108(11):2013–2031, 2020.

507

508 Saar Barkai, Ido Hakimi, and Assaf Schuster. Gap aware mitigation of gradient staleness. *arXiv*
 509 *preprint arXiv:1909.10802*, 2019.

510 Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar.
 511 signsgd: Compressed optimisation for non-convex problems. In *International Conference on*
 512 *Machine Learning*, pp. 560–569. PMLR, 2018.

513

514 Matt Beton, Matthew Reed, Seth Howes, Alex Cheema, and Mohamed Baoumy. Improving the
 515 efficiency of distributed training using sparse parameter averaging. In *ICLR 2025 Workshop on*
 516 *Modularity for Collaborative, Decentralized, and Continual Deep Learning, 2025*.

517 Léon Bottou, Frank E Curtis, and Jorge Nocedal. Optimization methods for large-scale machine
 518 learning. *SIAM review*, 60(2):223–311, 2018.

519

520 Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, Jonathan Eckstein, et al. Distributed optimiza-
 521 tion and statistical learning via the alternating direction method of multipliers. *Foundations and*
 522 *Trends® in Machine learning*, 3(1):1–122, 2011.

523 Chi-Chung Chen, Chia-Lin Yang, and Hsiang-Yun Cheng. Efficient and robust parallel dnn training
 524 through model parallelism on multi-gpu platform. *arXiv preprint arXiv:1809.02839*, 2018.

525

526 Arthur Douillard, Qixuan Feng, Andrei A Rusu, Rachita Chhaparia, Yani Donchev, Adhiguna
 527 Kuncoro, Marc’Aurelio Ranzato, Arthur Szlam, and Jiajun Shen. Diloco: Distributed low-
 528 communication training of language models. *arXiv preprint arXiv:2311.08105*, 2023.

529 Arthur Douillard, Yanislav Donchev, Keith Rush, Satyen Kale, Zachary Charles, Zachary Garrett,
 530 Gabriel Teston, Dave Lacey, Ross McIlroy, Jiajun Shen, et al. Streaming diloco with overlapping
 531 communication: Towards a distributed free lunch. *arXiv preprint arXiv:2501.18512*, 2025.

532

533 Timothy Dozat. Incorporating nesterov momentum into adam. *ICLR Workshop*, 2016.

534 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
 535 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
 536 *arXiv preprint arXiv:2407.21783*, 2024.

537

538 Louis Fournier, Adel Nabli, Masih Aminbeidokhti, Marco Pedersoli, Eugene Belilovsky, and
 539 Edouard Oyallon. Wash: Train your ensemble with communication-efficient weight shuffling,
 then average. *arXiv preprint arXiv:2405.17517*, 2024.

- 540 Aaron Gokaslan, Vanya Cohen, Ellie Pavlick, and Stefanie Tellex. Openwebtext corpus. [http://](http://SkyLion007.github.io/OpenWebTextCorpus)
541 SkyLion007.github.io/OpenWebTextCorpus, 2019.
- 542
- 543 P Goyal. Accurate, large minibatch sgd: training imagenet in 1 hour. *arXiv preprint*
544 *arXiv:1706.02677*, 2017.
- 545 Lei Guan, Wotao Yin, Dongsheng Li, and Xicheng Lu. Xpipe: Efficient pipeline model parallelism
546 for multi-gpu dnn training. *arXiv preprint arXiv:1911.04610*, 2019.
- 547
- 548 Lei Guan, Dong-Sheng Li, Ji-Ye Liang, Wen-Jian Wang, Ke-Shi Ge, and Xi-Cheng Lu. Advances of
549 pipeline model parallelism for deep learning training: an overview. *Journal of Computer Science*
550 *and Technology*, 39(3):567–584, 2024.
- 551 Ido Hakimi, Saar Barkai, Moshe Gabel, and Assaf Schuster. Taming momentum in a distributed
552 asynchronous environment. *arXiv preprint arXiv:1907.11612*, 2019.
- 553
- 554 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza
555 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Train-
556 ing compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- 557
- 558 Yanping Huang, Youlong Cheng, Ankur Bapna, Orhan Firat, Dehao Chen, Mia Chen, Hyoungho
559 Lee, Jiquan Ngiam, Quoc V Le, Yonghui Wu, et al. Gpipe: Efficient training of giant neural
560 networks using pipeline parallelism. *Advances in neural information processing systems*, 32,
561 2019.
- 562 Satyen Kale, Arthur Douillard, and Yanislav Donchev. Eager updates for overlapped communication
563 and computation in diloco. *arXiv preprint arXiv:2502.12996*, 2025.
- 564 Andrej Karpathy. NanoGPT. <https://github.com/karpathy/nanoGPT>, 2022.
- 565
- 566 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convo-
567 lutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017.
- 568
- 569 Sunwoo Lee, Anit Kumar Sahu, Chaoyang He, and Salman Avestimehr. Partial model averaging in
570 federated learning: Performance guarantees and benefits. *Neurocomputing*, 556:126647, 2023.
- 571 Shen Li, Yanli Zhao, Rohan Varma, Omkar Salpekar, Pieter Noordhuis, Teng Li, Adam Paszke, Jeff
572 Smith, Brian Vaughan, Pritam Damania, et al. Pytorch distributed: Experiences on accelerating
573 data parallel training. *arXiv preprint arXiv:2006.15704*, 2020.
- 574
- 575 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,
576 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint*
577 *arXiv:2412.19437*, 2024a.
- 578 Bo Liu, Rachita Chhaparia, Arthur Douillard, Satyen Kale, Andrei A Rusu, Jiajun Shen, Arthur
579 Szlam, and Marc’Aurelio Ranzato. Asynchronous local-sgd training for language modeling.
580 *arXiv preprint arXiv:2401.09135*, 2024b.
- 581
- 582 I Loshchilov. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- 583
- 584 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguerre y Arcas.
585 Communication-efficient learning of deep networks from decentralized data. In *Artificial intelli-*
586 *gence and statistics*, pp. 1273–1282. PMLR, 2017.
- 587
- 588 Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture
589 models, 2016.
- 589
- 590 Konstantin Mishchenko, Francis Bach, Mathieu Even, and Blake Woodworth. Asynchronous sgd
591 beats minibatch sgd under arbitrary delays. *URL <https://arxiv.org/abs/2206.07638>*, 2(6):7, 2022.
- 592
- 593 Ioannis Mitliagkas, Ce Zhang, Stefan Hadjis, and Christopher Ré. Asynchrony begets momentum,
with an application to deep learning. In *2016 54th Annual Allerton Conference on Communica-*
tion, Control, and Computing (Allerton), pp. 997–1004. IEEE, 2016.

- 594 Deepak Narayanan, Aaron Harlap, Amar Phanishayee, Vivek Seshadri, Nikhil R Devanur, Gre-
595 gory R Ganger, Phillip B Gibbons, and Matei Zaharia. Pipedream: Generalized pipeline par-
596 allelism for dnn training. In *Proceedings of the 27th ACM symposium on operating systems*
597 *principles*, pp. 1–15, 2019.
- 598
599 Deepak Narayanan, Amar Phanishayee, Kaiyu Shi, Xie Chen, and Matei Zaharia. Memory-efficient
600 pipeline-parallel dnn training. In *International Conference on Machine Learning*, pp. 7937–7947.
601 PMLR, 2021a.
- 602
603 Deepak Narayanan, Mohammad Shoeybi, Jared Casper, Patrick LeGresley, Mostofa Patwary, Vi-
604 jay Korthikanti, Dmitri Vainbrand, Prethvi Kashinkunti, Julie Bernauer, Bryan Catanzaro, et al.
605 Efficient large-scale language model training on gpu clusters using megatron-lm. In *Proceed-*
606 *ings of the International Conference for High Performance Computing, Networking, Storage and*
607 *Analysis*, pp. 1–15, 2021b.
- 608
609 Yurii Nesterov. A method for solving the convex programming problem with convergence rate o
610 $(1/k^2)$. In *Dokl akad nauk Sssr*, volume 269, pp. 543, 1983.
- 611
612 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
613 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-
614 performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- 615
616 Guilherme Penedo, Hynek Kydlíček, Anton Lozhkov, Margaret Mitchell, Colin A Raffel, Leandro
617 Von Werra, Thomas Wolf, et al. The fineweb datasets: Decanting the web for the finest text data
618 at scale. *Advances in Neural Information Processing Systems*, 37:30811–30849, 2024.
- 619
620 Penghui Qi, Xinyi Wan, Guangxing Huang, and Min Lin. Zero bubble pipeline parallelism. *arXiv*
621 *preprint arXiv:2401.10241*, 2023.
- 622
623 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language
624 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 625
626 Sameera Ramasinghe, Thalaiyasingam Ajanthan, Gil Avraham, Yan Zuo, and Alexander Long. Pro-
627 tocol models: Scaling decentralized training with communication-efficient model parallelism,
628 2025. URL <https://arxiv.org/abs/2506.01260>.
- 629
630 Benjamin Recht, Christopher Re, Stephen Wright, and Feng Niu. Hogwild!: A lock-free approach
631 to parallelizing stochastic gradient descent. *Advances in neural information processing systems*,
632 24, 2011.
- 633
634 Sashank Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný,
635 Sanjiv Kumar, and H Brendan McMahan. Adaptive federated optimization. *arXiv preprint*
636 *arXiv:2003.00295*, 2020.
- 637
638 Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathemati-*
639 *cal statistics*, pp. 400–407, 1951.
- 640
641 Herbert Robbins and David Siegmund. A convergence theorem for non negative almost super-
642 martingales and some applications. In *Optimizing methods in statistics*, pp. 233–257. Elsevier,
643 1971.
- 644
645 Max Ryabinin, Eduard Gorbunov, Vsevolod Plokhotnyuk, and Gennady Pekhimenko. Moshpit sgd:
646 Communication-efficient decentralized training on heterogeneous unreliable devices. *Advances*
647 *in Neural Information Processing Systems*, 34:18195–18211, 2021.
- 648
649 Max Ryabinin, Tim Dettmers, Michael Diskin, and Alexander Borzunov. Swarm parallelism: Train-
650 ing large models can be surprisingly communication-efficient. In *International Conference on*
651 *Machine Learning*, pp. 29416–29440. PMLR, 2023.
- 652
653 Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan
654 Catanzaro. Megatron-lm: Training multi-billion parameter language models using model par-
655 allelism. *arXiv preprint arXiv:1909.08053*, 2019.

- 648 Sebastian U Stich and Sai Praneeth Karimireddy. The error-feedback framework: Better rates for
649 sgd with delayed gradients and compressed communication. *arXiv preprint arXiv:1909.05350*,
650 2019.
- 651 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
652 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-*
653 *tion processing systems*, 30, 2017.
- 654 Thijs Vogels, Sai Praneeth Karimireddy, and Martin Jaggi. Powersgd: Practical low-rank gradient
655 compression for distributed optimization. *Advances in Neural Information Processing Systems*,
656 32, 2019.
- 657 Jialei Wang, Mladen Kolar, Nathan Srebro, and Tong Zhang. Efficient distributed learning with
658 sparsity. In *International conference on machine learning*, pp. 3636–3645. PMLR, 2017.
- 659 Jue Wang, Yucheng Lu, Binhang Yuan, Beidi Chen, Percy Liang, Christopher De Sa, Christopher
660 Re, and Ce Zhang. Cocktailsgd: Fine-tuning foundation models over 500mbps networks. In
661 *International Conference on Machine Learning*, pp. 36058–36076. PMLR, 2023.
- 662 Jianqiao Wangni, Jialei Wang, Ji Liu, and Tong Zhang. Gradient sparsification for communication-
663 efficient distributed optimization. *Advances in Neural Information Processing Systems*, 31, 2018.
- 664 Bowen Yang, Jian Zhang, Jonathan Li, Christopher Ré, Christopher Aberger, and Christopher De Sa.
665 Pipemare: Asynchronous pipeline parallel dnn training. *Proceedings of Machine Learning and*
666 *Systems*, 3:269–296, 2021.
- 667 Hao Yu, Rong Jin, and Sen Yang. On the linear speedup analysis of communication efficient mo-
668 mentum sgd for distributed non-convex optimization. In *International Conference on Machine*
669 *Learning*, pp. 7184–7193. PMLR, 2019.
- 670 Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuandong
671 Tian. Galore: Memory-efficient llm training by gradient low-rank projection. *arXiv preprint*
672 *arXiv:2403.03507*, 2024.
- 673 Yanli Zhao, Andrew Gu, Rohan Varma, Liang Luo, Chien-Chin Huang, Min Xu, Less Wright,
674 Hamid Shojanazeri, Myle Ott, Sam Shleifer, et al. Pytorch fsdp: experiences on scaling fully
675 sharded data parallel. *arXiv preprint arXiv:2304.11277*, 2023.
- 676 Shuxin Zheng, Qi Meng, Taifeng Wang, Wei Chen, Nenghai Yu, Zhi-Ming Ma, and Tie-Yan Liu.
677 Asynchronous stochastic gradient descent with delay compensation. In *International conference*
678 *on machine learning*, pp. 4120–4129. PMLR, 2017.
- 679 Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and
680 Sanja Fidler. Aligning books and movies: Towards story-like visual explanations by watching
681 movies and reading books. In *The IEEE International Conference on Computer Vision (ICCV)*,
682 December 2015.
- 683 Yan Zuo, Gil Avraham, Thalaiyasingam Ajanthan, Sameera Ramasinghe, and Alexander Long. Ex-
684 ploring asynchronism in swarm parallelism. In *ICLR Workshop on Modularity for Collaborative,*
685 *Decentralized, and Continual Deep Learning*, 2025.
- 686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

A THEORETICAL INSIGHTS

We restate the problem setup, introduce some notations, and then turn to the proofs.

A.1 SETUP AND NOTATIONS

We consider the DP setup without PP for simplified theoretical analysis. In this, the consensus objective takes the following form:

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}; \mathcal{D}) &:= \min_{\mathbf{w}_i \in \mathbb{R}^d} \sum_{i=1}^m f(\mathbf{w}_i; \mathcal{D}_i), \\ \text{s.t. } \mathbf{w}_i &= \mathbf{w}, \quad \forall i \in \{1, \dots, m\}. \end{aligned} \quad (11)$$

Here, $i \in \{1, \dots, m\}$ denotes the workers, $\mathbf{w}_i, \mathcal{D}_i$ denote the model weights and data chunk for worker i , and f is the loss function. Typically, \mathcal{D}_i is an i.i.d. subset of \mathcal{D} .

Let us first define some notations:

$$\begin{aligned} \bar{\mathbf{w}}^t &:= \frac{1}{m} \sum_{i=1}^m \mathbf{w}_i^t, && \text{averaged weights,} \\ \mathbf{g}_i^t &:= \nabla f_i(\mathbf{w}_i^t; \mathcal{B}_i^t), && \text{minbatch gradient,} \\ \bar{\mathbf{g}}^t &:= \frac{1}{m} \sum_{i=1}^m \mathbf{g}_i^t, && \text{average of gradients,} \\ \|\Delta^t\|^2 &:= \sum_{i=1}^m \|\mathbf{w}_i^t - \bar{\mathbf{w}}^t\|^2, && \text{consensus error.} \end{aligned} \quad (12)$$

Here, $\|\Delta^t\|^2$ denotes the post-averaging consensus error at iteration $t-1$, and the corresponding pre-averaging error is denoted as $\|\hat{\Delta}^t\|^2$. Analogously, $\Delta_i^t := \mathbf{w}_i^t - \bar{\mathbf{w}}^t$.

We are now ready to prove that sparse averaging ensures consensus across the weights of all workers on expectation.

A.2 SPARSE AVERAGING ENSURES CONSENSUS

Lemma 1. *Let p be the probability of an element $\mu \in \mathcal{S}^t$ independent of others, then sparse averaging shrinks the consensus error by a factor of $(1-p)$ on expectation, i.e.,*

$$\mathbb{E}[\|\Delta^t\|^2] = (1-p) \mathbb{E}[\|\hat{\Delta}^t\|^2].$$

Proof. By definition of sparse averaging, $w_{i;\mu}^{t+1} = \frac{1}{m} \sum_i \hat{w}_{i;\mu}^t$ for all $\mu \in \mathcal{S}^t$. Therefore,

$$\begin{aligned} \mathbb{E}[\|\Delta^t\|^2] &= \sum_{i=1}^m \sum_{\mu=1}^d \mathbb{E}[(\Delta_{i;\mu}^t)^2], && \text{linearity of expectation,} \\ &= \sum_{i=1}^m \sum_{\mu=1}^d (\hat{\Delta}_{i;\mu}^t)^2 \mathbb{E}[\mathbb{I}[\mu \notin \mathcal{S}^t]], && \mathbb{I}[\cdot] \text{ is the indicator,} \\ &= (1-p) \mathbb{E}[\|\hat{\Delta}^t\|^2]. && P[\mu \notin \mathcal{S}^t] = 1-p. \end{aligned} \quad (13)$$

□

Theorem 3. *Let f be a L -smooth function, the stochastic gradient \mathbf{g}_i^t be an unbiased estimate of ∇f and have bounded variance σ^2 , and $p > 0$ be the averaging probability for an element μ , then, for an appropriate choice of learning rate $\eta_t > 0$ the consensus error diminishes on expectation, i.e.,*

$$\lim_{t \rightarrow \infty} \mathbb{E}[\|\Delta^t\|^2] = 0.$$

756 *Proof.* Let us expand $\hat{\Delta}_i^{t+1}$:

$$\begin{aligned}
757 \hat{\Delta}_i^{t+1} &= \hat{\mathbf{w}}_i^{t+1} - \frac{1}{m} \sum_i \hat{\mathbf{w}}_i^{t+1}, & (14) \\
758 & \\
759 &= \mathbf{w}_i^t - \eta_t \mathbf{g}_i^t - \frac{1}{m} \sum_i (\mathbf{w}_i^t - \eta_t \mathbf{g}_i^t), & \text{local update,} \\
760 & \\
761 &= \Delta_i^t - \eta_t (\mathbf{g}_i^t - \bar{\mathbf{g}}^t). \\
762 & \\
763 & \\
764 &
\end{aligned}$$

765 Now, the pre-averaging consensus error can be written as:

$$\begin{aligned}
766 \|\hat{\Delta}^{t+1}\|^2 &= \sum_i \|\Delta_i^t - \eta_t (\mathbf{g}_i^t - \bar{\mathbf{g}}^t)\|^2, & (15) \\
767 & \\
768 &= \sum_i \left(\|\Delta_i^t\|^2 - 2\eta_t \Delta_i^t \cdot (\mathbf{g}_i^t - \bar{\mathbf{g}}^t) + \eta_t^2 \|\mathbf{g}_i^t - \bar{\mathbf{g}}^t\|^2 \right), & \cdot \text{ is the dot product,} \\
769 & \\
770 &= \sum_i \left(\|\Delta_i^t\|^2 - 2\eta_t \Delta_i^t \cdot \mathbf{g}_i^t + \eta_t^2 \|\mathbf{g}_i^t - \bar{\mathbf{g}}^t\|^2 \right), & \sum_i \Delta_i^t = 0. \\
771 & \\
772 & \\
773 & \\
774 &
\end{aligned}$$

775 We now bound each of the terms following techniques similar to that of (Yu et al., 2019). Consider

776 $\|\mathbf{g}_i^t - \bar{\mathbf{g}}^t\|^2$:

$$\begin{aligned}
777 \|\mathbf{g}_i^t - \bar{\mathbf{g}}^t\|^2 &= \|\mathbf{g}_i^t - \nabla f(\mathbf{w}_i^t) + \nabla f(\mathbf{w}_i^t) - \bar{\nabla} f(\mathbf{w}_i^t) + \bar{\nabla} f(\mathbf{w}_i^t) - \bar{\mathbf{g}}^t\|^2, & (16) \\
778 & \\
779 &\leq 3 \|\mathbf{g}_i^t - \nabla f(\mathbf{w}_i^t)\|^2 + 3 \|\nabla f(\mathbf{w}_i^t) - \bar{\nabla} f(\mathbf{w}_i^t)\|^2 + 3 \|\bar{\nabla} f(\mathbf{w}_i^t) - \bar{\mathbf{g}}^t\|^2, \\
780 & \\
781 &
\end{aligned}$$

782 where $\bar{\nabla} f(\mathbf{w}_i^t) = \frac{1}{m} \sum_i \nabla f(\mathbf{w}_i^t)$, and the second step is due to $(a + b + c)^2 \leq 3(a^2 + b^2 + c^2)$.

783 Each of these terms can be bounded as follows: 1) using bounded variance,

$$\begin{aligned}
784 \sum_i \|\mathbf{g}_i^t - \nabla f(\mathbf{w}_i^t)\|^2 &\leq m \sigma^2. & (17) \\
785 & \\
786 & \\
787 &
\end{aligned}$$

788 2) using L -smoothness and triangle inequality,

$$\begin{aligned}
789 \sum_i \|\nabla f(\mathbf{w}_i^t) - \bar{\nabla} f(\mathbf{w}_i^t)\|^2 &= \sum_i \|\nabla f(\mathbf{w}_i^t) - \nabla f(\bar{\mathbf{w}}^t) + \nabla f(\bar{\mathbf{w}}^t) - \bar{\nabla} f(\mathbf{w}_i^t)\|^2, & (18) \\
790 & \\
791 &= \sum_i 2 \|\nabla f(\mathbf{w}_i^t) - \nabla f(\bar{\mathbf{w}}^t)\|^2 + 2 \left\| \frac{1}{m} \sum_i \nabla f(\bar{\mathbf{w}}^t) - \frac{1}{m} \sum_i \nabla f(\mathbf{w}_i^t) \right\|^2, \\
792 & \\
793 &\leq 2L^2 \sum_i \|\mathbf{w}_i^t - \bar{\mathbf{w}}^t\|^2 + \frac{2L^2}{m} \sum_k \sum_i \|\mathbf{w}_i^t - \bar{\mathbf{w}}^t\|^2, \\
794 & \\
795 &= 4L^2 \|\Delta^t\|^2. \\
796 & \\
797 & \\
798 &
\end{aligned}$$

799 3) using triangle inequality and bounded variance,

$$\begin{aligned}
800 \sum_i \|\bar{\nabla} f(\mathbf{w}_i^t) - \bar{\mathbf{g}}^t\|^2 &\leq \frac{1}{m} \sum_i \|\mathbf{g}_i^t - \nabla f(\mathbf{w}_i^t)\|^2 \leq \sigma^2. & (19) \\
801 & \\
802 & \\
803 &
\end{aligned}$$

804 Altogether for some constant $C > 0$, we can write:

$$\begin{aligned}
805 \mathbb{E} \left[\sum_i \|\mathbf{g}_i^t - \bar{\mathbf{g}}^t\|^2 \right] &\leq C \left(m \sigma^2 + L^2 \mathbb{E} [\|\Delta^t\|^2] \right). & (20) \\
806 & \\
807 & \\
808 & \\
809 &
\end{aligned}$$

Now, consider the term $\mathbb{E}[\sum_i \Delta_i^t \cdot \mathbf{g}_i^t]$:

$$\begin{aligned}
\mathbb{E}\left[\sum_i \Delta_i^t \cdot \mathbf{g}_i^t\right] &= \mathbb{E}\left[\sum_i \Delta_i^t \cdot \nabla f(\mathbf{w}_i^t)\right], & \mathbb{E}[\mathbf{g}_i^t \mid \Delta_i^t] &= \mathbb{E}[\mathbf{g}_i^t \mid \mathbf{w}_i^t] = \nabla f(\mathbf{w}_i^t), \\
&= \mathbb{E}\left[\sum_i \Delta_i^t \cdot (\nabla f(\mathbf{w}_i^t) - \nabla f(\bar{\mathbf{w}}^t))\right], & \sum_i \Delta_i^t &= 0, \\
&\leq L \mathbb{E}\left[\sum_i \|\Delta_i^t\| \cdot \|\Delta_i^t\|\right], & L\text{-smooth, } \Delta_i^t &= \mathbf{w}_i^t - \bar{\mathbf{w}}^t, \\
&= L \mathbb{E}\left[\|\Delta^t\|^2\right].
\end{aligned} \tag{21}$$

Putting everything together,

$$\|\hat{\Delta}^{t+1}\|^2 \leq (1 + 2\eta_t L) \mathbb{E}\left[\|\Delta^t\|^2\right] + \eta_t^2 C \left(m\sigma^2 + L^2 \mathbb{E}\left[\|\Delta^t\|^2\right]\right). \tag{22}$$

From Lemma 1,

$$\begin{aligned}
\mathbb{E}\left[\|\Delta^{t+1}\|^2\right] &= (1 - p) \mathbb{E}\left[\|\hat{\Delta}^{t+1}\|^2\right], \\
&\leq (1 - p) (1 + 2\eta_t L) \mathbb{E}\left[\|\Delta^t\|^2\right] + \mathcal{O}(\eta_t^2).
\end{aligned} \tag{23}$$

Note, η_t can be chosen such that the quadratic term is vanishes, *i.e.*, $\eta_t > 0$, $\sum_t \eta_t = \infty$, and $\sum_t \eta_t^2 < \infty$ (Robbins & Monro, 1951), and the coefficient $(1 - p) (1 + 2\eta_t L)$ is strictly less than 1, *i.e.*, $\eta_t < \frac{p}{2(1-p)L}$.

This yields a contraction and ensures $\lim_{t \rightarrow \infty} \mathbb{E}\left[\|\Delta^t\|^2\right] = 0$. \square

This proves that sparse averaging can lead to consensus among workers, in the sense that, on expectation, all weight vectors converge to their average. This, together with the standard convergence proof of SGD (Bottou et al., 2018) guarantees that sparse averaging converges to a fixed point of Eq. (11).

A.3 EMA BASED DELAY CORRECTION ENSURES CONSENSUS

We consider a *homogeneous setup* where all workers are initialized to the same point, the data chunks are i.i.d., and the optimizer parameters are identical. In this, we first show that the expected staleness $\mathbb{E}[\hat{\mathbf{w}}_i^t - \hat{\mathbf{w}}_i^{t-\tau}]$ can be independently estimated in the worker i . Then, we prove that the drift of \mathbf{d}_i^t between different workers diminishes, ensuring consensus. In this section, with a slight abuse of notation, we define $\bar{\mathbf{w}}^t = \frac{1}{m} \sum_i \hat{\mathbf{w}}_i^t$.

Lemma 2. *In a homogeneous setup as defined above, the expected value of the weight drift is independent of the worker, *i.e.*, $\mathbb{E}[\hat{\mathbf{w}}_i^t - \hat{\mathbf{w}}_i^{t-\tau}] = \mathbf{D}^t$.*

Proof. Considering a single local update:

$$\mathbb{E}[\hat{\mathbf{w}}_i^{k+1} - \mathbf{w}_i^k] = -\mathbb{E}[\eta_k \mathbf{g}_i^k] = -\eta_k \nabla f(\mathbf{w}_i^k), \quad \mathbf{g}_i^k \text{ is unbiased.} \tag{24}$$

This is independent and identical for each worker if $\mathbf{w}_i^k = \mathbf{w}^k$ for all i . This argument can be extended to multiple steps for a homogeneous setup due to the same initialization, i.i.d. data samples, and identical optimizer hyperparameters. Intuitively, one can see that a particular weight trajectory is equally probable for all workers in a homogeneous setup.

Additionally, the EMA update parameters (\mathbf{d}_i^0 and λ_t) are also identical across workers. Therefore, the average weight drift $\mathbb{E}[\hat{\mathbf{w}}_i^t - \hat{\mathbf{w}}_i^{t-\tau}] = \mathbb{E}[\hat{\mathbf{w}}_i^t - \mathbf{w}_i^{t-\tau}] + \mathbb{E}[\mathbf{w}_i^{t-\tau} - \hat{\mathbf{w}}_i^{t-\tau}]$ is independent of the worker i . \square

Theorem 4. Consider a homogeneous setup where the average staleness \mathbf{D}^t is bounded and its drift $\boldsymbol{\alpha}_t := \mathbf{D}^t - \mathbf{D}^{t-1}$ is diminishing, i.e., $\lim_{t \rightarrow \infty} \|\boldsymbol{\alpha}_t\| = 0$. Then, the EMA based delay correction with λ_t satisfying $\sum_t \lambda_t = \infty$, $\sum_t \lambda_t^2 < \infty$, and $\sum_t \frac{\|\boldsymbol{\alpha}_t\|^2}{\lambda_t} < \infty$ ensures consensus, i.e., $\lim_{t \rightarrow \infty} \mathbb{E}[\|\boldsymbol{\Delta}^t\|^2] = 0$.

Proof. Now consider $\boldsymbol{\Delta}_i^t$:

$$\begin{aligned} \boldsymbol{\Delta}_i^t &= \mathbf{w}_i^t - \bar{\mathbf{w}}_i^t, \\ &= \bar{\mathbf{w}}^{t-\tau} + \mathbf{d}_i^t - \frac{1}{m} \sum_i (\bar{\mathbf{w}}^{t-\tau} + \mathbf{d}_i^t), \quad \text{EMA}, \\ &= \mathbf{d}_i^t - \bar{\mathbf{d}}_i^t. \end{aligned} \tag{25}$$

Note that, $\bar{\mathbf{d}}_i^t = \frac{1}{m} \sum_i \mathbf{d}_i^t$ is the empirical estimate of $\mathbb{E}[\mathbf{d}_i^t]$. Following the arguments from Lemma 2, we can see that $\mathbb{E}[\mathbf{d}_i^t] = \mathbf{D}^t$ for all workers. If \mathbf{D}^t is a constant, then the standard stochastic approximation theory (Robbins & Monro, 1951) yields the desired result. However, since \mathbf{D}^t varies with time, we need an additional assumption on the interplay between the drift and the EMA coefficient, that $\frac{\|\boldsymbol{\alpha}_t\|^2}{\lambda_t} \rightarrow 0$ and carefully bound the consensus error.

Substituting $\bar{\mathbf{d}}_i^t = \mathbf{D}^t$, together with the EMA update:

$$\begin{aligned} \boldsymbol{\Delta}_i^t &= \mathbf{d}_i^t - \mathbf{D}^t, \\ &= (1 - \lambda_t) \mathbf{d}_i^{t-1} + \lambda_t (\hat{\mathbf{w}}_i^t - \hat{\mathbf{w}}_i^{t-\tau}) - \mathbf{D}^t, \\ &= (1 - \lambda_t) (\mathbf{d}_i^{t-1} - \mathbf{D}^{t-1} + \mathbf{D}^{t-1} - \mathbf{D}^t) + \lambda_t (\hat{\mathbf{w}}_i^t - \hat{\mathbf{w}}_i^{t-\tau} - \mathbf{D}^t), \\ &= (1 - \lambda_t) \boldsymbol{\Delta}_i^{t-1} - (1 - \lambda_t) \boldsymbol{\alpha}_t + \lambda_t \boldsymbol{\xi}_i^t, \quad \boldsymbol{\xi}_i^t = \hat{\mathbf{w}}_i^t - \hat{\mathbf{w}}_i^{t-\tau} - \mathbf{D}^t. \end{aligned} \tag{26}$$

Now square both sides:

$$\begin{aligned} \|\boldsymbol{\Delta}_i^t\|^2 &= (1 - \lambda_t)^2 \|\boldsymbol{\Delta}_i^{t-1}\|^2 + (1 - \lambda_t)^2 \|\boldsymbol{\alpha}_t\|^2 + \lambda_t^2 \|\boldsymbol{\xi}_i^t\|^2 \\ &\quad + 2(1 - \lambda_t)\lambda_t \boldsymbol{\Delta}_i^{t-1} \cdot \boldsymbol{\xi}_i^t - 2(1 - \lambda_t)^2 \boldsymbol{\Delta}_i^{t-1} \cdot \boldsymbol{\alpha}_t - 2(1 - \lambda_t)\lambda_t \boldsymbol{\alpha}_t \cdot \boldsymbol{\xi}_i^t, \end{aligned} \tag{27}$$

where \cdot is the dot-product. Consider the first cross term:

$$\mathbb{E}[\boldsymbol{\Delta}_i^{t-1} \cdot \boldsymbol{\xi}_i^t] = \mathbb{E}[\boldsymbol{\Delta}_i^{t-1} \cdot \mathbb{E}[\boldsymbol{\xi}_i^t]] = 0, \quad \mathbb{E}[\boldsymbol{\xi}_i^t] = 0 \text{ due to Lemma 2.} \tag{28}$$

Using Young's inequality for the other cross terms, i.e., $2ab \leq \epsilon a^2 + \frac{1}{\epsilon} b^2$ for any $\epsilon > 0$, we can write:

$$\begin{aligned} \mathbb{E}[\|\boldsymbol{\Delta}_i^t\|^2] &\leq (1 - \lambda_t)^2 \mathbb{E}[\|\boldsymbol{\Delta}_i^{t-1}\|^2] + (1 - \lambda_t)^2 \|\boldsymbol{\alpha}_t\|^2 + \lambda_t^2 \mathbb{E}[\|\boldsymbol{\xi}_i^t\|^2] \\ &\quad + (1 - \lambda_t)^2 \left(\epsilon \mathbb{E}[\|\boldsymbol{\Delta}_i^{t-1}\|^2] + \frac{\|\boldsymbol{\alpha}_t\|^2}{\epsilon} \right) + (1 - \lambda_t)\lambda_t \left(\epsilon \mathbb{E}[\|\boldsymbol{\xi}_i^t\|^2] + \frac{\|\boldsymbol{\alpha}_t\|^2}{\epsilon} \right). \end{aligned} \tag{29}$$

The term $\mathbb{E}[\|\boldsymbol{\xi}_i^t\|^2]$ can be bounded due to bounded delay τ . By setting $\epsilon = \lambda_t$, the above can be simplified as:

$$\begin{aligned} \mathbb{E}[\|\boldsymbol{\Delta}_i^t\|^2] &\leq (1 - \lambda_t)^2 (1 + \lambda_t) \mathbb{E}[\|\boldsymbol{\Delta}_i^{t-1}\|^2] + \frac{A}{\lambda_t} \|\boldsymbol{\alpha}_t\|^2 \\ &\quad + B (\|\boldsymbol{\alpha}_t\|^2 + \lambda_t^2), \quad \text{for some constants } A, B > 0, \\ &\leq (1 - \lambda_t) \mathbb{E}[\|\boldsymbol{\Delta}_i^{t-1}\|^2] + \frac{A}{\lambda_t} \|\boldsymbol{\alpha}_t\|^2 + B (\|\boldsymbol{\alpha}_t\|^2 + \lambda_t^2), \quad 0 < \lambda_t < 1. \end{aligned} \tag{30}$$

Since $\frac{\|\boldsymbol{\alpha}_t\|^2}{\lambda_t}$, $\|\boldsymbol{\alpha}_t\|^2$, and λ_t^2 are diminishing, the above can be shown to be a contraction and therefore, $\mathbb{E}[\|\boldsymbol{\Delta}_i^t\|^2] \rightarrow 0$ following the classical stochastic approximation theory (Robbins & Monro, 1951; Robbins & Siegmund, 1971). Consequently, the consensus error vanishes:

$$\mathbb{E}[\|\boldsymbol{\Delta}^t\|^2] = \mathbb{E}\left[\sum_i \|\boldsymbol{\Delta}_i^t\|^2\right] = \sum_i \mathbb{E}[\|\boldsymbol{\Delta}_i^t\|^2] \rightarrow 0. \tag{31}$$

□

Note that the EMA update and the arguments in the theorem are elementwise, hence, they naturally extend to the PP with sparse averaging setup. This, together with Theorem 3 and the convergence proof of SGD, provides a theoretical justification for convergence for our method in the PP setup with sparse averaging on expectation, despite a fixed delay.

B EXPERIMENTS

B.1 EXPERIMENTAL SETUP

We evaluate on four large-scale language modeling datasets, namely, WikiText (WT) (Merity et al., 2016), BookCorpus (BC) (Zhu et al., 2015), OpenWebText (OWT) (Gokaslan et al., 2019), and FineWeb (FW) (Penedo et al., 2024). For WikiText, we utilize the predefined training and validation splits; for BookCorpus and OpenWebText, we randomly select 10% of the training set as the held-out validation set; and for FineWeb, we use the streaming feature in Huggingface datasets and hold out 10k samples in the stream for validation. Our architecture is based on NanoGPT (Karpathy, 2022) with no dropout. The base model has a context length of 1024, an embedding dimension of 768, 12 attention heads, and 12 layers, with approximately 163M parameters. We use the GPT2 tokenizer (Radford et al., 2019) and train the model from scratch. For configurations with 2 and 4 pipeline stages, equal number of layers are assigned to each stage, unless specified otherwise. For 8 stages, stages 2 – 5 are assigned 2 layers, and others have 1 layer each.

Across all experiments, we maintain a microbatch size of 8 per DP replica, a learning rate of $3e-4$, a weight decay of 0.01, and gradient clipping norm of 1. For experiments with asynchronous PP, NAdamW optimizer (Dozat, 2016) with momentum 0.99 is used as per (Ajanthan et al., 2025a). For synchronous PP experiments, GPipe (Huang et al., 2019) with AdamW optimizer (Loshchilov, 2017) is used, and the number of microbatches is set to 2. Each experiment is run for 30k iterations, with a linear warmup of 3k iterations starting from a learning rate of $1e-7$. Then, it is decayed to $3e-5$ following a cosine decay schedule. For our method, the EMA variable d_t^{\dagger} is initialized to zero.

For DiLoCo, the outer learning rate of 1 performs better (*i.e.*, averaging instead of outer optimization step) in our 2D mesh with AsyncPP, and the outer-update interval is set to 10 steps in our experiments.

1B Model. We maintain the number of stages at 4 with stage assignment of [1, 3, 4, 4] number of layers, but increase the embedding dimension to 2304, with 24 attention heads. The learning rate warmup step is adjusted to 6k for all methods and run for 100k iterations. All other hyperparameters are the same as the base model.

Varying Configurations. When changing a criterion, all other criteria are kept to the default values: `subset-size = 5%`, `async-delay = 10`, `avg-interval = 1`, `DP-replicas = 4`, `PP-stages = 4`. While varying the averaging interval, the asynchronous delay is set to 1, such that the effective delay is equal to the averaging interval.

B.2 ADDITIONAL RESULTS

We provide additional validation loss trajectories for various methods with synchronous PP in Fig. 6, on multiple datasets corresponding to the results in the main paper in Fig. 7, and for compute optimal training for the base model in Fig. 8. Furthermore, we provide consensus error plots for our method confirming the theory in Fig. 9, validation loss vs time plot for the 1B model in Fig. 10.

Gain in Wall-clock Time. Since we simulate the asynchronous DP setup via buffering, and due to implementation differences between our method and FullSync, their practical wall-clock times are not comparable. However, the theoretical gain in wall-clock time per iteration due to asynchronous sparse averaging is $O(d^2/B)$, where d is the embedding dimension, and B is the bandwidth, as we fully eliminate DP overhead. Since the staleness is in update steps, the ratio between allowed time and the data transfer volume is $O(Bd^3\tau/pd^2) = O(Bd\tau/p)$, where τ is the allowed delay, p is the subset size, and note the compute time is approximately cubic in d (Ryabinin et al., 2023). Therefore, our method scales favourably for large models. Note, AsyncPP (Ajanthan et al., 2025a) further improves this wall-clock time gain.

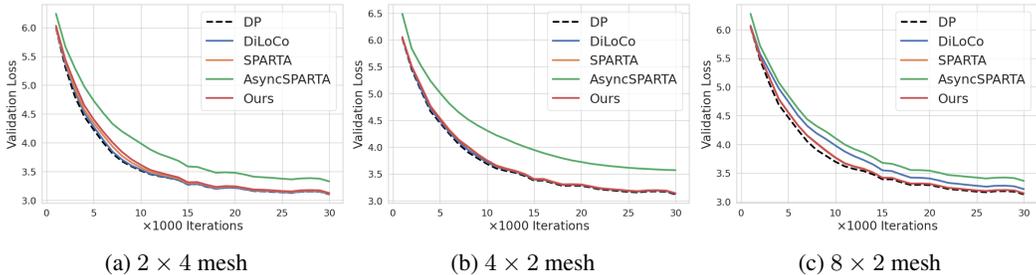


Figure 6: Results on WikiText with varying mesh configurations with synchronous PP. In all scenarios, our method matches the performance of the fully synchronous method, same as the case with AsyncPP.

For completeness, we show the empirical time savings that can be achieved by our asynchronous sparse averaging in Table 2. This clearly shows just by making sparse averaging asynchronous with 1-step delay, we can achieve a speed-up of $1.5\times - 3.7\times$ depending on the mesh configuration, with larger meshes yielding better speed-ups.

Heterogeneous Setup So far, our experiments have been on a homogeneous setup where the compute capability of the devices is the same. To stress test our method, we simulate a heterogeneous setup for a 44 mesh by varying device speeds, namely, $H2\times: [2, 1, 2, 3]$, $H5\times: [5, 3, 1, 11]$, and $H10\times: [10, 8, 1, 21]$, along with the homogeneous setup: $H1\times: [1, 1, 1, 1]$, with the provided relative device speeds. To ensure each device initiates the averaging operation (*i.e.*, all-reduce) at approximately the same time, we set the averaging interval proportional to the device speed, *i.e.*, faster devices will perform more iterations between an averaging step. This means, model replicas at different training stages are sparsely averaged with a delay due to asynchronous DP. This is analogous to dynamic local updates in (Liu et al., 2024b). For fair comparison, we fix the total number of iterations across all replicas to be 120k. As shown in Fig. 11, even with up to $10\times$ difference in device speeds (denoted $H10\times$), and delayed sparse averaging of replicas at different training stages, the degradation of validation loss remains small. This degradation is negligible when accounting for the wall-clock speed-up, which scales with device speeds due to the absence of DP communication overhead.

Comparison with Other Communication Efficient Methods. We compare with some existing methods that compress the DP communication via quantization and TopK sample in Fig. 12. The quantization results show sparse averaging methods (even with delay) are more robust to quantization, allowing further reduction in DP communication requirements. Finally, we compare with the concurrent work of eager updates for DiLoCo (Kale et al., 2025) in Fig. 13, showing that our method strictly generalizes it.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

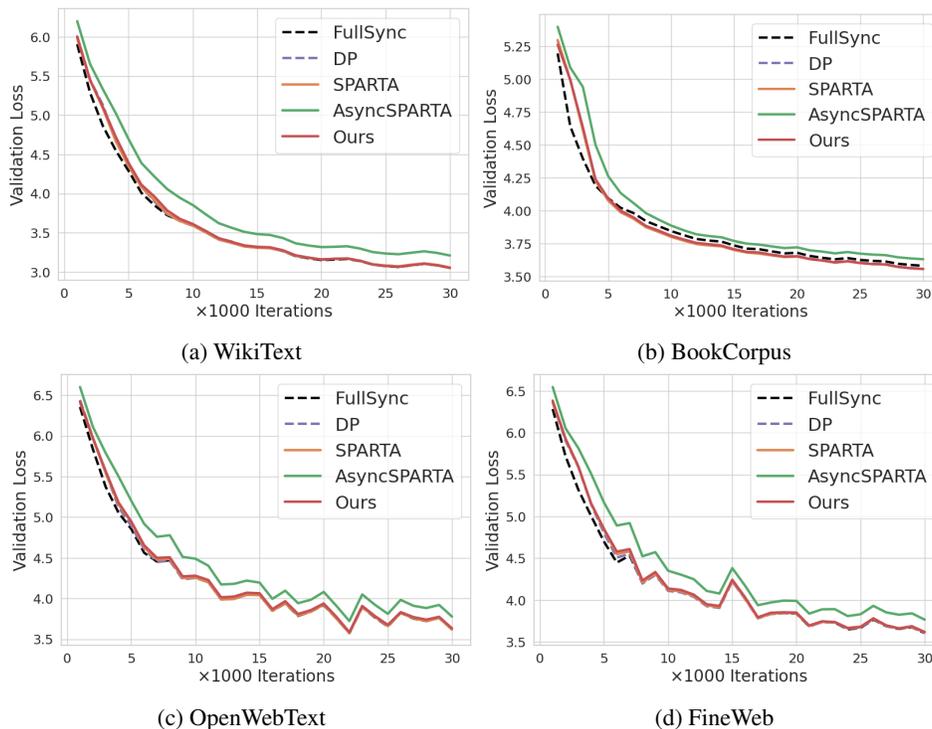


Figure 7: Results on different datasets for 4×2 mesh. Our method performs similarly to FullSync in all scenarios, demonstrating virtually no performance degradation due to staleness or sparse averaging.

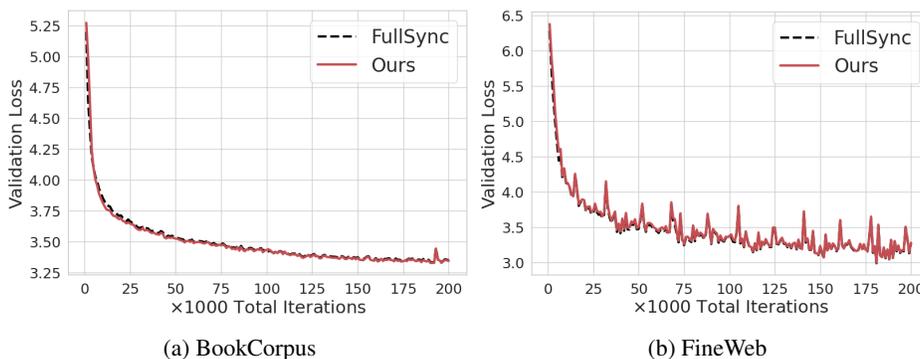


Figure 8: Compute optimal training (Hoffmann et al., 2022) for BookCorpus and FineWeb for the base model with 4×2 mesh. Our method is nearly identical to FullSync for longer training as well, validating its merits. The final validation perplexities are, for BookCorpus, FullSync: 28.02, and Ours: 27.86 and for FineWeb, FullSync: 19.92, and Ours: 20.10. The validation curve for FineWeb is noisy for both methods, probably due to the way the validation set is selected from the Huggingface stream.

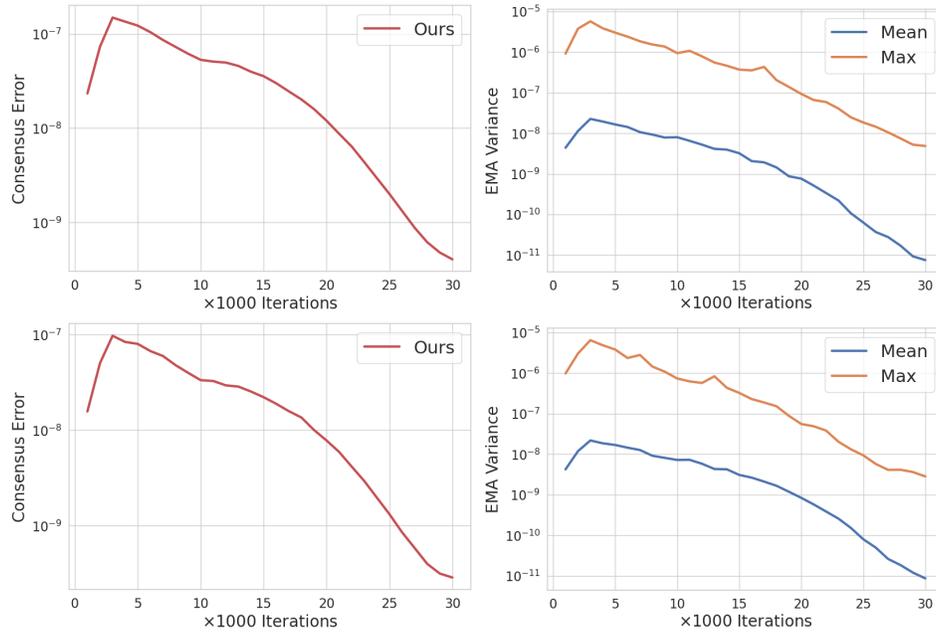


Figure 9: Mean consensus error $\frac{1}{md} \|\Delta^t\|^2$ as in Eq. (12) and variance between EMA estimates (\mathbf{d}_i^t) in each replica, for our method on the 2×4 (top) and 4×2 mesh (bottom). For EMA, the mean and max across the model dimension are shown. Results perfectly align with the theory that independent EMA estimates in each replica converge to the expected value (i.e., variance vanishes), and the consensus error for our method diminishes to zero.

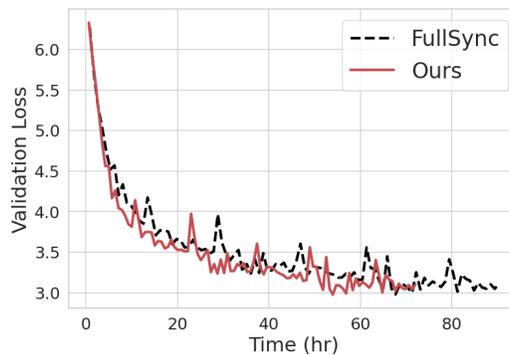


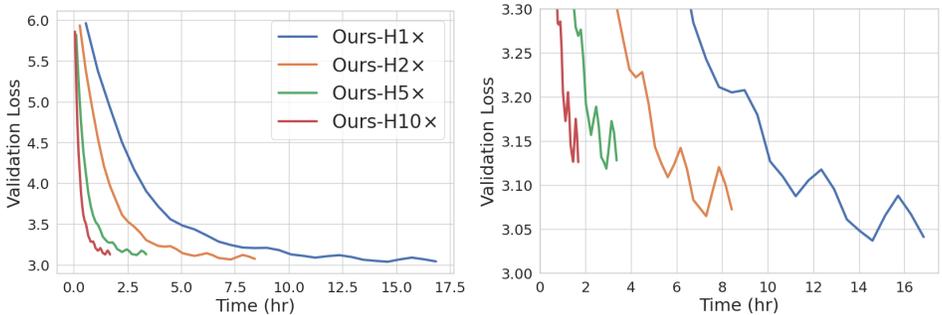
Figure 10: Validation loss vs time for 1B model. Even with suboptimal implementation, fast interconnects, and **not** considering the time gains due to asynchronous updates, FullSync is about 20% slower than our method.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146

PP × DP Mesh	AWS Instances	Communication Time for tail (ms)	Forward-Backward time for tail (ms)	Speed-up
4 × 2	1 × p4d.24	515 ± 35	320 ± 25	~ 2.6×
4 × 4	1 × p4d.24	920 ± 150	590 ± 35	~ 2.6×
4 × 6	2 × p4d.24	925 ± 40	580 ± 15	~ 2.6×
4 × 8	2 × p4d.24	1240 ± 70	600 ± 30	~ 3.1×
4 × 12	3 × p4d.24	1610 ± 80	595 ± 25	~ 3.7×
2 × 4	1 × p4d.24	540 ± 55	470 ± 35	~ 2.2×
4 × 4	1 × p4d.24	920 ± 150	590 ± 35	~ 2.6×
6 × 4	2 × p4d.24	660 ± 120	480 ± 75	~ 2.4×
8 × 4	2 × p4d.24	400 ± 65	815 ± 35	~ 1.5×
12 × 4	3 × p4d.24	735 ± 65	355 ± 85	~ 3.1×

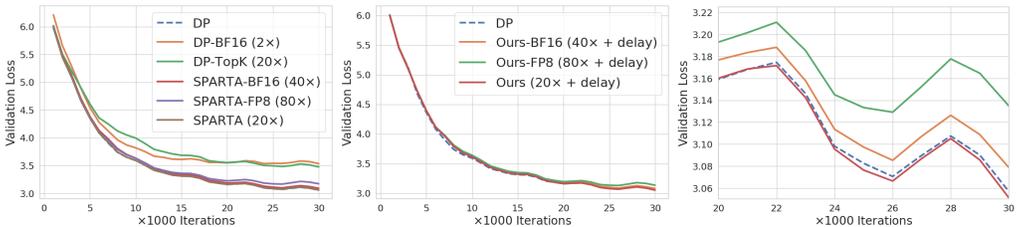
1147 Table 2: We report the SPARTA communication time (i.e., averaging 5% of the parameters) for
1148 tail stage for the 8-layer base model, that would be masked by our asynchronous sparse averaging
1149 for various configurations above and compare it with the forward-backward times. This shows the
1150 empirical speed-up that could be obtained by our method, which ranges from 1.5× – 3.7× and
1151 improves with larger mesh. This speed-up only considers asynchronous DP and any benefits due to
1152 AsyncPP is additional to this.

1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165



1166 Figure 11: Our method in a heterogeneous setup with drastically varying device speeds. The per-
1167 formance degradation is negligible compared to the gain in wall-clock time.

1168
1169
1170
1171
1172
1173
1174
1175
1176
1177



1178 Figure 12: Effect of compression for different methods, for the base model with 4 × 2 mesh on Wiki-
1179 Text. **Left:** Both quantisation and TopK sampling based on weight magnitude (instead of random)
1180 degrade the performance for DP (DP-FP8 did not converge). However, SPARTA is robust to quanti-
1181 zation. **Right:** Similar to SPARTA, our method is robust to quantization even with a 10-step delay,
1182 and the degradation is minimal. TopK sampling did not converge for our method, aligning with our
1183 insight that the sampling needs to be unbiased. DP is even more sensitive to quantization with delay,
1184 and DP-BF16 with delay did not converge. The robustness of sparse averaging (even with delay) to
1185 quantization is intriguing, and it may be explained by the fact that since the quantization error is
1186 introduced only for a small subset (5% in this case) at each iteration, the effect of quantization on
1187 training is negligible. However, this warrants further study, which is beyond the scope of this work.

1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241

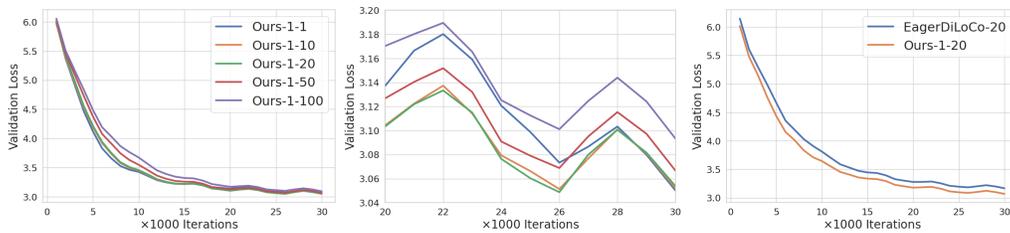


Figure 13: Results in an equivalent setup to EagerDiLoCo (Kale et al., 2025). Ours-1-X denotes our method with subset size 1 (full averaging) and varying averaging interval (i.e., delay for asynchronous DP). On the right, we compare against EagerDiLoCo for 20 inner-steps (i.e., delay). Our EMA approach outperforms eager updates, confirming the strict generality, and the performance varies only slightly with different numbers of inner steps.