
Decentralized Training of Foundation Models in Heterogeneous Environments

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Training foundation models, such as GPT-3 and PaLM, can be extremely expensive,
2 often involving tens of thousands of GPUs running continuously for months. These
3 models are typically trained in specialized clusters featuring fast, homogeneous
4 interconnects and using carefully designed software systems that support both
5 data parallelism and model/pipeline parallelism. Such dedicated clusters can be
6 costly and difficult to obtain. *Can we instead leverage the much greater amount of*
7 *decentralized, heterogeneous, and lower-bandwidth interconnected compute?* Pre-
8 vious works examining the heterogeneous, decentralized setting focus on relatively
9 small models that can be trained in a purely data parallel manner. State-of-the-art
10 schemes for model parallel foundation model training, such as Megatron, only
11 consider the homogeneous data center setting. In this paper, we present the first
12 study of training large foundation models with model parallelism in a decentralized
13 regime over a heterogeneous network. Our key technical contribution is a schedul-
14 ing algorithm that allocates different computational “tasklets” in the training of
15 foundation models to a group of decentralized GPU devices connected by a slow
16 heterogeneous network. We provide a formal cost model and further propose
17 an efficient evolutionary algorithm to find the optimal allocation strategy. We
18 conduct extensive experiments that represent different scenarios for learning over
19 geo-distributed devices simulated using real-world network measurements. In the
20 most extreme case, across 8 different cities spanning 3 continents, our approach is
21 $4.8\times$ faster than prior state-of-the-art training systems (Megatron).

22 1 Introduction

23 Recent years have witnessed the rapid development of deep learning models, particularly foundation
24 models (FMs) [1] such as GPT-3 [2] and PaLM [3]. Along with these rapid advancements, however,
25 comes computational challenges in training these models: the training of these FMs can be very
26 expensive — a single GPT3-175B training run takes 3.6K Petaflops-days [2]— this amounts to \$4M
27 on today’s AWS on demand instances, even assuming 50% device utilization (V100 GPUs peak at
28 125 TeraFLOPS)! Even the smaller scale language models, e.g., GPT3-XL (1.3 billion parameters),
29 on which this paper evaluates, require 64 Tesla V100 GPUs to run for one week, costing \$32K on
30 AWS. As a result, speeding up training and decreasing the cost of FMs have been active research
31 areas. Due to their vast number of model parameters, state-of-the-art systems (e.g., Megatron[4],
32 Deepspeed[5], FairScale[6]) leverage multiple forms of parallelism [4, 7, 8, 9, 10, 11]. However, their
33 design is only tailored to *fast, homogeneous* data center networks.

34 On the other hand, *decentralization* is a natural and promising direction. Jon Peddie Research reports
35 that the PC and AIB GPU market shipped 101 million units in Q4 2021 alone [12]. Furthermore,

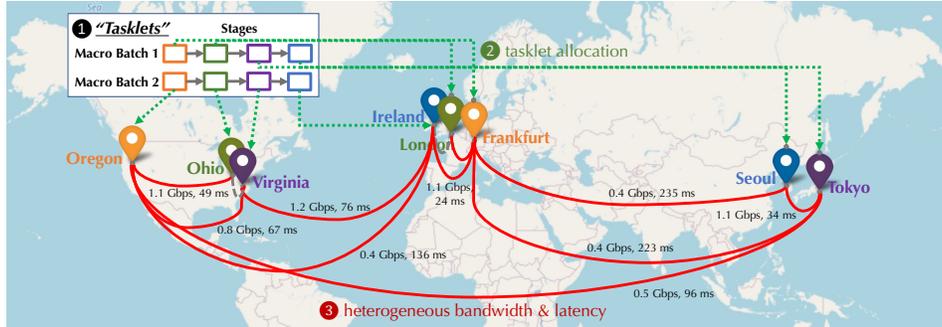


Figure 1: Given ① a set of computation tasklets involved in training foundation models (corresponding to different micro-batches and layers), and ② a heterogeneous network between devices, the goal is to find the optimal ③ allocation of tasklets to devices.

36 many of these GPUs are underutilized. Leveraging this fact, volunteer computing projects such
 37 as Folding@Home [13] have sourced upwards of 40K Nvidia and AMD GPUs continuously [14].
 38 Moreover, The incremental electricity and HVAC costs of running a V100 GPU for a volunteer are
 39 50–100× lower than the spot prices for an equivalent device on AWS [15]. If we could make use
 40 of these devices in a decentralized open-volunteering paradigm for foundation model training, this
 41 would be a revolutionary alternative to the expensive solutions offered by data centers.

42 This vision inspired many recent efforts in decentralized learning, including both those that are
 43 theoretical and algorithmic [16, 17, 18], as well as recent prototypes such as Learning@Home [19]
 44 and DeDLOC [20]. However, efforts to-date in decentralized training either focus solely on *data*
 45 *parallelism* [16, 17, 18, 20], which alone is insufficient for FMs whose parameters exceed the
 46 capacity of a single device, or orient around alternative architectures, e.g., mixture of experts [19].
 47 These alternative architectures provide promising directions for decentralized learning, however, they
 48 are currently only trained and evaluated on smaller datasets and at a smaller computational scale (e.g.,
 49 MNIST and WikiText-2 in [19]) than their state-of-the-art counterparts, e.g., GLaM [21]. In this
 50 paper, we focus on a standard GPT-style architecture, without considering any changes that might
 51 alter the model architecture or the convergence behaviour during training.

52 To fulfill the potential of decentralization for the training of FMs, we need to be able to (1) take
 53 advantage of computational devices connected via heterogeneous networks with limited bandwidth
 54 and significant latency, and (2) support forms of parallelism beyond pure data parallelism. In this
 55 paper, we tackle one fundamental aspect of this goal — how can we assign different computational
 56 “tasklets”, corresponding to a macro-batch and a subset of layers, to a collection of geo-distributed
 57 devices connected via heterogeneous, slow networks? This is not an easy task — even for fast and
 58 homogeneous data center networks, such assignments are still an open ongoing research challenge [22,
 59 23, 24, 25, 26]. For the heterogeneous setting, it becomes even more challenging as the size of the
 60 search space increases dramatically. In the homogeneous setting, the homogeneity of the edges in the
 61 communication graph reduces the search space into many equivalent classes representing allocation
 62 strategies with the same communication costs, enabling efficient polynomial runtime algorithms [23,
 63 24, 22, 25, 26]; however, in the heterogeneous setting, one has to consider potentially exponentially
 64 many more distinct allocation strategies — as we will see later, because of the heterogeneity of the
 65 communication matrix, even the sub-problem of finding the best pipeline parallelism strategy equates
 66 to a hard open loop travelling salesman problem [27].

67 In this paper, we focus on this challenging scheduling problem of decentralized training of FMs over
 68 slow, heterogeneous networks, and make the following contributions:

- 69 • We study the problem of allocating distributed training jobs over a group of decentralized GPU
 70 devices connected via a slow heterogeneous network. More specifically:
 - 71 – To capture the complex communication cost for training FMs, we propose a natural, but
 72 novel, formulation involving decomposing the *cost model* into two levels: the first level is
 73 a *balanced graph partitioning* problem corresponding to the communication cost of data

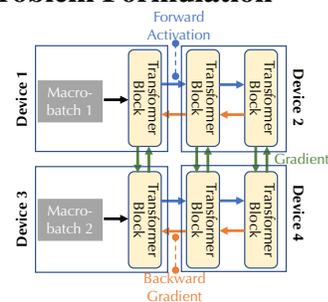
- 74 parallelism, whereas the second level is a joint *graph matching* and *traveling salesman*
 75 problem corresponding to the communication cost of pipeline parallelism.
- 76 – We propose a novel *scheduling algorithm* to search for the optimal allocation strategy given
 77 our cost model. Developing a direct solution to this optimization problem is hard; thus, we
 78 propose an efficient evolutionary algorithm based on a collection of novel heuristics, going
 79 beyond the traditional heuristics used in standard graph partitioning methods [28].
 - 80 • We carefully designed and implemented a collection of *system optimizations* to hide communi-
 81 cation within the computation to further reduce the impact of slow connections.
 - 82 • We conduct extensive experiments that represent different scenarios of collaborative decentral-
 83 ized learning, simulated by using network measurements from different geographical regions of
 84 AWS. In the worldwide setting with 64 GPUs across 8 regions (Oregon, Virginia, Ohio, Tokyo,
 85 Seoul, London, Frankfurt, Ireland), we show that our system is $3.8\text{-}4.8\times$ faster, in end-to-end
 86 runtime, than the state-of-the-art system, Megatron, for training GPT3-XL, *without any differ-*
 87 *ence in what is computed or convergence dynamics*. In addition, we also provide careful ablation
 88 studies to show the individual effectiveness of the scheduler and system optimizations.
 - 89 • We shed light on the potential of decentralized learning — our prototype in the global heteroge-
 90 neous setting is only $1.7\text{-}3.5\times$ slower than Megatron in data centers even though its network can
 91 be $100\times$ slower. We hope this paper can inspire future explorations of decentralized learning
 92 for FMs, over geo-distributed servers, desktops, laptops, or even mobile devices.

93 **Limitations and Moving Forward.** In this paper, we tackle one foundational aspect of decentralized
 94 learning but leave as future work many problems that are important for a practical system. We assume
 95 that communication between devices is relatively stable for a reasonable amount of time and that all
 96 devices are always online without failure or eviction. Note that we also do not train a full system to
 97 full convergence, instead running partial training to confirm intermediate result equivalence across
 98 regimes. Scheduling over a dynamic, heterogeneous environment and providing fault tolerance,
 99 potentially with checkpointing, while training to convergence are directions for future exploration.

100 2 Decentralized Training of Foundation Models: Problem Formulation

101 We first introduce concepts, technical terms, and the procedure of
 102 decentralized training. Then we formally define the scheduling
 103 problem this paper tackles.

104 **Decentralized setting.** We assume a group of devices (GPUs)
 105 participating in collaborative training of a foundation model. Each
 106 pair of devices has a connection with potentially different delay and
 107 bandwidth. These devices can be geo-distributed, as illustrated in
 108 Figure 1, with vastly different pairwise communication bandwidth
 109 and latency. In decentralized training, all layers of a model are



110 split into multiple *stages*, where each device handles a consecutive sequence of layers, e.g., several
 111 transformer blocks [29]. In addition, since the input for foundation model pre-training is huge, e.g., a
 112 few millions of tokens, it is also split into multiple *macro-batches* that can be handled in parallel.

113 **Problem definition.** We define *tasklets* as a collection of computational tasks in foundation model
 114 training — Tasklet $t_{i,j}$ is the forward and backward computation for a stage j with a macro-batch i
 115 of training data in a training iteration. We aim to design an effective scheduler to assign each tasklet
 116 to a particular device so that the training throughput is maximized in decentralized training.

117 **Parallelism.** The above setting involves two forms of parallelism, *pipeline* and *data*. In pipeline
 118 parallelism, the compute in multiple stages is parallelized — each device handles activation or
 119 gradient computation for different macro-batches in parallel and the results can be communicated or
 120 passed to subsequent stages. Data parallelism means that devices compute the gradient for different
 121 macro-batches independently, but need to synchronize these gradients through communication. In
 122 a decentralized environment, the training procedure is *communication-bounded*. The scheduling
 123 problem is to accelerate the communication procedure by allocating tasklets that require high
 124 communication volumes between them to devices with faster connections.

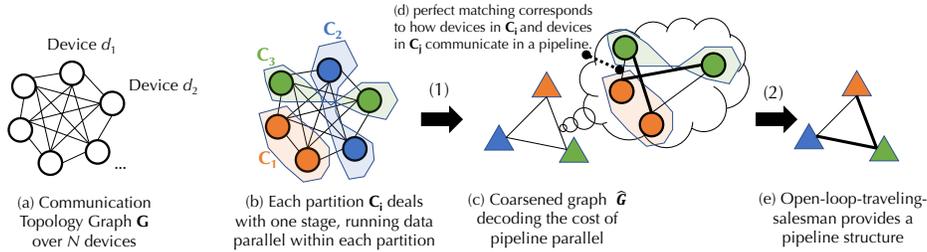


Figure 2: (a) Communication graph \mathbf{G} ; and (b, c, d, e) an illustration of the cost model given \mathbf{G} .

125 **Formalization of the scheduling problem.** Formally, our scheduling problem is as follows.

- 126 • Let $\mathbf{D} = \{d_1 \dots d_N\}$ be a set of N devices; $\mathbf{A} \in \mathbb{R}_+^{N \times N}$ and $\mathbf{B} \in \mathbb{R}_+^{N \times N}$ be the communication
127 matrix between these devices describing the delay and bandwidth respectively, where the delay
128 and bandwidth between device d and d' is $\alpha_{d,d'}$ and $\beta_{d,d'}$.
129 • Given the communication matrix \mathbf{A} and \mathbf{B} , we construct a communication graph \mathbf{G} (Figure 2(a))
130 — each device corresponds to a node in \mathbf{G} and each edge between d and d' is labeled with the
131 average latency and bandwidth between d and d' : $((\alpha_{d,d'} + \alpha_{d',d})/2, (\beta_{d,d'} + \beta_{d',d})/2)$. Even
132 though \mathbf{A} and \mathbf{B} are asymmetric (i.e., upload and download speed might be different), the
133 communication graph \mathbf{G} is symmetric because in our workloads all communications between
134 two devices happen to involve the same amount of upload and download.
135 • The number of stages that a macro-batch needs to go through is D_{pp} (noted as pipeline parallel
136 degree); the number of batch partition that needs to run model gradient synchronization is D_{DP}
137 (noted as data parallel degree); we have $D_{DP} \times D_{pp} = N$, i.e., the total number of devices.
138 • c_{pp} (resp. c_{DP}) represent the number of bytes of activations for a macro-batch (resp. param-
139 eters/gradients for a stage) communicated in pipeline parallelism (resp. data parallelism).
140 • We denote a training tasklet as $t_{i,j}$, where $i \in \{1, \dots, D_{DP}\}$ and $j \in \{1, \dots, D_{pp}\}$, each of which
141 corresponds to one specific macro-batch i and pipeline stage j .
142 • An assignment strategy $\sigma \in \mathbf{D}^{D_{DP} \times D_{pp}}$ assigns, for each tasklet $t_{i,j}$, a device $\sigma_{i,j} \in \mathbf{D}$, which
143 means that device $\sigma_{i,j}$ runs the training tasklet $t_{i,j}$. An valid assignment needs to be unique, i.e.,
144 $\forall (i, j) \neq (i', j'): \sigma_{i,j} \neq \sigma_{i',j'}$. We use Σ to denote the set of all valid assignments.
145 • An *optimal assignment strategy* is an assignment σ that minimizes the communication cost

$$\sigma^* = \arg \min_{\sigma \in \Sigma} \text{COMM-COST}(\sigma)$$

146 **Challenges and Goals.** Our goal is to find the optimal assignment strategy, which involves two
147 challenges: (1) How to effectively model the communication cost $\text{COMM-COST}(\sigma)$ for a given
148 assignment σ under a heterogeneous network environment? and (2) How to effectively search for the
149 optimal assignment strategy that minimizes such a cost? We tackle these two questions in Section 3.

150 3 Scheduling in Heterogeneous Environments

151 Scheduling in the heterogeneous setting is a challenging task, as the size of the search space increases
152 dramatically compared to that of the homogeneous case. In the homogeneous data-center case, the
153 network delay can be usually ignored (e.g., $\mathbf{A} = \mathbf{0}$) and the bandwidth \mathbf{B} are assumed to be formed
154 by just a few constants — e.g., the communication bandwidths between different machines on the
155 same rack are assumed to be same [23, 24, 22, 22, 26]. This significantly constrains the search
156 space — one can ignore the influence of communication given uniform connections [23, 24, 22], or
157 organize the device with a hierarchical structure [22, 26], making the scheduling problem solvable in
158 polynomial time in terms of the number of machines.

159 In contrast, in the fully heterogeneous scenario the communication matrix \mathbf{A} and \mathbf{B} consists of
160 distinct values, which can make the search space grows exponentially. In this section, we describe
161 our scheduler that searches for an optimal strategy in the complex search space.

162 3.1 Overview of the scheduler

163 We carefully design a bi-level scheduling algorithm based on extended *balanced graph partition*
164 problem (see Figure 2), and solve this problem by an evolutionary algorithm with a carefully designed
165 local search strategy. Given an assignment strategy $\sigma = \{\sigma_{i,j}\}$ for all tasklets $\{t_{i,j}\}$, we first

166 model its communication cost. During the training of FMs, the communication costs come from
 167 two different sources: (1) *Data parallel*: All devices that are assigned with the tasklets dealing
 168 with the same stage j (handling the same layers) of different macro-batches need to communicate
 169 within themselves to exchange gradients of these layers. For layer j , we call these devices its
 170 *data parallel group*: $\mathbf{C}_j = \{\sigma_{i,j} \mid \forall i \in [D_{DP}]\}$. We can implement the communication using
 171 different primitives, e.g., AllReduce [30], ScatterGather [31], or other decentralized average
 172 protocols [16]. (2) *Pipeline parallel*: All devices that are assigned with the tasklets dealing with the
 173 same macro-batch i of different stages need to form a pipeline, communicating within themselves to
 174 exchange activations and backward gradients. For macro-batch i , these devices are $\mathbf{P}_i = \{\sigma_{i,j} \mid \forall j \in$
 175 $[D_{PP}]\}$. Because these devices need to form a linear pipeline, any *permutation* over \mathbf{P}_i corresponds
 176 one strategy of how these machines can conduct pipeline parallelism within them.

177 **Scheduling Problem.** The goal of our scheduler is to minimize both costs. One design decision that
 178 we made is to decompose this complex optimization problem into two levels. At the first level, we
 179 consider the best way of forming \mathbf{C}_j 's, incurring data parallel communication costs within them. At
 180 the second level, we consider the cost of pipeline parallelism *given* an layout from the first level:

$$\begin{aligned} \min_{\mathbf{C}_1 \dots \mathbf{C}_{D_{PP}}} \text{COMM-COST}(\mathbf{C}_1 \dots \mathbf{C}_{D_{PP}}) &:= \text{DATAP-COST}(\mathbf{C}_1 \dots \mathbf{C}_{D_{PP}}) + \text{PIPELINEP-COST}(\mathbf{C}_1 \dots \mathbf{C}_{D_{PP}}) \\ \text{s.t. } |\mathbf{C}_1| &= \dots = |\mathbf{C}_{D_{PP}}| = D_{DP}, \forall j, j' : \mathbf{C}_j \cap \mathbf{C}_{j'} = \emptyset, \mathbf{C}_1 \cup \dots \cup \mathbf{C}_{D_{PP}} = \mathbf{D} \end{aligned} \quad (1)$$

181 where computing PIPELINEP-COST($\mathbf{C}_1 \dots \mathbf{C}_{D_{PP}}$) involves finding the optimal pipeline structure.

182 In Section 3.2 and Section 3.3, we provide details on COMM-COST($\mathbf{C}_1 \dots \mathbf{C}_{D_{PP}}$). Notice that this
 183 modified objective makes our problem different from the textbook graph partition problem; thus, we
 184 need a carefully designed evolutionary algorithm for finding such a solution introduced in Section 3.4.

185 3.2 Modelling data parallel communication cost

186 Given the communication graph \mathbf{G} forming data parallel groups $\mathbf{C}_1 \dots \mathbf{C}_{D_{PP}}$ corresponds to a *partition*
 187 of \mathbf{G} — In Figure 2(b), different colors correspond to devices in the same \mathbf{C}_j . The data parallel
 188 cost within \mathbf{C}_j only relies on all communication channels (edges in the communication graph)
 189 connecting devices in \mathbf{C}_j . If we assume a colocated sharded parameter server [31] implementation
 190 for communicating within each \mathbf{C}_j , and recall that c_{DP} represents the total amount of data (in bytes)
 191 that needs to be exchanged during gradient aggregation — each device in \mathbf{C}_j needs to manage
 192 c_{DP}/D_{DP} bytes of parameter shard. Once the gradient is ready, each device needs to send each of its
 193 local shards to the corresponding device; next, each device can aggregate the gradients it receives
 194 from all other devices in \mathbf{C}_j ; and finally, each device will send the aggregated gradient shard to all
 195 other devices. Therefore, we can model the data parallel cost for \mathbf{C}_j as follows:

$$\text{DATAP-COST}(\mathbf{C}_j) = \max_{d \in \mathbf{C}_j} \sum_{d' \in \mathbf{C}_j - \{d\}} 2 \cdot \left(\alpha_{d,d'} + \frac{c_{dp}}{D_{DP} \beta_{d,d'}} \right). \quad (2)$$

196 Here, the total cost is bounded by the *slowest* device ($\max_{d \in \mathbf{C}_j}$), which needs to exchange data with
 197 all other machines ($\sum_{d' \in \mathbf{C}_j - \{d\}}$). Because the communication of these different data parallel groups
 198 $\mathbf{C}_1 \dots \mathbf{C}_{D_{PP}}$ can be conducted in parallel and we are only bounded by the slowest data parallel group.
 199 This allows us to model the total communication cost for data parallelism as:

$$\text{DATAP-COST}(\mathbf{C}_1 \dots \mathbf{C}_{D_{PP}}) = \max_{j \in [D_{PP}]} \text{DATAP-COST}(\mathbf{C}_j)$$

200 3.3 Modeling pipeline parallel communication cost

201 Given $\mathbf{C}_1 \dots \mathbf{C}_{D_{PP}}$, to model the communication cost of pipeline parallelism, we need to consider
 202 two factors: (1) each *permutation* π of $\{\mathbf{C}_1 \dots \mathbf{C}_{D_{PP}}\}$ corresponds to a specific pipeline strategy —
 203 devices in \mathbf{C}_{π_j} and devices in $\mathbf{C}_{\pi_{j+1}}$ communicates to exchange activations (during forward pass) and
 204 gradients on activations (during backward pass); and (2) devices in \mathbf{C}_{π_j} and devices in $\mathbf{C}_{\pi_{j+1}}$ need
 205 to be “matched” — only devices that are dealing with the same macro-batch needs to communicate.
 206 This makes modeling the cost of pipeline parallel communication more complex.

207 To model the cost of pipeline parallel communication, we first consider the best possible way that
 208 devices in \mathbf{C}_j and $\mathbf{C}_{j'}$ can be matched. We do this by creating a *coarsened communication graph*
 209 (Figure 2(c)). A coarsened communication graph $\hat{\mathbf{G}}_{\mathbf{C}_1 \dots \mathbf{C}_{D_{PP}}}$ is a fully connected graph, and each
 210 partition \mathbf{C}_j in the original communication graph \mathbf{G} corresponds to a node in $\hat{\mathbf{G}}_{\mathbf{C}_1 \dots \mathbf{C}_{D_{PP}}}$.

211 In the coarsened graph $\widehat{\mathbf{G}}$, the weight on an edge between \mathbf{C}_j and $\mathbf{C}_{j'}$ corresponds to the following
 212 — if \mathbf{C}_j and $\mathbf{C}_{j'}$ need to communicate in a pipeline, what is the communicate cost of the optimal
 213 matching strategy between devices in \mathbf{C}_j and devices in $\mathbf{C}_{j'}$? Recall that c_{pp} represents the amount
 214 of data between two devices for pipeline parallel communication, we can model this cost by

$$\min_{\mathcal{M}} \max_{(d,d') \in \mathcal{M}} 2 \left(\alpha_{d,d'} + \frac{c_{pp}}{\beta_{d,d'}} \right) \quad (3)$$

215 where \mathcal{M} is a perfect matching between \mathbf{C}_j and $\mathbf{C}_{j'}$ — $(d, d') \in \mathcal{M}$ means that device $d \in \mathbf{C}_j$ will
 216 communicate with device $d' \in \mathbf{C}_{j'}$ (i.e., they deal with the same macro-batch). Computing this value
 217 is similar to the classical minimal sum weight perfect matching problem (MinSumWPM) in bipartite
 218 graphs [32], with the only difference being that we compute the *max* instead of the *sum*. As we will
 219 show in the supplementary material, similar to MinSumWPM, Eq 3 can also be solved in PTIME.

220 The coarsened communication graph captures the pipeline parallel communication cost between two
 221 groups of devices, *assuming* they become neighbors in the pipeline. Given this, we need to find an
 222 optimal permutation of $\mathbf{C}_1 \dots \mathbf{C}_{D_{pp}}$, corresponds to the structure of the pipeline. This becomes the
 223 *open-loop traveling salesman problem* [27] over this condensed graph (Figure 2(e)). Formally, we
 224 have the following definition of the pipeline parallel cost:

$$\text{PIPELINE-COST}(\mathbf{C}_1 \dots \mathbf{C}_{D_{pp}}) = \text{OPENLOOPTSP}(\widehat{\mathbf{G}}_{\mathbf{C}_1 \dots \mathbf{C}_{D_{pp}}}) \quad (4)$$

225 where $\widehat{\mathbf{G}}_{\mathbf{C}_1 \dots \mathbf{C}_{D_{pp}}}$ is the coarsened graph defined above.

226 3.4 Searching via hybrid generic algorithm

227 The scheduling problem solves the optimization problem in Eq 1, which corresponds to a *balanced*
 228 *graph partition* problem with a complex objective corresponding to the communication cost. Balanced
 229 graph partition problem is a challenging NP-hard problem [33]. Over the years, researchers have been
 230 tackling this problem via different ways [34, 35, 36]. We follow the line of research that uses hybrid
 231 genetic algorithm [37, 28] since it provides us the flexibility in dealing with complex objective.

232 **Hybrid Genetic Algorithm.** A hybrid genetic algorithm for balanced graph partition usually follows
 233 a structure as follows. The input is a set of candidate balanced graph partitions which serves as the
 234 initial population. The algorithm generates the next generation as follows. It first generates a new
 235 “offspring” o given two randomly selected “parents” p_1 and p_2 . One popular way is to randomly swap
 236 some nodes between these two parents (we follow [28]). Given this offspring o , we then conduct local
 237 search starting at o to find a new balanced partitioning strategy o^* that leads to better cost. We then
 238 add o^* to the population and remove the worst partition candidate in the population if o^* has a better
 239 cost. As suggested by [37], the combination of heuristic-based local search algorithms and generic
 240 algorithm can accelerate convergence by striking the balance between local and global optimum.

241 **Existing Local Search Strategy.** The key in designing this algorithm is to come up with a good local
 242 search strategy. For traditional graph partitioning task, one popular choice is to use the Kernighan-Lin
 243 Algorithm [38]. Which, at each iteration, tries to find a pair of nodes: d in partition \mathbf{C}_j and d' in
 244 partition $\mathbf{C}_{j'}$, to swap. To find such a pair to swap, it uses the following “gain” function:

$$\text{GAIN}_{KL}((d, \mathbf{C}_j) \leftrightarrow (d', \mathbf{C}_{j'})) = \sum_{d'' \in \mathbf{C}_{j'}} w_{d,d''} - \sum_{d'' \in \mathbf{C}_j - \{d\}} w_{d,d''} + \sum_{d'' \in \mathbf{C}_j} w_{d',d''} - \sum_{d'' \in \mathbf{C}_{j'} - \{d'\}} w_{d',d''} - 2w_{d,d'}$$

245 where $w_{i,j}$ corresponds to the weight between node i and j in the graph. However, directly applying
 246 this local search strategy, as we will also show in the experiment (Section 4) does not work well.
 247 Greedily following GAIN_{KL} does not decrease the communication cost of foundation model training.
 248 Therefore, we have to design a new local search strategy tailored to our cost model.

249 **Improving Local Search Strategy.** Our local search strategy is inspired by two observations:

- 250 1. Removing the device d_1 with a *fast* connection (say with d_2) within partition \mathbf{C}_j will not tend to
 251 change the data parallel cost within \mathbf{C}_j , since it is only bounded by the slowest connections.
- 252 2. Once d_1 is moved to $\mathbf{C}_{j'}$, highly likely the pipeline parallel matching between \mathbf{C}_j and $\mathbf{C}_{j'}$ will
 253 consist of the link $d_1 \leftrightarrow d_2$, since it is a fast connection.

254 Therefore, in our local search strategy we only consider the fastest connection within \mathbf{C}_j : $d_1 \leftrightarrow d_2$
 255 and the fastest connection within $\mathbf{C}_{j'}$: $d'_1 \leftrightarrow d'_2$ and generate four swap candidates: $d_1 \leftrightarrow d'_1$,
 256 $d_1 \leftrightarrow d'_2$, $d_2 \leftrightarrow d'_1$, $d_2 \leftrightarrow d'_2$. We use the following gain function (take $d_1 \leftrightarrow d'_1$ as an example):

$$\text{GAIN}((d, \mathbf{C}_j) \leftrightarrow (d', \mathbf{C}_{j'})) = \frac{1}{|\mathbf{C}_{j'}|} \sum_{d'' \in \mathbf{C}_{j'}} w_{d_1, d''} - w_{d_1, d_2} + \frac{1}{|\mathbf{C}_j|} \sum_{d'' \in \mathbf{C}_j} w_{d'_1, d''} - w_{d'_1, d'_2}$$

257 where $\frac{1}{|\mathbf{C}_{j'}|} \sum_{d'' \in \mathbf{C}_{j'}} w_{d_1, d''}$ measures the *expected* pipeline parallel cost of connecting d_1 with
 258 other devices in $\mathbf{C}_{j'}$ *before the swap*, and w_{d_1, d_2} is the cost of connecting d_1 with other devices in
 259 $\mathbf{C}_{j'}$ *after the swap*, assuming this fast link $d_1 \leftrightarrow d_2$ will now be used for pipeline parallelism.

260 Just like how Kernighan-Lin Algorithm [38] can be extended to a circular version [28] to swap
 261 multiple nodes beyond a pair, we can also extend our method into a circular one, following procedure
 262 as circular KL with our new gain function.

263 3.5 Other System Optimizations

264 We also have some system optimizations to further improve the performance. The most important
 265 optimization involves pipelining of communications and computations. We divide each stage in the
 266 pipeline into three slots: a receiving slot, a computation slot, and a sending slot. The receiving slot of
 267 stage j needs to build connections to receive activations from the stage $j - 1$ in forward propagation
 268 and to receive gradients of activations from stage $j + 1$. The computation slot handles the computation
 269 in forward and backward propagation. Symmetric to the receiving slot, the sending slot of stage j
 270 needs to build connections to send activations to stage $j + 1$ in the forward propagation and send
 271 gradients of activations to stage $j - 1$ in the backward propagations. These three slots are assigned to
 272 three CUDA streams so that they will be further pipelined efficiently; as a result, communication will
 273 overlap with computation. In the decentralized scenario (communication bound), computation can be
 274 fully hidden inside the communication time.

275 4 Evaluation

276 We demonstrate that our system can speed up foundation model training in decentralized setting.
 277 Specifically, (1) We show that our system is $4.8\times$ faster, in end-to-end runtime, than the state-of-the-
 278 art system Megatron training GPT3-XL in world-wide geo-distributed setting. Surprisingly, it is only
 279 $1.7 - 2.3\times$ slower than Megatron in data centers. This indicates that we can bridge the gap between
 280 decentralized and data center training (up to $100\times$ slower networks) through scheduling and system
 281 optimization; (2) We demonstrate the necessity of our scheduler through an ablation study. We show
 282 that with the scheduler, our system is $2.7\times$ faster in world-wide geo-distributed setting.

283 **Experimental Setup** To simulate the decentralized setting, we use 8 different AWS regions (Oregon,
 284 Virginia, Ohio, Tokyo, Seoul, London, Frankfurt, and Ireland) and measure the latency and bandwidth
 285 between these regions (we consider the bandwidth that we can realistically obtain using NCCL and
 286 UDP hole punching between these regions). Given these measurements, we use 64 Tesla V100 GPUs
 287 and control their pairwise communication latency and bandwidth for five different cases:

288 Case 1. Data center on demand. This is a standard setting that a user can obtain to train foundation
 289 models. we use 8 AWS p3.16xlarge nodes (each with 8 V100 GPUs); the intra-node connection
 290 has a bandwidth of 100 Gbps, and the inter-node connection has a bandwidth of 25 Gbps. We do not
 291 manually control latency and bandwidth in this case.

292 Case 2. Data center spot instances. Spot GPUs are cheaper in a data center, but can be located on
 293 different types of machine. In this case, we rent 4 AWS p3.8xlarge nodes (each with 4 V100) and
 294 32 p3.2xlarge nodes (each with 1 V100); the intra- p3.8xlarge node connection has a bandwidth
 295 of 100 Gbps, and the inter-node connection has a bandwidth of 10 Gbps. We do not manually control
 296 latency and bandwidth in this case.

297 Case 3. Multiple Data Centers. We consider two organizations, one in Ohio and another in Virginia,
 298 each organization contributes 32 V100 GPUs; within each organization, the bandwidth is 10 Gbps,
 299 and connections cross different campuses have a delay of 10 ms and bandwidth of 1.12 Gbps.

300 Case 4. Regional geo-distributed. We consider individual GPUs cross four different regions in US
 301 (California, Ohio, Oregon, and Virginia); within each region, the delay is 5 ms and bandwidth is 2
 302 Gbps; cross different regions, the delay is $10\sim 70\text{ms}$ and the bandwidth is $1.0\sim 1.3$ Gbps.

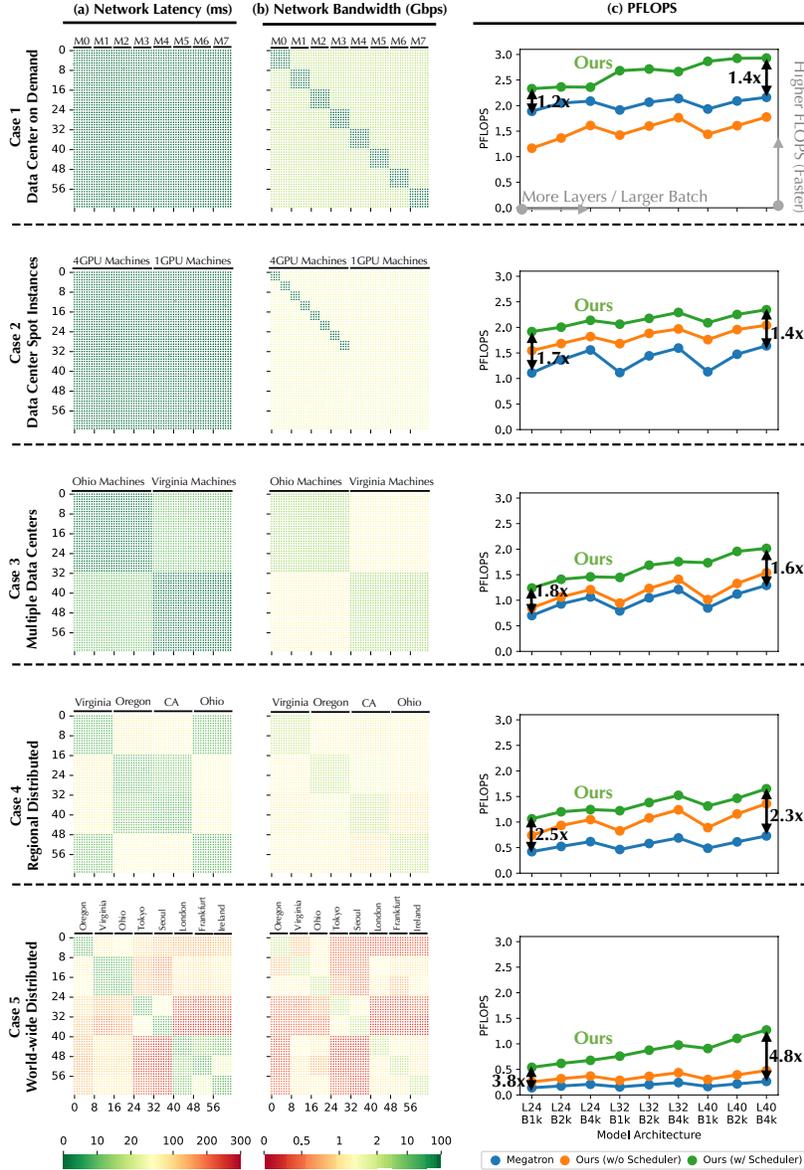


Figure 3: End to end comparison of our system with Megatron in 5 different scenarios. Column (a) and (b) visualize the delay and bandwidth of 5 scenarios respectively; Column (c) illustrates the comparison of Megatron and our system with and without scheduler.

303 *Case 5. World-wide geo-distributed.* We consider individual GPUs across eight different regions
 304 world-wide (Oregon, Virginia, Ohio, Tokyo, Seoul, London, Frankfurt, and Ireland); within each
 305 region, the delay is 5 ms and bandwidth is 2 Gbps; across different regions, the delay is 10~250ms
 306 and the bandwidth is 0.3~1.3Gbps.

307 **Metrics and Model Architecture.** Since we do not introduce any optimizations that might change
 308 the computation or convergence, we can compare all methods by their throughput, we can compare
 309 all systems by the total number of floating point operations per second (PFLOPS), which is inverse
 310 proportional to the *runtime of each iteration* (which we show in Appendix). We use the standard
 311 GPT3-XL architecture [2], while also benchmarked different numbers of layers {24, 32, 40}, and
 312 batch sizes {1024, 2048, 4096}. **Tuning of Megatron.** When comparing with Megatron, we did a
 313 careful grid search of different parallelism settings in Megatron and report the optimal results in each
 314 case—in Case 1, the optimal setting includes tensor model parallelism; all other cases the optimal
 315 settings are based on pipeline and data parallelism. We discuss more details in Appendix.

316 4.1 Results

317 **End-to-end Comparison.** Figure 3(c) shows the end-to-end comparison between Megatron and our
318 system in terms of averaged PFLOPS achieved across different settings and different batch sizes and
319 number of layers. From the world-wide geo-distributed cases, we observe that our system achieves an
320 $4.8\times$ speedup, as compared with Megatron. While in all other cases, our system can be $1.2 - 2.5\times$
321 faster. If we compare our system in Case 5 (world-wide geo-distributed) and Megatron in Case 1 (data
322 center on demand), it is exciting to see that the performance slowdown caused by decentralization is
323 only $1.7 - 3.5\times$! This illustrates the great potential of decentralized training for foundation models.
324 Additionally, Figure 3(c) illustrates another interesting behavior pattern. As increasing the batch size
325 does not increase the communication cost of data parallelism and increasing # layers per device does
326 not increase the communication cost of pipeline parallelism, with a larger batch size and a deeper
327 model, the gap between centralized Megatron and our decentralized system is even smaller.

328 **Effectiveness of Scheduler.** To evaluate the effectiveness
329 of the scheduler, we disable it and use a random
330 assignment in all cases and the results are also illustrated
331 in Figure 3(c). We see that with our scheduler provides up
332 to $2.7\times$ speeds up. To evaluate our local search strategy,
333 we also compare our scheduler with a scheduler that uses
334 the standard Kernighan-Lin algorithm for local search,
335 illustrated in Figure 4 (a - e). We see that, while both
336 outperform random, our carefully designed local search
337 strategy significantly outperforms Kernighan-Lin.

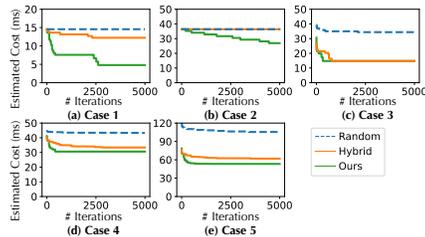


Figure 4: Comparison of Different Local Search Strategies.

338 5 Related Work

339 **Foundation models.** Foundation models[1] refer to models that are trained on large-scale data and
340 can be adapted (e.g., fine-tuned) to a wide range of downstream tasks. Current examples include
341 BERT [39], GPT-3 [2], and CLIP[40]. Foundation models are usually trained in a data center, where
342 the connection between GPUs is fast and homogeneous. ML infrastructures such as Megatron[4]
343 and ZeRO[10, 11] have been proposed to distribute the training of these foundation models in a data
344 center. Megatron uses `AllReduce` to synchronize activations in tensor model parallelism; ZeRO
345 adopts `ScatterGather` to dispatch sharded parameters for layer-wise data parallelism. However,
346 such collective communication paradigms would cause serious performance problems with slow and
347 heterogeneous connections (see Appendix for detailed discussions).

348 **Decentralized optimization.** Decentralized training is first proposed within the scope of data paral-
349 lelism, where each worker only synchronizes gradients with its neighbors (instead of all workers) to
350 remove the latency bottleneck [17, 41, 16, 42, 43, 44]. Recently, [45] has also modified the imple-
351 mentation of data parallelism to support training in an open collaborative environment. Varuna [46]
352 is released by Microsoft to support the training of GPT models in spot instances from a cloud service
353 provider, which has the potential to be extended to the open collective scenario, but there is limited
354 consideration with respect to the challenges of heterogeneous connections.

355 **Volunteer computing.** Distributing computationally intensive tasks over an open collaborative
356 environment has been advocated for a few decades since the development of BOINC[47]; for
357 example, the folding@home project [13] has been running simulations about protein dynamics on
358 volunteers’ personal computers for more than 20 years. Recently, the learning@home project[19]
359 starts to consider training of mixture-of-expert transformers in such a volunteer computing setting.

360 6 Conclusion

361 In this paper, we probe the opportunity to train foundation models via a decentralized training regime
362 with devices connected over a heterogeneous network. We propose an effective scheduling algorithm
363 to assign tasklets from the foundation model pre-train computation. Empirical studies suggest that, in
364 the worldwide geo-distributed scenario, our proposed scheduling algorithm enables a $4.8\times$ speed-
365 up compared to prior state-of-the-art training systems. We believe that the decentralization and
366 democratization of the training of FMs can shift the balance of power positively, but also necessitate
367 new governance structures to help ensure the responsible development and deployment of FMs.

References

- [1] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- [2] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [3] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- [4] Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. Megatron-lm: Training multi-billion parameter language models using model parallelism. *arXiv preprint arXiv:1909.08053*, 2019.
- [5] Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 3505–3506, 2020.
- [6] Mandeep Baines, Shruti Bhosale, Vittorio Caggiano, Naman Goyal, Siddharth Goyal, Myle Ott, Benjamin Lefaudeux, Vitaliy Liptchinsky, Mike Rabbat, Sam Sheiffer, et al. FairScale: A general purpose modular pytorch library for high performance and large scale training, 2021.
- [7] Yuanzhong Xu, HyoukJoong Lee, Dehao Chen, Blake Hechtman, Yanping Huang, Rahul Joshi, Maxim Krikun, Dmitry Lepikhin, Andy Ly, Marcello Maggioni, et al. Gspmd: general and scalable parallelization for ml computation graphs. *arXiv preprint arXiv:2105.04663*, 2021.
- [8] Lianmin Zheng, Zhuohan Li, Hao Zhang, Yonghao Zhuang, Zhifeng Chen, Yanping Huang, Yida Wang, Yuanzhong Xu, Danyang Zhuo, Joseph E Gonzalez, et al. Alpa: Automating inter- and intra-operator parallelism for distributed deep learning. *arXiv preprint arXiv:2201.12023*, 2022.
- [9] Zhuohan Li, Siyuan Zhuang, Shiyuan Guo, Danyang Zhuo, Hao Zhang, Dawn Song, and Ion Stoica. Terapipe: Token-level pipeline parallelism for training large-scale language models. In *International Conference on Machine Learning*, pages 6543–6552. PMLR, 2021.
- [10] Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: Memory optimizations toward training trillion parameter models. In *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–16. IEEE, 2020.
- [11] Jie Ren, Samyam Rajbhandari, Reza Yazdani Aminabadi, Olatunji Ruwase, Shuangyan Yang, Minjia Zhang, Dong Li, and Yuxiong He. {ZeRO-Offload}: Democratizing {Billion-Scale} model training. In *2021 USENIX Annual Technical Conference (USENIX ATC 21)*, pages 551–564, 2021.
- [12] Q4’21 sees a nominal rise in gpu and pc shipments quarter-to-quarter. <https://www.jonpeddie.com/press-releases/q421-sees-a-nominal-rise-in-gpu-and-pc-shipments-quarter-to-quarter>.
- [13] Michael Shirts and Vijay S Pande. Screen savers of the world unite! *Science*, 290(5498):1903–1904, 2000.
- [14] OS Statistics. <https://stats.foldingathome.org/os>, 2022. [Online; accessed 15-May-2022].

- 412 [15] Gpu economics cost analysis. [https://venturebeat.com/2018/02/25/
413 the-real-cost-of-mining-ethereum/](https://venturebeat.com/2018/02/25/the-real-cost-of-mining-ethereum/).
- 414 [16] Xiangru Lian, Ce Zhang, Huan Zhang, Cho-Jui Hsieh, Wei Zhang, and Ji Liu. Can decentralized
415 algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic
416 gradient descent. *Advances in Neural Information Processing Systems*, 30, 2017.
- 417 [17] Anastasia Koloskova, Sebastian Stich, and Martin Jaggi. Decentralized stochastic optimization
418 and gossip algorithms with compressed communication. In *International Conference on
419 Machine Learning*, pages 3478–3487. PMLR, 2019.
- 420 [18] Anastasia Koloskova, Nicolas Loizou, Sadra Boreiri, Martin Jaggi, and Sebastian Stich. A
421 unified theory of decentralized sgd with changing topology and local updates. In *International
422 Conference on Machine Learning*, pages 5381–5393. PMLR, 2020.
- 423 [19] Max Ryabinin and Anton Gusev. Towards crowdsourced training of large neural networks
424 using decentralized mixture-of-experts. *Advances in Neural Information Processing Systems*,
425 33:3659–3672, 2020.
- 426 [20] Michael Diskin, Alexey Bukhtiyarov, Max Ryabinin, Lucile Saulnier, Anton Sinitin, Dmitry
427 Popov, Dmitry V Pyrkun, Maxim Kashirin, Alexander Borzunov, Albert Villanova del Moral,
428 et al. Distributed deep learning in open collaborations. *Advances in Neural Information
429 Processing Systems*, 34:7879–7897, 2021.
- 430 [21] Nan Du, Yanping Huang, Andrew M. Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu,
431 Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, Barret Zoph, Liam Fedus, Maarten
432 Bosma, Zongwei Zhou, Tao Wang, Yu Emma Wang, Kellie Webster, Marie Pellat, Kevin
433 Robinson, Kathy Meier-Hellstern, Toju Duke, Lucas Dixon, Kun Zhang, Quoc V. Le, Yonghui
434 Wu, Zhifeng Chen, and Claire Cui. Glam: Efficient scaling of language models with mixture-
435 of-experts. *CoRR*, abs/2112.06905, 2021.
- 436 [22] Deepak Narayanan, Aaron Harlap, Amar Phanishayee, Vivek Seshadri, Nikhil R Devanur,
437 Gregory R Ganger, Phillip B Gibbons, and Matei Zaharia. Pipedream: generalized pipeline
438 parallelism for dnn training. In *Proceedings of the 27th ACM Symposium on Operating Systems
439 Principles*, pages 1–15, 2019.
- 440 [23] Jakub M Tarnawski, Amar Phanishayee, Nikhil Devanur, Divya Mahajan, and Fanny Nina Par-
441 avecino. Efficient algorithms for device placement of dnn graph operators. *Advances in Neural
442 Information Processing Systems*, 33:15451–15463, 2020.
- 443 [24] Jakub M Tarnawski, Deepak Narayanan, and Amar Phanishayee. Piper: Multidimensional
444 planner for dnn parallelization. *Advances in Neural Information Processing Systems*, 34, 2021.
- 445 [25] Shiqing Fan, Yi Rong, Chen Meng, Zongyan Cao, Siyu Wang, Zhen Zheng, Chuan Wu, Guoping
446 Long, Jun Yang, Lixue Xia, et al. Dapple: A pipelined data parallel approach for training large
447 models. In *Proceedings of the 26th ACM SIGPLAN Symposium on Principles and Practice of
448 Parallel Programming*, pages 431–445, 2021.
- 449 [26] Jay H Park, Gyeongchan Yun, M Yi Chang, Nguyen T Nguyen, Seungmin Lee, Jaesik Choi,
450 Sam H Noh, and Young-ri Choi. {HetPipe}: Enabling large {DNN} training on (whimpy)
451 heterogeneous {GPU} clusters through integration of pipelined model parallelism and data
452 parallelism. In *2020 USENIX Annual Technical Conference (USENIX ATC 20)*, pages 307–321,
453 2020.
- 454 [27] Christos H Papadimitriou. The euclidean travelling salesman problem is np-complete. *Theoreti-
455 cal computer science*, 4(3):237–244, 1977.

- 456 [28] So-Jin Kang and Byung-Ro Moon. A hybrid genetic algorithm for multiway graph partitioning.
457 In *Proceedings of the 2nd Annual Conference on Genetic and Evolutionary Computation*, pages
458 159–166. Citeseer, 2000.
- 459 [29] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
460 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information*
461 *processing systems*, 30, 2017.
- 462 [30] Alexander Sergeev and Mike Del Balso. Horovod: fast and easy distributed deep learning in
463 tensorflow. *arXiv preprint arXiv:1802.05799*, 2018.
- 464 [31] Yimin Jiang, Yibo Zhu, Chang Lan, Bairen Yi, Yong Cui, and Chuanxiong Guo. A unified
465 architecture for accelerating distributed {DNN} training in heterogeneous {GPU/CPU} clusters.
466 In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*,
467 pages 463–479, 2020.
- 468 [32] Michel X Goemans. Lecture notes on bipartite matching. *Massachusetts Institute of Technology*,
469 2009.
- 470 [33] Michael R Garey and David S Johnson. *Computers and intractability*, volume 174. freeman
471 San Francisco, 1979.
- 472 [34] Konstantin Andreev and Harald Racke. Balanced graph partitioning. *Theory of Computing*
473 *Systems*, 39(6):929–939, 2006.
- 474 [35] Peter Sanders and Christian Schulz. Think locally, act globally: Highly balanced graph
475 partitioning. In *International Symposium on Experimental Algorithms*, pages 164–175. Springer,
476 2013.
- 477 [36] Aydın Buluç, Henning Meyerhenke, Ilya Safro, Peter Sanders, and Christian Schulz. Recent
478 advances in graph partitioning. *Algorithm engineering*, pages 117–158, 2016.
- 479 [37] Tarek A El-Mihoub, Adrian A Hopgood, Lars Nolle, and Alan Battersby. Hybrid genetic
480 algorithms: A review. *Eng. Lett.*, 13(2):124–137, 2006.
- 481 [38] Thang Nguyen Bui and Byung Ro Moon. Genetic algorithm and graph partitioning. *IEEE*
482 *Transactions on computers*, 45(7):841–855, 1996.
- 483 [39] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of
484 deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*,
485 2018.
- 486 [40] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
487 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
488 models from natural language supervision. In *International Conference on Machine Learning*,
489 pages 8748–8763. PMLR, 2021.
- 490 [41] Youjie Li, Mingchao Yu, Songze Li, Salman Avestimehr, Nam Sung Kim, and Alexander
491 Schwing. Pipe-sgd: a decentralized pipelined sgd framework for distributed deep net training.
492 In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*,
493 pages 8056–8067, 2018.
- 494 [42] Xiangru Lian, Wei Zhang, Ce Zhang, and Ji Liu. Asynchronous decentralized parallel stochastic
495 gradient descent. In *International Conference on Machine Learning*, pages 3043–3052. PMLR,
496 2018.
- 497 [43] Hanlin Tang, Shaoduo Gan, Ce Zhang, Tong Zhang, and Ji Liu. Communication compression
498 for decentralized training. In *Proceedings of the 32nd International Conference on Neural*
499 *Information Processing Systems*, pages 7663–7673, 2018.

- 500 [44] Hanlin Tang, Xiangru Lian, Ming Yan, Ce Zhang, and Ji Liu. D2: Decentralized training
501 over decentralized data. In *International Conference on Machine Learning*, pages 4848–4856.
502 PMLR, 2018.
- 503 [45] G Tse Edwin, Dana M Klug, and Matthew H Todd. Open science approaches to covid-19.
504 *F1000Research*, 9, 2020.
- 505 [46] Sanjith Athlur, Nitika Saran, Muthian Sivathanu, Ramachandran Ramjee, and Nipun Kwatra.
506 Varuna: scalable, low-cost training of massive deep learning models. In *Proceedings of the*
507 *Seventeenth European Conference on Computer Systems*, pages 472–487, 2022.
- 508 [47] David P Anderson. Boinc: A system for public-resource computing and storage. In *Fifth*
509 *IEEE/ACM international workshop on grid computing*, pages 4–10. IEEE, 2004.

510 Checklist

511 The checklist follows the references. Please read the checklist guidelines carefully for information on
512 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
513 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
514 the appropriate section of your paper or providing a brief inline description. For example:

- 515 • Did you include the license to the code and datasets? **[N/A]**
- 516 • Did you include the license to the code and datasets? **[N/A]**
- 517 • Did you include the license to the code and datasets? **[N/A]**

518 Please do not modify the questions and only use the provided macros for your answers. Note that the
519 Checklist section does not count towards the page limit. In your paper, please delete this instructions
520 block and only keep the Checklist section heading above along with the questions/answers below.

521 1. For all authors...

- 522 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
523 contributions and scope? **[Yes]**
- 524 (b) Did you describe the limitations of your work? **[Yes]** In lines 122-124 of the Problem
525 Formulation, we highlight the core limitations and simplifying assumptions.
- 526 (c) Did you discuss any potential negative societal impacts of your work? **[Yes]** In lines
527 364-366 of the Conclusion, we discuss the potential negative impact of decentralization
528 and democratization, noting that while there are positives, this can also lead to acceler-
529 ation and lack of control. We look forward to actively engaging with the community
530 on questions of governance.
- 531 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
532 them? **[Yes]**

533 2. If you are including theoretical results...

- 534 (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
- 535 (b) Did you include complete proofs of all theoretical results? **[N/A]**

536 3. If you ran experiments...

- 537 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
538 mental results (either in the supplemental material or as a URL)? **[Yes]**
- 539 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
540 were chosen)? **[Yes]** We fully follow the standard settings for GPT3-XL training.
- 541 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
542 ments multiple times)? **[Yes]** We repeated all the experiments for at least three times
543 and reported the average in the paper. Details are reported in Appendix.
- 544 (d) Did you include the total amount of compute and the type of resources used (e.g., type
545 of GPUs, internal cluster, or cloud provider)? **[Yes]**

546 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- 547 (a) If your work uses existing assets, did you cite the creators? **[Yes]**
- 548 (b) Did you mention the license of the assets? **[Yes]** Our systems was built on PyTorch and
549 CuPy and was evaluated using publicly available abstracts. All of them use open-source
550 licenses and can be used for non-commercial educational purposes.
- 551 (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]**
- 552 (d) Did you discuss whether and how consent was obtained from people whose data you’re
553 using/curating? **[N/A]**
- 554 (e) Did you discuss whether the data you are using/curating contains personally identifiable
555 information or offensive content? **[N/A]**

- 556 5. If you used crowdsourcing or conducted research with human subjects...
- 557 (a) Did you include the full text of instructions given to participants and screenshots, if
- 558 applicable? [N/A]
- 559 (b) Did you describe any potential participant risks, with links to Institutional Review
- 560 Board (IRB) approvals, if applicable? [N/A]
- 561 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 562 spent on participant compensation? [N/A]