
Bayesian Optimization with Early Trial Termination for Speeding Up Parallel Neural Network Training

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Training of large neural networks (NNs) is often done in parallel on multiple
2 GPUs. While existing parallel training frameworks easily allow NN training
3 using multi-dimensional parallelism, the challenge remains in finding the optimal
4 hyperparameters, such as the best balance between the sizes of various parallelism
5 dimensions, which would result in the highest training throughput. Due to the
6 large number of possible parallelism configurations (PCs) for a given training
7 scenario, an exhaustive search over them is prohibitively costly. Existing PC
8 optimization methods either require running training trials on a large number of
9 PCs, each of which is costly, or rely on an approximate cost model which may be
10 inaccurate and hardware-specific. To overcome these issues, this paper presents
11 OPPA that can boost the efficiency of Bayesian optimization for optimizing the
12 PC by novelly exploiting (a) the domain knowledge of parallel NN training via
13 parallelism-informed prior beliefs that are general in catering to a variety of NN
14 training scenarios, and (b) early termination of trials involving suboptimal PCs.
15 Despite incorporating these nontrivial efficiency tricks, OPPA is still theoretically
16 guaranteed to achieve sublinear regret. We empirically show that OPPA finds
17 optimal PCs more efficiently than existing methods for parallel training of NNs
18 with varying architectures, training frameworks, and multi-GPU hardware setups.

19 1 Introduction

20 Modern advances in deep learning have arisen from the ability to scale neural networks (NNs) to
21 larger sizes. In natural language processing, for example, transformer models [5, 35], large language
22 models (LLMs) [23, 34], and multimodal LLMs [20, 24], which comprise millions to billions of
23 parameters, have shown tremendous success in tasks such as text classification, text generation, and
24 language understanding. These large NNs often cannot be trained on standard machines with a single
25 processor. To scale up, it is necessary to distribute the NN training workload across a cluster of
26 machines and parallelize the training process. Different parallelism techniques for NN training have
27 been proposed, such as that of data [26, 43], pipeline [10, 22], tensor [31], and their combinations
28 which are also referred to as *multi-dimensional parallelism* [17, 28, 31].

29 To fully utilize the given hardware for reducing the time of NN training, one can maximize its
30 *throughput*, i.e., average number of training steps being run per unit time. The training throughput
31 depends on the selected *parallelism configuration* (PC), which, in large-scale parallel training
32 frameworks [13, 17, 28, 31], is specified by several hyperparameters such as the sizes of various
33 parallelism dimensions. In practice, it is **challenging to accurately determine how the choice of PC**
34 **affects the training throughput** since such a complex relationship depends on the NN architecture,
35 training data, compute and communications hardware, and the exact implementation of the parallel
36 training framework. *Existing works* that approximate (and optimize) the training throughput of a PC

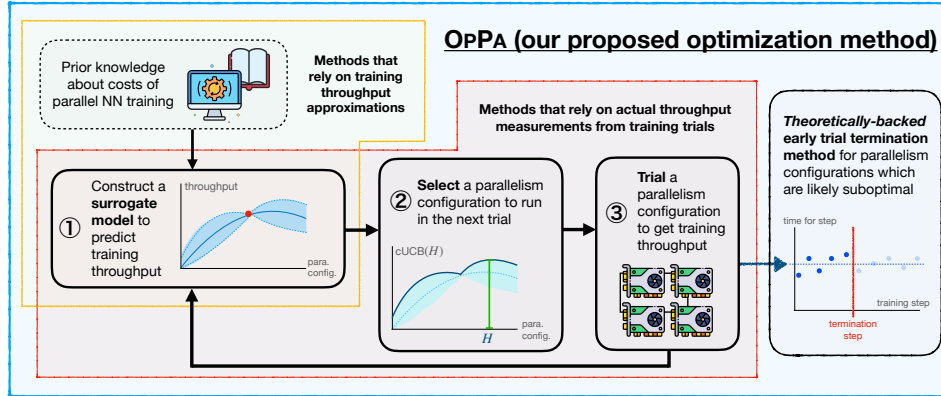


Figure 1: The key idea of OPQA is to exploit domain knowledge of parallel NN training & early trial termination to rapidly find the parallelism configuration (PC) with the highest training throughput.

37 [15, 41, 44] require strong assumptions on the hardware and parallel training implementation, and
 38 may not capture all the complex intricacies of parallel training. So, it may not be reliable to rely on
 39 any of these works alone to recover the actual optimal PC directly.

40 In practice, **the most reliable method to consider all possible factors and their interactions during**
 41 **parallel training is to run real training trials** with each PC on the real hardware. Unfortunately,
 42 due to the large number of possible PCs, an exhaustive search over them is prohibitively costly. To
 43 circumvent this, existing parallel training frameworks have selected a subset of candidate PCs to trial
 44 using simplistic optimization algorithms, but *cannot utilize measured throughputs nor prior domain*
 45 *knowledge*, hence inefficiently needing excessive trials to find the optimal PC (Sec. 2).

46 To efficiently search for the PC with the highest throughput, it is desirable to be able to adaptively
 47 select potentially good PCs for trialing and filter out poor candidates using information from trialed
 48 PCs. Given these considerations, black-box optimization methods such as Bayesian optimization
 49 (BO) (Sec. 2) would be well-suited for the task. However, *a naive use of vanilla BO is still inefficient*
 50 *as it does not exploit the inherent characteristics of PC trialing nor that of the parallel training process.*

51 This paper presents the OPTIMIZER FOR PARALLELISM CONFIGURATIONS (OPQA) that can adaptively
 52 and rapidly optimize the PC for efficient parallel NN training. OPQA (Fig. 1) boosts the
 53 efficiency of BO by novelly exploiting *early termination of trials involving suboptimal PCs* and
 54 *parallelism-informed prior beliefs* that are general in catering to a variety of NN training scenarios.
 55 Despite incorporating these nontrivial efficiency tricks, OPQA is still theoretically guaranteed to
 56 achieve sublinear regret. In Sec. 3, we first formulate the problem of finding the optimal PC as one of
 57 black-box optimization with black-box constraints. Then, we develop surrogate models representing
 58 *parallelism-informed prior beliefs* (based on domain knowledge of parallel NN training) that are
 59 general in catering to a variety of hardware setups and training frameworks (Sec. 4.1) and used by
 60 BO to select promising PCs for trialing (Sec. 4.2). We also describe the process of trialing a PC and
 61 *propose a novel technique to early terminate trials involving suboptimal PCs* with both theoretical
 62 and empirical justifications (Sec. 4.3). In Sec. 5, we empirically show that OPQA can be easily used
 63 on top of existing parallel training frameworks to find a good PC for training NNs more efficiently
 64 than existing methods and vanilla BO without modifications.

65 2 Background and Related Works

66 **Parallel NN training on GPUs.** A large NN can be effectively trained by distributing its workload
 67 across multiple GPUs. Different parallelism dimensions split the workload differently, which affects
 68 the amounts of computation per GPU, memory required in each GPU, and communication between
 69 GPUs. The most basic parallelism technique is data parallelism (DP) [16] where a batch of training
 70 data is split and distributed to each device to be separately processed by its local NN replica before
 71 gathering the gradients from all devices to update the NN weights. While DP is simple, it requires
 72 replicating the NN onto each device, thus incurring additional memory. To resolve this, techniques,

73 such as the Zero Redundancy Optimizer (ZERO) [26] in the DEEPSPEED package or Fully-Shared
74 Data Parallel (FSDP) [43], have been proposed to perform some sharding of NN weights or gradients
75 to avoid full NN replication. Furthermore, tensor parallelism (TP) [2, 31] and pipeline parallelism
76 (PP) [10, 22] have been proposed to shard the tensors in the NN and the NN execution pipelines onto
77 multiple devices, respectively. The specific implementations of DP, TP, and PP, such as which tensors
78 (e.g., NN weights, gradients) are sharded or how many shards they are partitioned into, are often
79 manually specified by the user to control the overall throughput. App. A discusses DP, TP, and PP
80 more. Besides, other types of parallelism (e.g., sequence [18], expert [27]) have also been developed.

81 **Multi-dimensional parallelism.** Several frameworks [17, 28, 31] have also been developed to
82 combine different types of parallelism into the same training process. These frameworks provide
83 simple interfaces for users to specify the desired *parallelism configuration* (PC) that comprises
84 hyperparameters controlling the execution of the parallel training process, a subset of which specifies
85 the sizes of various parallelism dimensions. These frameworks then automatically handle tensor
86 sharding and execute the parallel training pipeline as per the specified PC. These frameworks may
87 also manage training on multi-node or even heterogeneous hardware. While these frameworks allow
88 practitioners to easily execute parallel NN training, *finding the PC with the highest throughput is*
89 *challenging* since it depends non-trivially on the NN architecture, training data, GPU specifications,
90 communication bandwidth between GPUs, and the exact implementation of the parallel training
91 framework [17, 19, 36]. For example, DP is ineffective for large NNs since excessive NN replications
92 can cause out-of-memory errors. Also, PP is less effective on smaller NNs as the communication
93 costs between pipeline stages may dominate the computation of the fragmented pipeline.

94 **Optimizing PC.** The most accurate way to find the optimal PC would be to trial all possible PCs on
95 the actual training hardware and determine which one yields the highest training throughput. However,
96 this is prohibitively costly since there can be a large number of possible PCs and trialing each PC can
97 be computationally costly. To circumvent this, frameworks such as NEMO [13] and DEEPSPEED
98 [28] have implemented methods for automatic PC tuning¹ based on running NN training trials for
99 a few training steps on a number of PCs. The PCs trialed are often either selected non-adaptively
100 (e.g., based on random selection) or adaptively based on a simple surrogate function. However,
101 these methods *cannot efficiently use the measured throughputs of trialed PCs* to perform informed
102 optimization, and therefore still require a large number of training trials to obtain a good PC.

103 Since training trials are costly, we can consider constructing a surrogate model to approximate the
104 computation and communication costs for different PCs [15, 41, 44]. This would allow us to use
105 domain knowledge to filter out suboptimal PCs and run fewer or even no trials at all. These methods,
106 however, have *implicitly assumed that the surrogate model correctly predicts the actual training*
107 *throughputs*. This is not possible in practice since the surrogate model cannot capture all the complex
108 intricacies of parallel training. Furthermore, a fixed surrogate model cannot be easily extended
109 to include new hyperparameters into the PC, which is important especially with the ever-growing
110 parallel training literature.

111 **Overview of BO.** To efficiently search for the PC with the highest throughput, we utilize BO
112 [7, 8]. BO aims to maximize an expensive-to-query black-box function \mathcal{R} (representing the training
113 throughput over the space of all possible PCs) whose derivative is unknown. The black-box function
114 \mathcal{R} is modeled by a Gaussian process (GP) prior belief specified by its prior mean and covariance
115 [29]. Given the measured throughputs of trialed PCs, we can obtain a predictive distribution of the
116 throughputs (of untried PCs) in the form of a GP posterior belief specified by its posterior mean
117 and covariance. Then, the BO algorithm selects a PC for trialing that maximizes a given acquisition
118 function (e.g., expected improvement [11], upper confidence bound [33]) based on the GP posterior
119 belief. Such an acquisition function trades off between exploring unique PCs that have not been
120 queried to predict \mathcal{R} better vs. exploiting PCs likely to have high values of \mathcal{R} so as to efficiently find
121 the optimal PC. App. B gives a more technical overview of GP modeling and BO.

122 BO is used in a wide range of black-box optimization problems, such as experimental design [14, 25]
123 and material design [42]. More relevant to our work here, BO is also commonly used to optimize
124 the NN architecture to achieve the best performance in a given task [32]. Finding the optimal PC,
125 however, differs in two aspects: (a) The hyperparameters in a PC affecting the training throughput
126 have better-defined mechanics (even if not completely known), which can be partially described

¹<https://docs.nvidia.com/nemo-framework/user-guide/latest/usingautoconfigurator.html>,
<https://www.deepspeed.ai/tutorials/autotuning/>

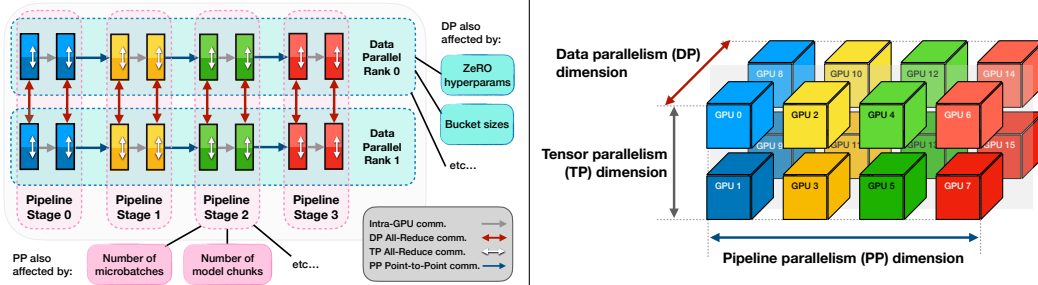


Figure 2: (Left) Visualization of a parallelism configuration (PC) with tunable hyperparameters. (Right) Visualization of GPU allocation for 3D parallelism according to the dimension sizes.

127 based on domain knowledge. Modeling via a GP allows us to incorporate this knowledge through
 128 an informed choice of the GP prior belief, which can reduce the number of trials required. (b)
 129 Trialing a PC allows for repeated measurements of the time taken per training step, which are used to
 130 calculate the throughput. Running many training steps is costly and may be unnecessary, so an early
 131 termination of trials can be adaptively applied to reduce the number of training steps needed to be
 132 run. We will elaborate more in the sections to follow.

133 3 Problem Setting

134 A parallelism configuration (PC) comprises several tunable hyperparameters found in typical parallel
 135 training frameworks and a subset of them determines the *sizes of various parallelism dimensions*, as
 136 visualized in Fig. 2. In our work here, we will mainly consider 3D parallelism where dp , tp , and
 137 pp denote the sizes of the data, tensor, and pipeline parallelism dimensions, respectively; note that
 138 their product $dp \cdot tp \cdot pp$ equals the number n_gpus of available GPUs. Also, as demonstrated in our
 139 experiments (Sec. 5), other types of parallelism can be included in the PC as well. The remaining
 140 hyperparameters determine the *specific implementations of the parallel training framework*, including
 141 (but not limited to) hyperparameters in the ZER0 optimizer, the number of model chunks for PP,
 142 whether communication overlaps are used, and whether gradient checkpointing is used. App. C.1
 143 lists all of the hyperparameters considered in the PC for our problem, including additional parallelism
 144 dimensions that we consider in the experiments.

145 Let \mathcal{H} be the set of all possible PCs. The goal of our problem
 146 is to find the optimal PC $H^* \in \mathcal{H}$ with the *highest training*
 147 *throughput* (i.e., largest average number of training steps being
 148 run per unit time) and its *maximum memory usage not exceeding*
 149 M_0 (i.e., fitting on the given GPUs). For any PC $H \in \mathcal{H}$, let
 150 $\mathcal{R}(H)$ and $\mathcal{M}(H)$ be the throughput and the maximum memory
 151 usage, respectively. Then, our problem of finding the optimal PC
 152 $H^* \in \mathcal{H}$ can be formulated as one of constrained maximization:

$$\max_{H \in \mathcal{H}} \mathcal{R}(H) \quad \text{subject to} \quad \mathcal{M}(H) \leq M_0.$$

153 When querying a PC H , we run a short training trial (instead
 154 of training the NN till convergence) to obtain sufficiently accurate
 155 estimates of the throughput and maximum memory usage.
 156 To estimate the throughput of H , we measure the time taken
 157 $t_1, t_2, \dots, t_{q_{\max}}$ for q_{\max} consecutive training steps and use them
 158 to estimate the throughput: $\mathcal{R}(H) \approx \bar{r}_{q_{\max}} \triangleq q_{\max}^{-1} \sum_{j=1}^{q_{\max}} t_j^{-1}$. As
 159 shown in Fig. 3, the time taken for each training step fluctuates
 160 due to factors like additional overheads during process initial-
 161 ization, fluctuating communication speeds, and unpredictable system behaviors occurring during
 162 training. So, it is more reliable to estimate the throughput by averaging the time taken over multiple
 163 (rather than a few) training steps. While we can run the full trial with q_{\max} training steps, we can also
 164 terminate the trial early after fewer than q_{\max} steps at the cost of a higher variance of the estimator \bar{r}_q .

165 When designing our method, we note two other characteristics of the PC optimization problem:

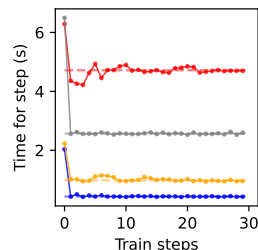


Figure 3: Time taken for each training step for a few training trials with various PCs. Each color represents a different PC used, and the dashed line represents the throughput estimate for that PC. Note the fluctuations in the time taken for each training step.

166 [A] Since \mathcal{R} and \mathcal{M} depend on many factors that may be difficult to model or be known exactly, in
 167 practice, **their exact forms are only partially known** given the domain knowledge. Therefore,
 168 desirable surrogate models for \mathcal{R} and \mathcal{M} **should be able to consider how the actual measure-**
 169 **ments may deviate from the domain knowledge.** Additionally, the surrogate models should
 170 *also be able to quantify its confidence* of any prediction it makes.

171 [B] Though a PC can be (and should be) trialed on real hardware, *running a full trial incurs a high*
 172 *cost.* This is especially true with suboptimal PCs since the same number of training steps on a
 173 suboptimal PC would incur more time. This motivates us to design an optimization algorithm
 174 such that *promising PCs are prioritized for trialing* while **PCs which are likely to be suboptimal**
 175 **are not trialed or are only trialed for a shorter duration.**

176 4 Method

177 In this section, we describe OPPA, which ex-
 178 ploits parallelism-informed GP prior beliefs and
 179 early trial termination to boost the performance
 180 of BO in optimizing the PC for parallel train-
 181 ing. As summarized in Algorithm 1, OPPA alter-
 182 nates between 3 steps: ① modeling the training
 183 throughput and maximum memory usage based
 184 on measurements via *parallelism-informed GP*
 185 *prior beliefs*, ② selecting the next PC to trial
 186 via BO, and ③ running an NN training trial
 187 (and *early terminating it if suboptimal*) to obtain
 188 measurements for the throughput and maximum
 189 memory usage for the selected PC.

Algorithm 1 OPPA (shortened version of Algorithm 4)

```

1: Generate  $\mathcal{H}$  of all valid PCs
2: for  $i = 1, 2, \dots$  do
3:   // Step ① – Modeling  $\mathcal{R}$  and  $\mathcal{M}$  (Sec. 4.1)
4:   Obtain GP predictive/posterior beliefs of  $\mathcal{R}$  and  $\mathcal{M}$ 
   (22) from parallelism-informed GP prior beliefs using
   measurements of trialed PCs
5:   // Step ② – Selecting the next PC to trial (Sec. 4.2)
6:   Select  $H_i \in \mathcal{H}$  with cUCB criterion (27)
7:   // Step ③ – Trialing the selected PC (Sec. 4.3)
8:   for training steps  $q = 1, \dots, q_{\max}$  do
9:     Measure training step time  $t_{i,q}$ 
10:    if  $I_{i,q} = 0$  then terminate trial // From (3)
11:    if budget exhausted then break
12: return  $H_i$  with the best measured throughput

```

190 4.1 Constructing Surrogate Models

191 In Step ①, we construct surrogate models to predict the throughput and the maximum memory usage.
 192 As suggested in [A], to explicitly capture the deviations from the domain knowledge, we decompose
 193 the actual throughput \mathcal{R} and maximum memory usage \mathcal{M} into

$$194 \mathcal{R}(H) = \widehat{\mathcal{R}}(H; \theta_{\mathcal{R}}) + f_{\mathcal{R}}(H) \quad \text{and} \quad \mathcal{M}(H) = \widehat{\mathcal{M}}(H; \theta_{\mathcal{M}}) + f_{\mathcal{M}}(H) \quad (1)$$

195 where the functions $\widehat{\mathcal{R}}$ and $\widehat{\mathcal{M}}$ (with hyperparameters $\theta_{\mathcal{R}}$ and $\theta_{\mathcal{M}}$) capture the domain knowledge of
 196 parallel NN training, and $f_{\mathcal{R}}$ and $f_{\mathcal{M}}$ capture the unknown deviations from the domain knowledge.

196 **Domain knowledge.** The functions $\widehat{\mathcal{R}}$ and $\widehat{\mathcal{M}}$ aim to, respectively, *predict* \mathcal{R} and \mathcal{M} using only the
 197 domain knowledge of parallel NN training. Hence, *they are not expected to give completely accurate*
 198 *predictions but instead capture general trends in any given NN training scenario;* any factors not
 199 accounted for would be captured in the unknown deviation terms anyway, to be described below.

200 In practice, OPPA allows the predictions of $\widehat{\mathcal{R}}$ and $\widehat{\mathcal{M}}$ to be improved by incorporating additional
 201 domain knowledge from practitioners with more expertise in parallel NN training. Consequently, this
 202 would enable OPPA to more rapidly find the optimal PC.

203 For OPPA, we provide default choices of $\widehat{\mathcal{R}}$ and $\widehat{\mathcal{M}}$ by considering parallel NN training under
 204 idealized conditions: For $\widehat{\mathcal{R}}$, we consider the time taken per training step for the computation
 205 $\widehat{\mathcal{T}}_{\text{comp}}$ and for the communication $\widehat{\mathcal{T}}_{\text{comm}}$, which can be combined to approximate the throughput as
 206 $\widehat{\mathcal{R}}(H; \theta_{\mathcal{R}}) = [\widehat{\mathcal{T}}_{\text{comp}}(H; t_c) + \widehat{\mathcal{T}}_{\text{comm}}(H; \mathbf{C})]^{-1}$ where $\theta_{\mathcal{R}} = \{t_{\text{comp}}, \mathbf{C}\}$ are learned hyperparameters.
 207 For $\widehat{\mathcal{T}}_{\text{comp}}$, we consider the additional computation time arising from the pipeline bubble in PP [22].
 208 For $\widehat{\mathcal{T}}_{\text{comm}}$, inspired by [38], we consider an idealized training scenario visualized in Fig. 2, which
 209 assumes the NN to contain roughly identical blocks, and considers the communication from the
 210 All-Reduce operations involved in DP and TP, and point-to-point communications involved in PP.
 211 We separately model the intra- and inter-node connections, collapsing the constants for specific
 212 communication types into \mathbf{C} ; further details are in App. D.1. For $\widehat{\mathcal{M}}$, we consider the memory
 213 required per GPU to store the NN parameters and gradient values for backpropagation, as detailed in
 214 App. D.2.

215 The reason for our choices of $\widehat{\mathcal{R}}$ and $\widehat{\mathcal{M}}$ is that they are general in catering to NN training scenarios
 216 with varying model, hardware, and network configurations, while exhibiting a simple analytical form
 217 whose forward pass is cheap. We also *use learnable parameters* (as opposed to fixed constants) to
 218 capture approximate cost multipliers which depend on the model and training data sizes, and intra-
 219 and inter-node network communications, each of which vary across different training scenarios. The
 220 design choice allows the prior beliefs to capture the general trends of the throughput and maximum
 221 memory usage well enough especially when combined with the learned unknown deviation, while
 222 not incurring too much additional computational costs during the PC optimization process.

223 **Unknown deviations from domain knowledge.** Due to incomplete domain knowledge to fully
 224 describe a parallel training process, we learn the deviations $f_{\mathcal{R}}$ and $f_{\mathcal{M}}$ of the actual measurements
 225 from the domain knowledge using real training trials by modeling them as Gaussian processes (GPs).
 226 The benefit of using a GP is twofold. First, a GP is typically flexible enough to model unknown
 227 functions that may not have an analytical form. Second, a GP can quantify its uncertainty, which
 228 allows the surrogate model to determine how much it knows about the throughput of a certain PC.
 229 This allows OPPA to potentially trial PCs whose throughput it is more uncertain about.

230 To model $f_{\mathcal{R}}$ and $f_{\mathcal{M}}$, we use a GP with zero mean such that the resulting GPs will be “centered”
 231 around the prior means. For their prior covariances, we concatenate each hyperparameter value
 232 in H , then scale each dimension to be between 0 and 1 to form an embedding $e(H)$. Given the
 233 embedding, we then use the Matern kernel [29] with $\nu = 5/2$ where the distance between two PCs
 234 is the Euclidean distance of their corresponding embeddings with kernel hyperparameters θ_k . In
 235 App. D.3, we provide the equation for the Matern kernel, and describe the informativeness of the
 236 distance between embeddings in distinguishing between PCs as reflected by the empirical results.

237 The additive decomposition of the domain knowledge and the unknown deviation in (1) implies that
 238 \mathcal{R} and \mathcal{M} are samples from **parallelism-informed prior beliefs**, given by the GPs

$$\mathcal{R} \sim \mathcal{GP}(\widehat{\mathcal{R}}(\cdot; \theta_{\mathcal{R}}), k(\cdot, \cdot; \theta_k)) \quad \text{and} \quad \mathcal{M} \sim \mathcal{GP}(\widehat{\mathcal{M}}(\cdot; \theta_{\mathcal{M}}), k(\cdot, \cdot; \theta_k)). \quad (2)$$

239 Using measurements of trialed PCs, we find the optimal hyperparameters $\{\theta_{\mathcal{R}}, \theta_{\mathcal{M}}, \theta_k\}$ by maximiz-
 240 ing the marginal log-likelihood [29] and the subsequent GP posterior beliefs for the throughput and
 241 the memory usage for any PC $H \in \mathcal{H}$, which would be normal distributions with their respective
 242 mean and variance functions, as stated in App. D.4.

243 4.2 Selecting the Next PC to Trial

244 In Step ②, the next PC to trial is chosen based on the surrogate models constructed in ① using BO.
 245 The next PC $H_i \in \mathcal{H}$ to trial in round i is chosen to be the PC which maximizes the constrained
 246 upper confidence bound (cUCB) [33, 37], which we detail further in App. E.1. The cUCB criterion
 247 considers a balance between *exploration* of PCs which have not been trialed, and *exploitation* of PCs
 248 which are similar to those with already high throughputs [8, 11]. With this balance, OPPA is able to
 249 trial enough PCs to construct reasonable surrogate models for \mathcal{R} and \mathcal{M} , while still being able to
 250 trial enough promising PC candidates to find one with the maximum training throughput. This allows
 251 the optimization to be more guided and more efficient, satisfying the requirement in B.

252 4.3 Trialing the Next PC, while Terminating Suboptimal Trials Earlier

253 Finally, in Step ③, we run a training trial on the PC H_i chosen in Step ② for q_{\max} training steps,
 254 and measure the time taken for each training step as $t_{i,1}, \dots, t_{i,q_{\max}}$. We then use these measurements
 255 to obtain an estimate of the throughput $\bar{r}_{i,q_{\max}}$ where $\bar{r}_{i,q} = q^{-1} \sum_{j=1}^q t_{i,j}^{-1}$ and whose variance $\sigma_{\bar{r}_{i,q}}^2$
 256 we detail in App. F.1. In practice, as demonstrated earlier in Fig. 3, the actual measured time for a
 257 sequence of training steps may contain outliers which can skew the throughput estimates. As we
 258 discuss in App. F.2, to make our estimate more accurate, we remove outlier values of $t_{i,j}$ before
 259 computing the throughput estimate. Meanwhile, the maximum memory usage m_i across all training
 260 steps is retrieved by calling `torch.cuda.max_memory_allocated()`.

261 **Early trial termination.** In practice, some PCs do not need to be trialed for the full q_{\max} steps since
 262 fewer training steps are sufficient to determine that the PC is suboptimal. This is because $\bar{r}_{i,q}$ is an
 263 estimator for the throughput $\mathcal{R}(H_i)$ whose variance $\sigma_{\bar{r}_{i,q}}^2$ decreases as q gets larger. Intuitively, this
 264 means that after a certain number of steps, $\bar{r}_{i,q}$ will estimate $\mathcal{R}(H_i)$ with low enough variance that
 265 we can determine that it would not improve upon the best PC found so far.

266 To save time on these trials, we consider early trial termination where the training trial only continues
 267 if the estimated throughput is above some threshold. Formally, define an indicator variable

$$I_{i,q} = \mathbb{1}[(q \leq q_{\min}) \vee (\bar{r}_{i,q} \geq \max_{l \in \{1, \dots, i-1\}} \bar{r}_{l, \hat{q}_l} + \tau_q)]. \quad (3)$$

268 In words, $I_{i,q} = 1$ when fewer than q_{\min} training steps have been ran, or when the estimated
 269 throughput of H_i is likely to be higher than the throughputs found so far. We can continue the q th
 270 training step as long as $I_{i,q-1} = 1$, and terminate the training trial at the step \hat{q}_i when $I_{i, \hat{q}_i} = 0$, which
 271 is when H_i is unlikely to improve the best PC found so far. As we intuitively describe in App. F.3.1,
 272 this also corresponds to when more training steps are unlikely to provide additional information about
 273 the optimal PC or optimal throughput. This saves computation time in practice while still allowing
 274 the BO algorithm to recover the optimal PC, as formalized below:

275 **Theorem 4.1** (Informally stated in terms of \mathcal{R}). *Consider a simplified case where $\mathcal{M}(H) \leq M_0$ for
 276 all $H \in \mathcal{H}$. Let N be the number of PCs trialed so far. Then, there exists some choice of $\{\tau_q\}_{q=1}^{q_{\max}}$
 277 such that, with high probability, the cumulative regret is $\sum_{i=1}^N (\mathcal{R}(H^*) - \mathcal{R}(H_i)) = \tilde{O}(\sqrt{N/q_{\min}})$,
 278 and for all $i = 2, \dots, N$, if $\mathcal{R}(H_i) < \max_{k \in \{1, \dots, i-1\}} \mathcal{R}(H_k)$, then $\hat{q}_i < q_{\max}$.*

279 In App. F.3.2, we prove a more general version of Thm. 4.1 (where \mathcal{R} can instead be any function
 280 with a bounded RKHS norm), and provide empirical verification for early termination. Note that
 281 while Thm. 4.1 assumes $\mathcal{M}(H) \leq M_0$ for all $H \in \mathcal{H}$, the throughput cannot be measured when
 282 $\mathcal{M}(H) > M_0$ anyway (since the run would encounter an out-of-memory error), and so can be
 283 omitted in the theoretical analysis for convenience. Thm. 4.1 shows that sublinear regret can be
 284 achieved even with early termination, which means that OPPA will be able to recover the PC with
 285 the best throughput while allowing efficiency gains in practice. Furthermore, it also shows that PCs
 286 whose throughput is smaller than those of PCs already trialed will likely have their trials terminated
 287 early, allowing OPPA to save resources from trialing suboptimal PCs as mentioned in B. We present
 288 the complete pseudocode for OPPA incorporating steps ①, ②, and ③ in App. G.

289 5 Experiments and Discussion

290 In this section, we present the results for OPPA when used to find the optimal PC for training NNs on
 291 multi-GPU systems. We consider optimizing the PC for training NNs with different architectures and
 292 sizes, on various hardware configurations with varying number of GPUs, and on different parallel
 293 training frameworks. Detailed setups for the training scenarios are found in App. H.1.

294 We compare OPPA against RANDOM (random selection), XGBOOST (adaptive selection based
 295 on XGBOOST surrogate model [3]) as used by DEEPSPEED [28]), COST-MODEL (method which
 296 solely relies on the cost model of the throughput), and VANILLA-BO (BO without any additional
 297 modifications). App. H.2 provides detailed descriptions of these methods.

298 In each case, we plot the best obtained throughput (in training steps per second) vs. how long the
 299 optimization has been run rather than vs. the number of PCs that have been trialed. We do so since
 300 how long the optimization has been run is a better reflection of the costs involved in optimizing the
 301 PC, as it captures both the time required to run the optimization algorithm and the time to run the
 302 training trials. The latter also includes the time to initialize the NN training, to run the training steps,
 303 and to switch the training to a different PC, and will also be higher if the algorithm trials more PCs or
 304 trials suboptimal PCs. Nonetheless, plots for the achieved throughput vs. the number of trials run are
 305 in App. I.1. The model loss are independent of the chosen PC and thus are not reported.

306 **Training on single-node setups.** We first consider NN training on a single node. Fig. 4a shows the
 307 results for finding the optimal PC for training the BERT model [5]. We see that OPPA, which exploits
 308 a parallelism-informed prior and early termination mechanism, is able to outperform the non-adaptive
 309 and adaptive methods, and even outperforms vanilla BO. We find that OPPA automatically prioritizes
 310 PCs with only DP which matches our intuition that DP should be adequate for smaller NNs. Fig. 4b
 311 considers the Qwen model [39] where OPPA again finds a better PC compared to the other methods.
 312 Due to a larger model, OPPA now prefers a mix of DP and PP to reduce memory use while incurring
 313 minimal additional inter-GPU communications. We present the PCs selected by OPPA in App. I.2.

314 To visualize the efficiency gains of OPPA, in Fig. 5a, we see that OPPA is able to trial many more PCs
 315 in a given time period when compared to other methods due to the early termination of suboptimal
 316 trials to avoid wasting time. While OPPA will perform more switching between PCs to trial, the cost

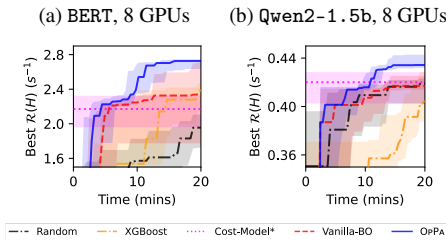


Figure 4: The best obtained throughput (higher is better) vs. the duration each method has been run for, when training on a system with a single node on COLOSSAL-AI [17]. For each method, the line shows the median of repeated trials, while the error band shows the quartiles.

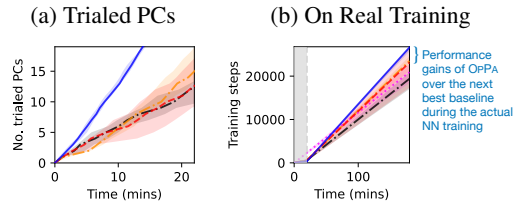


Figure 5: Efficiency of OPQA. The plots show statistics for the experiments in Fig. 4a. Fig. 5a shows the no. of PCs trialed by each method during its run. Fig. 5b shows the no. of training steps that can be run by each method during optimization (the shaded gray region on the left) and during subsequent NN training using the optimized PC (unshaded region on the right). The legend is the same as in Fig. 4.

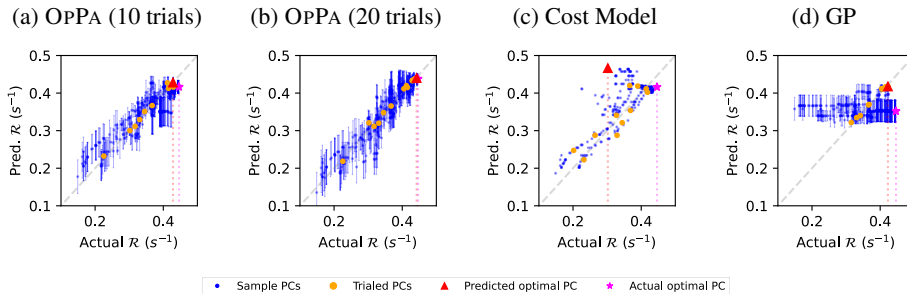


Figure 6: Predicted throughputs from different surrogate models vs. the measured throughputs for the Qwen2-1.5b example (Fig. 4b). Figs. 6a and 6b show the predictions from OPQA after 10 and 20 trials respectively, with error bars for the predictive variance. Fig. 6c shows the predictions from the cost model alone, while Fig. 6d shows the predictions from using a GP alone (both after 20 trials).

317 incurred by this switching is less than the time saved by not running all training trials in full, allowing
 318 OPQA to trial more PCs overall. Along with the parallelism-informed GP prior beliefs to efficiently
 319 filter out suboptimal PCs, OPQA returns a PC with a higher throughput, which in turn allows more
 320 training steps to be run in the subsequent training even after just 20 minutes of PC optimization, as
 321 shown in Fig. 5b. This demonstrates the necessity of OPQA to achieve faster parallel training.

322 **Flexibility across different parallel training frameworks and NN architectures.** In App. I.3 we
 323 run similar experiments as that in Fig. 4, except that we perform training on the NEMO framework
 324 [13] where we additionally optimize for the size of the sequence parallelism dimension [18] in
 325 the PC. Meanwhile, in App. I.4, we compared OPQA with other methods in optimizing the PC for
 326 Vision Transformers and Mixture of Experts (MoE). Optimizing PC for MoE training is particularly
 327 interesting since the PC now includes the size of expert parallelism dimension [27]. In all of these
 328 scenarios, OPQA outperforms the other baselines. This demonstrates that the surrogate models used
 329 by OPQA will not overfit to a specific training scenario, but instead can cater to a wide variety of
 330 training frameworks, various NN architectures, and even unseen hyperparameters in a PC.

331 **Accuracy of surrogate models.** Figs. 6a and 6b compare the throughput predicted by OPQA with the
 332 actual measured throughput. We see that even after 10 trials, the predictions made by our surrogate
 333 models already correlate well with the measured throughput. As we progress, the prediction becomes
 334 more accurate, especially among PCs with high throughput where more trials are being run, which
 335 allows OPQA to better identify the optimal PCs. In contrast, a cost model alone (Fig. 6c) can capture
 336 general trends of $\mathcal{R}(H)$ but not all complex intricacies of the throughput, while a GP alone (Fig. 6d) is
 337 unable to learn meaningful interpolations for $\mathcal{R}(H)$. The shortcomings of these two surrogate models
 338 create a mismatch between the predicted optimal PC and the actual optimal PC. We further discuss
 339 the quality of the surrogate models for the throughput and maximum memory usage in App. I.5.

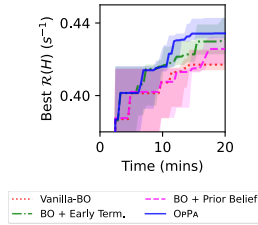


Figure 7: Effects of early termination and parallelism-informed GP prior beliefs on the ability of OPPA to optimize the throughput for Qwen2-1.5b training.

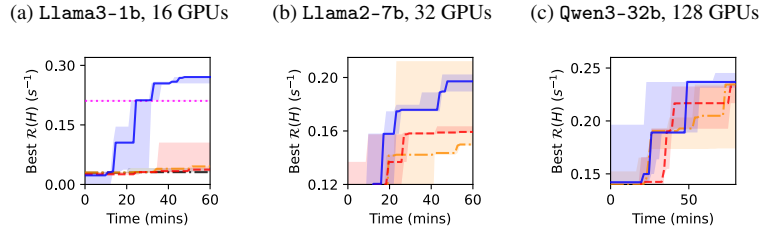


Figure 8: The best obtained throughput vs. the duration each method has been run for, when training NNs on multiple nodes. Figs. 8a and 8b was ran on COLOSSAL-AI [17], while Fig. 8c was ran on NEMO [13] with sizes of sequence and context parallelism dimensions also being tuned. The legend is the same as in Fig. 4. Some methods were not ran for results in Figs. 8b and 8c due to the available computational budget.

340 **Effects of each component in OP**PA. In Fig. 7, we performed ablation studies to isolate the effects
 341 from each proposed component in OPPA. We see that without early termination, BO would spend
 342 more time on suboptimal trials, resulting in a slower search process. Similarly, without the GP
 343 prior beliefs, we would be less informed about PCs which may be optimal, requiring more time to find
 344 the optimal PC. Additional results are presented in App. L.7. We also show the effects of q_{\min} on the
 345 performance of OPPA in App. L.8, demonstrating minimal degradation for small q_{\min} .

346 **Training on multi-node setups.** In addition to training on single-node setups, we also tested OPPA
 347 on optimizing the PC on multi-node setups. In these cases, the additional communication costs makes
 348 the throughput computation less straightforward, while larger number of feasible PCs complicates
 349 the optimization problem. In Fig. 8a, we show the results for optimizing the PC for training a smaller
 350 model [9] with 1 billion parameters on a commodity cluster with 16 GPUs, which is a typical setup
 351 found by practitioners with existing hardware in practice. Meanwhile, Figs. 8b and 8c are the results
 352 for optimizing the PC to train larger models [34, 40] on high-performance computing clusters with 32
 353 and 128 GPUs, respectively. In all of these cases, OPPA outperforms other baselines, and also does
 354 so more consistently (i.e., smaller variance in the resulting throughput), demonstrating its robustness.
 355 This trend is observed across different scales of hardware setup, and even across training frameworks
 356 or when including new parallelism dimensions in the PC (e.g., the experiment for Fig. 8c which
 357 includes six types of parallelism). We also see that vanilla BO outperforms other non-BO methods
 358 due to its ability to balance exploration and exploitation, however remaining worse than OPPA.

359 In Fig. 9 we also demonstrate that OPPA is additionally able to
 360 outperform cost model-based methods such as AMP [15] and
 361 NNSCALER [19]. This shows the superiority of adaptive meth-
 362 ods which can effectively incorporate the measured throughput
 363 from actual training trials as opposed to solely using cost
 364 models, and the non-trivial components exploited by OPPA that
 365 allow for this boost in performance.

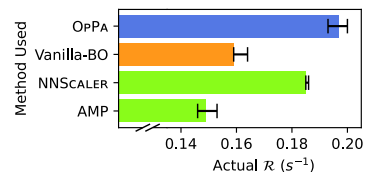


Figure 9: Obtained \mathcal{R} from OPPA compared to those from cost model-based methods for Llama2-7b training with 32 GPUs. The results reported for BO and OPPA are the same from Fig. 8b.

366 6 Conclusion

367 We have presented OPPA, which improves upon constrained
 368 Bayesian optimization by exploiting parallelism-informed GP
 369 prior beliefs and an early termination mechanism to efficiently
 370 optimize the parallelism configuration and achieve the highest
 371 training throughput across a variety of NN architectures, training frameworks, and hardware setups.
 372 As demonstrated, OPPA can also be easily adapted to optimize other unseen hyperparameters in a PC
 373 as well. We believe the parallelism-informed GP prior beliefs can be embedded with more domain
 374 knowledge of specific parallel training frameworks or knowledge of specific NN architectures, which
 375 would boost the performance of OPPA further.

376 **References**

- 377 [1] M. Balandat, B. Karrer, D. R. Jiang, S. Daulton, B. Letham, A. G. Wilson, and E. Bakshy.
 378 BoTorch: A Framework for Efficient Monte-Carlo Bayesian Optimization. In *Proc. NeurIPS*,
 379 2020.
- 380 [2] Z. Bian, Q. Xu, B. Wang, and Y. You. Maximizing Parallelism in Distributed Training for Huge
 381 Neural Networks. 2021. arXiv:2105.14450.
- 382 [3] T. Chen and C. Guestrin. XGBoost: A Scalable Tree Boosting System. In *Proc. SIGKDD*, 2016.
- 383 [4] J. Couto. Kernel k-means for categorical data. In *Proc. IDA*, 2005.
- 384 [5] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional
 385 transformers for language understanding. In *Proc. NAACL-HLT*, 2019.
- 386 [6] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani,
 387 M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby. An Image is Worth 16x16
 388 Words: Transformers for Image Recognition at Scale. In *Proc. ICLR*, 2021.
- 389 [7] P. I. Frazier. A Tutorial on Bayesian Optimization. 2018. arXiv:1807.02811.
- 390 [8] M. A. Gelbart, J. Snoek, and R. P. Adams. Bayesian optimization with unknown constraints. In
 391 *Proc. UAI*, 2014.
- 392 [9] A. Grattafiori, A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur,
 393 A. Schelten, A. Vaughan, A. Yang, A. Fan, A. Goyal, A. Hartshorn, A. Yang, A. Mitra,
 394 A. Sravankumar, A. Korenev, A. Hinsvark, A. Rao, A. Zhang, A. Rodriguez, A. Gregerson,
 395 A. Spataru, B. Roziere, B. Biron, B. Tang, B. Chern, C. Caucheteux, C. Nayak, C. Bi, C. Marra,
 396 C. McConnell, C. Keller, C. Touret, C. Wu, C. Wong, C. C. Ferrer, C. Nikolaidis, D. Allonsius,
 397 D. Song, D. Pintz, D. Livshits, D. Wyatt, D. Esiobu, D. Choudhary, D. Mahajan, D. Garcia-
 398 Olano, D. Perino, D. Hupkes, E. Lakomkin, E. AlBadawy, E. Lobanova, E. Dinan, E. M. Smith,
 399 F. Radenovic, F. Guzmán, F. Zhang, G. Synnaeve, G. Lee, G. L. Anderson, G. Thattai, G. Nail,
 400 G. Mialon, G. Pang, G. Cucurell, H. Nguyen, H. Korevaar, H. Xu, H. Touvron, I. Zarov, I. A.
 401 Ibarra, I. Kloumann, I. Misra, I. Evtimov, J. Zhang, J. Copet, J. Lee, J. Geffert, J. Vranes, J. Park,
 402 J. Mahadeokar, J. Shah, J. van der Linde, J. Billock, J. Hong, J. Lee, J. Fu, J. Chi, J. Huang,
 403 J. Liu, J. Wang, J. Yu, J. Bitton, J. Spisak, J. Park, J. Rocca, J. Johnstun, J. Saxe, J. Jia, K. V.
 404 Alwala, K. Prasad, K. Upasani, K. Plawiak, K. Li, K. Heafield, K. Stone, K. El-Arini, K. Iyer,
 405 K. Malik, K. Chiu, K. Bhalla, K. Lakhotia, L. Rantala-Yearly, L. van der Maaten, L. Chen,
 406 L. Tan, L. Jenkins, L. Martin, L. Madaan, L. Malo, L. Blecher, L. Landzaat, L. de Oliveira,
 407 M. Muzzi, M. Pasupuleti, M. Singh, M. Paluri, M. Kardas, M. Tsimpoukelli, M. Oldham,
 408 M. Rita, M. Pavlova, M. Kambadur, M. Lewis, M. Si, M. K. Singh, M. Hassan, N. Goyal,
 409 N. Torabi, N. Bashlykov, N. Bogoychev, N. Chatterji, N. Zhang, O. DuChenne, O. Çelebi,
 410 P. Alrassy, P. Zhang, P. Li, P. Vasic, P. Weng, P. Bhargava, P. Dubal, P. Krishnan, P. S. Koura,
 411 P. Xu, Q. He, Q. Dong, R. Srinivasan, R. Ganapathy, R. Calderer, R. S. Cabral, R. Stojnic,
 412 R. Raileanu, R. Maheswari, R. Girdhar, R. Patel, R. Sauvestre, R. Polidoro, R. Sumbaly,
 413 R. Taylor, R. Silva, R. Hou, R. Wang, S. Hosseini, S. Chennabasappa, S. Singh, S. Bell, S. S.
 414 Kim, S. Edunov, S. Nie, S. Narang, S. Raparthy, S. Shen, S. Wan, S. Bhosale, S. Zhang,
 415 S. Vandenhende, S. Batra, S. Whitman, S. Sootla, S. Collot, S. Gururangan, S. Borodinsky,
 416 T. Herman, T. Fowler, T. Sheasha, T. Georgiou, T. Scialom, T. Speckbacher, T. Mihaylov,
 417 T. Xiao, U. Karn, V. Goswami, V. Gupta, V. Ramanathan, V. Kerkez, V. Gonguet, V. Do,
 418 V. Vogeti, V. Albiero, V. Petrovic, W. Chu, W. Xiong, W. Fu, W. Meers, X. Martinet, X. Wang,
 419 X. Wang, X. E. Tan, X. Xia, X. Xie, X. Jia, X. Wang, Y. Goldschlag, Y. Gaur, Y. Babaei, Y. Wen,
 420 Y. Song, Y. Zhang, Y. Li, Y. Mao, Z. D. Coudert, Z. Yan, Z. Chen, Z. Papakipos, A. Singh,
 421 A. Srivastava, A. Jain, A. Kelsey, A. Shajnfeld, A. Gangidi, A. Victoria, A. Goldstand, A. Menon,
 422 A. Sharma, A. Boesenberg, A. Baeviski, A. Feinstein, A. Kallet, A. Sangani, A. Teo, A. Yunus,
 423 A. Lupu, A. Alvarado, A. Caples, A. Gu, A. Ho, A. Poulton, A. Ryan, A. Ramchandani,
 424 A. Dong, A. Franco, A. Goyal, A. Saraf, A. Chowdhury, A. Gabriel, A. Bharambe, A. Eisenman,
 425 A. Yazdan, B. James, B. Maurer, B. Leonhardi, B. Huang, B. Loyd, B. D. Paola, B. Paranjape,
 426 B. Liu, B. Wu, B. Ni, B. Hancock, B. Wasti, B. Spence, B. Stojkovic, B. Gamido, B. Montalvo,
 427 C. Parker, C. Burton, C. Mejia, C. Liu, C. Wang, C. Kim, C. Zhou, C. Hu, C.-H. Chu, C. Cai,
 428 C. Tindal, C. Feichtenhofer, C. Gao, D. Civin, D. Beaty, D. Kreymer, D. Li, D. Adkins, D. Xu,
 429 D. Testuggine, D. David, D. Parikh, D. Liskovich, D. Foss, D. Wang, D. Le, D. Holland,
 430 E. Dowling, E. Jamil, E. Montgomery, E. Presani, E. Hahn, E. Wood, E.-T. Le, E. Brinkman,
 431 E. Arcaute, E. Dunbar, E. Smothers, F. Sun, F. Kreuk, F. Tian, F. Kokkinos, F. Ozgenel,

- 432 F. Caggioni, F. Kanayet, F. Seide, G. M. Florez, G. Schwarz, G. Badeer, G. Swee, G. Halpern,
433 G. Herman, G. Sizov, Guangyi, Zhang, G. Lakshminarayanan, H. Inan, H. Shojanazeri, H. Zou,
434 H. Wang, H. Zha, H. Habeeb, H. Rudolph, H. Suk, H. Aspegren, H. Goldman, H. Zhan,
435 I. Damlaj, I. Molybog, I. Tufanov, I. Leontiadis, I.-E. Veliche, I. Gat, J. Weissman, J. Geboski,
436 J. Kohli, J. Lam, J. Asher, J.-B. Gaya, J. Marcus, J. Tang, J. Chan, J. Zhen, J. Reizenstein,
437 J. Teboul, J. Zhong, J. Jin, J. Yang, J. Cummings, J. Carvill, J. Shepard, J. McPhie, J. Torres,
438 J. Ginsburg, J. Wang, K. Wu, K. H. U, K. Saxena, K. Khandelwal, K. Zand, K. Matosich,
439 K. Veeraraghavan, K. Michelena, K. Li, K. Jagadeesh, K. Huang, K. Chawla, K. Huang,
440 L. Chen, L. Garg, L. A. L. Silva, L. Bell, L. Zhang, L. Guo, L. Yu, L. Moshkovich, L. Wehrstedt,
441 M. Khabsa, M. Avalani, M. Bhatt, M. Mankus, M. Hasson, M. Lennie, M. Reso, M. Groshev,
442 M. Naumov, M. Lathi, M. Keneally, M. Liu, M. L. Seltzer, M. Valko, M. Restrepo, M. Patel,
443 M. Vyatskov, M. Samvelyan, M. Clark, M. Macey, M. Wang, M. J. Hermoso, M. Metanat,
444 M. Rastegari, M. Bansal, N. Santhanam, N. Parks, N. White, N. Bawa, N. Singhal, N. Egebo,
445 N. Usunier, N. Mehta, N. P. Laptev, N. Dong, N. Cheng, O. Chernoguz, O. Hart, O. Salpekar,
446 O. Kalinli, P. Kent, P. Parekh, P. Saab, P. Balaji, P. Rittner, P. Bontrager, P. Roux, P. Dollar,
447 P. Zvyagina, P. Ratanchandani, P. Yuvraj, Q. Liang, R. Alao, R. Rodriguez, R. Ayub, R. Murthy,
448 R. Nayani, R. Mitra, R. Parthasarathy, R. Li, R. Hogan, R. Battey, R. Wang, R. Howes, R. Rinott,
449 S. Mehta, S. Siby, S. J. Bondu, S. Datta, S. Chugh, S. Hunt, S. Dhillon, S. Sidorov, S. Pan,
450 S. Mahajan, S. Verma, S. Yamamoto, S. Ramaswamy, S. Lindsay, S. Lindsay, S. Feng, S. Lin,
451 S. C. Zha, S. Patil, S. Shankar, S. Zhang, S. Zhang, S. Wang, S. Agarwal, S. Sajuyigbe,
452 S. Chintala, S. Max, S. Chen, S. Kehoe, S. Satterfield, S. Govindaprasad, S. Gupta, S. Deng,
453 S. Cho, S. Virk, S. Subramanian, S. Choudhury, S. Goldman, T. Remez, T. Glaser, T. Best,
454 T. Koehler, T. Robinson, T. Li, T. Zhang, T. Matthews, T. Chou, T. Shaked, V. Vontimitta,
455 V. Ajayi, V. Montanez, V. Mohan, V. S. Kumar, V. Mangla, V. Ionescu, V. Poenaru, V. T.
456 Mihailescu, V. Ivanov, W. Li, W. Wang, W. Jiang, W. Bouaziz, W. Constable, X. Tang, X. Wu,
457 X. Wang, X. Wu, X. Gao, Y. Kleinman, Y. Chen, Y. Hu, Y. Jia, Y. Qi, Y. Li, Y. Zhang,
458 Y. Zhang, Y. Adi, Y. Nam, Yu, Wang, Y. Zhao, Y. Hao, Y. Qian, Y. Li, Y. He, Z. Rait, Z. DeVito,
459 Z. Rosnbrick, Z. Wen, Z. Yang, Z. Zhao, and Z. Ma. The Llama 3 Herd of Models. 2024.
- 460 [10] Y. Huang, Y. Cheng, A. Bapna, O. Firat, M. X. Chen, D. Chen, H. Lee, J. Ngiam, Q. V. Le,
461 Y. Wu, and Z. Chen. GPipe: efficient training of giant neural networks using pipeline parallelism.
462 In *Proc. NeurIPS*, 2019.
- 463 [11] D. R. Jones, M. Schonlau, and W. J. Welch. Efficient Global Optimization of Expensive
464 Black-Box Functions. *Journal of Global Optimization*, 13(4):455–492, 1998.
- 465 [12] J. Kirschner and A. Krause. Information Directed Sampling and Bandits with Heteroscedastic
466 Noise. In *Proc. COLT*, 2018.
- 467 [13] O. Kuchaiev, J. Li, H. Nguyen, O. Hrinchuk, R. Leary, B. Ginsburg, S. Kriman, S. Beliaev,
468 V. Lavrukhin, J. Cook, P. Castonguay, M. Popova, J. Huang, and J. M. Cohen. NeMo: a toolkit
469 for building AI applications using Neural Modules. 2019. arXiv:1909.09577.
- 470 [14] B. Lei, T. Q. Kirk, A. Bhattacharya, D. Pati, X. Qian, R. Arroyave, and B. K. Mallick. Bayesian
471 optimization with adaptive surrogate models for automated experimental design. *npj Computa-
472 tional Materials*, 7(1):1–12, Dec. 2021.
- 473 [15] D. Li, H. Wang, E. Xing, and H. Zhang. AMP: automatically finding model parallel strategies
474 with heterogeneity awareness. In *Proc. NeurIPS*, 2022.
- 475 [16] S. Li, Y. Zhao, R. Varma, O. Salpekar, P. Noordhuis, T. Li, A. Paszke, J. Smith, B. Vaughan,
476 P. Damania, and S. Chintala. PyTorch distributed: experiences on accelerating data parallel
477 training. In *Proc. VLDB*, 2020.
- 478 [17] S. Li, H. Liu, Z. Bian, J. Fang, H. Huang, Y. Liu, B. Wang, and Y. You. Colossal-AI: A Unified
479 Deep Learning System For Large-Scale Parallel Training. In *Proc. ICPP*, 2023.
- 480 [18] S. Li, F. Xue, C. Baranwal, Y. Li, and Y. You. Sequence parallelism: Long sequence training
481 from system perspective. In *Proc. ACL*, 2023.
- 482 [19] Z. Lin, Y. Miao, Q. Zhang, F. Yang, Y. Zhu, C. Li, S. Maleki, X. Cao, N. Shang, Y. Yang, W. Xu,
483 M. Yang, L. Zhang, and L. Zhou. nnScaler: constraint-guided parallelization plan generation
484 for deep learning training. In *Proc. OSDI*, 2024.
- 485 [20] H. Liu, C. Li, Q. Wu, and Y. J. Lee. Visual instruction tuning. In *Proc. NeurIPS*, 2023.

- 486 [21] A. Makarova, I. Usmanova, I. Bogunovic, and A. Krause. Risk-averse Heteroscedastic Bayesian
487 Optimization. In *Proc. NeurIPS*, 2021.
- 488 [22] D. Narayanan, A. Harlap, A. Phanishayee, V. Seshadri, N. R. Devanur, G. R. Ganger, P. B.
489 Gibbons, and M. Zaharia. PipeDream: generalized pipeline parallelism for DNN training. In
490 *Proc. SOSP*, 2019.
- 491 [23] OpenAI, J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida,
492 J. Altschmidt, S. Altman, S. Anadkat, R. Avila, I. Babuschkin, S. Balaji, V. Balcom, P. Bal-
493 tescu, H. Bao, M. Bavarian, J. Belgum, I. Bello, J. Berdine, G. Bernadett-Shapiro, C. Berner,
494 L. Bogdonoff, O. Boiko, M. Boyd, A.-L. Brakman, G. Brockman, T. Brooks, M. Brundage,
495 K. Button, T. Cai, R. Campbell, A. Cann, B. Carey, C. Carlson, R. Carmichael, B. Chan,
496 C. Chang, F. Chantzis, D. Chen, S. Chen, R. Chen, J. Chen, M. Chen, B. Chess, C. Cho,
497 C. Chu, H. W. Chung, D. Cummings, J. Currier, Y. Dai, C. Decareaux, T. Degry, N. Deutsch,
498 D. Deville, A. Dhar, D. Dohan, S. Dowling, S. Dunning, A. Ecoffet, A. Eleti, T. Eloundou,
499 D. Farhi, L. Fedus, N. Felix, S. P. Fishman, J. Forte, I. Fulford, L. Gao, E. Georges, C. Gibson,
500 V. Goel, T. Gogineni, G. Goh, R. Gontijo-Lopes, J. Gordon, M. Grafstein, S. Gray, R. Greene,
501 J. Gross, S. S. Gu, Y. Guo, C. Hallacy, J. Han, J. Harris, Y. He, M. Heaton, J. Heidecke, C. Hesse,
502 A. Hickey, W. Hickey, P. Hoeschele, B. Houghton, K. Hsu, S. Hu, X. Hu, J. Huizinga, S. Jain,
503 S. Jain, J. Jang, A. Jiang, R. Jiang, H. Jin, D. Jin, S. Jomoto, B. Jonn, H. Jun, T. Kaftan, L. Kaiser,
504 A. Kamali, I. Kanitscheider, N. S. Keskar, T. Khan, L. Kilpatrick, J. W. Kim, C. Kim, Y. Kim,
505 J. H. Kirchner, J. Kiros, M. Knight, D. Kokotajlo, L. Kondraciuk, A. Kondrich, A. Konstan-
506 tinidis, K. Kosic, G. Krueger, V. Kuo, M. Lampe, I. Lan, T. Lee, J. Leike, J. Leung, D. Levy,
507 C. M. Li, R. Lim, M. Lin, S. Lin, M. Litwin, T. Lopez, R. Lowe, P. Lue, A. Makanju, K. Mal-
508 facini, S. Manning, T. Markov, Y. Markovski, B. Martin, K. Mayer, A. Mayne, B. McGrew,
509 S. M. McKinney, C. McLeavey, P. McMillan, J. McNeil, D. Medina, A. Mehta, J. Menick,
510 L. Metz, A. Mishchenko, P. Mishkin, V. Monaco, E. Morikawa, D. Mossing, T. Mu, M. Murati,
511 O. Murk, D. Mély, A. Nair, R. Nakano, R. Nayak, A. Neelakantan, R. Ngo, H. Noh, L. Ouyang,
512 C. O’Keefe, J. Pachocki, A. Paino, J. Palermo, A. Pantuliano, G. Parascandolo, J. Parish,
513 E. Parparita, A. Passos, M. Pavlov, A. Peng, A. Perelman, F. d. A. B. Peres, M. Petrov, H. P. d. O.
514 Pinto, Michael, Pokorny, M. Pokrass, V. H. Pong, T. Powell, A. Power, B. Power, E. Proehl,
515 R. Puri, A. Radford, J. Rae, A. Ramesh, C. Raymond, F. Real, K. Rimbach, C. Ross, B. Rotsted,
516 H. Roussez, N. Ryder, M. Saltarelli, T. Sanders, S. Santurkar, G. Sastry, H. Schmidt, D. Schnurr,
517 J. Schulman, D. Selsam, K. Sheppard, T. Sherbakov, J. Shieh, S. Shoker, P. Shyam, S. Sidor,
518 E. Sigler, M. Simens, J. Sitkin, K. Slama, I. Sohl, B. Sokolowsky, Y. Song, N. Staudacher, F. P.
519 Such, N. Summers, I. Sutskever, J. Tang, N. Tezak, M. B. Thompson, P. Tillet, A. Tootoonchian,
520 E. Tseng, P. Tuggle, N. Turley, J. Tworek, J. F. C. Uribe, A. Vallone, A. Vijayvergiya, C. Voss,
521 C. Wainwright, J. J. Wang, A. Wang, B. Wang, J. Ward, J. Wei, C. J. Weinmann, A. Welihinda,
522 P. Welinder, J. Weng, L. Weng, M. Wiethoff, D. Willner, C. Winter, S. Wolrich, H. Wong,
523 L. Workman, S. Wu, J. Wu, M. Wu, K. Xiao, T. Xu, S. Yoo, K. Yu, Q. Yuan, W. Zaremba,
524 R. Zellers, C. Zhang, M. Zhang, S. Zhao, T. Zheng, J. Zhuang, W. Zhuk, and B. Zoph. GPT-4
525 Technical Report. 2024. arXiv:2303.08774.
- 526 [24] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,
527 P. Mishkin, J. Clark, G. Krueger, and I. Sutskever. Learning Transferable Visual Models From
528 Natural Language Supervision. In *Proc. ICML*, 2021.
- 529 [25] T. Rainforth, A. Foster, D. R. Ivanova, and F. B. Smith. Modern Bayesian Experimental Design.
530 2023. arXiv:2302.14545.
- 531 [26] S. Rajbhandari, J. Rasley, O. Ruwase, and Y. He. ZeRO: memory optimizations toward training
532 trillion parameter models. In *Proc. SC*, 2020.
- 533 [27] S. Rajbhandari, C. Li, Z. Yao, M. Zhang, R. Y. Aminabadi, A. A. Awan, J. Rasley, and Y. He.
534 DeepSpeed-moe: Advancing mixture-of-experts inference and training to power next-generation
535 ai scale. In *Proc. ICML*, 2022.
- 536 [28] J. Rasley, S. Rajbhandari, O. Ruwase, and Y. He. DeepSpeed: System Optimizations Enable
537 Training Deep Learning Models with Over 100 Billion Parameters. In *Proc. SIGKDD*, 2020.
- 538 [29] C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press,
539 2006. ISBN 978-0-262-18253-9.
- 540 [30] N. Shazeer, A. Mirhoseini, K. Maziarz, A. Davis, Q. Le, G. Hinton, and J. Dean. Outrageously
541 large neural networks: The sparsely-gated mixture-of-experts layer. In *Proc. ICLR*, 2017.

- 542 [31] M. Shoenybi, M. Patwary, R. Puri, P. LeGresley, J. Casper, and B. Catanzaro. Megatron-
543 LM: Training Multi-Billion Parameter Language Models Using Model Parallelism. 2020.
544 arXiv:1909.08053.
- 545 [32] J. Snoek, H. Larochelle, and R. P. Adams. Practical Bayesian optimization of machine learning
546 algorithms. In *Proc. NeurIPS*, 2012.
- 547 [33] N. Srinivas, A. Krause, S. M. Kakade, and M. Seeger. Gaussian Process Optimization in the
548 Bandit Setting: No Regret and Experimental Design. *IEEE Transactions on Information Theory*,
549 58:3250–3265, 2012.
- 550 [34] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra,
551 P. Bhargava, S. Bhosale, D. Bikel, L. Blecher, C. C. Ferrer, M. Chen, G. Cucurull, D. Esiobu,
552 J. Fernandes, J. Fu, W. Fu, B. Fuller, C. Gao, V. Goswami, N. Goyal, A. Hartshorn, S. Hosseini,
553 R. Hou, H. Inan, M. Kardas, V. Kerkez, M. Khabsa, I. Kloumann, A. Korenev, P. S. Koura, M.-A.
554 Lachaux, T. Lavril, J. Lee, D. Liskovich, Y. Lu, Y. Mao, X. Martinet, T. Mihaylov, P. Mishra,
555 I. Molybog, Y. Nie, A. Poulton, J. Reizenstein, R. Rungta, K. Saladi, A. Schelten, R. Silva, E. M.
556 Smith, R. Subramanian, X. E. Tan, B. Tang, R. Taylor, A. Williams, J. X. Kuan, P. Xu, Z. Yan,
557 I. Zarov, Y. Zhang, A. Fan, M. Kambadur, S. Narang, A. Rodriguez, R. Stojnic, S. Edunov, and
558 T. Scialom. Llama 2: Open Foundation and Fine-Tuned Chat Models. 2023.
- 559 [35] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and
560 I. Polosukhin. Attention is All you Need. In *Proc. NeurIPS*, 2017.
- 561 [36] M. Wagenländer, G. Li, B. Zhao, L. Mai, and P. Pietzuch. Tenplex: Dynamic Parallelism for
562 Deep Learning using Parallelizable Tensor Collections. In *Proc. SOSP*, 2024.
- 563 [37] J. T. Wilson, R. Moriconi, F. Hutter, and M. P. Deisenroth. The reparameterization trick for
564 acquisition functions. 2017.
- 565 [38] D. Xiong, L. Chen, Y. Jiang, D. Li, S. Wang, and S. Wang. Revisiting the Time Cost Model of
566 AllReduce. 2024. arXiv:2409.04202.
- 567 [39] A. Yang, B. Yang, B. Hui, B. Zheng, B. Yu, C. Zhou, C. Li, C. Li, D. Liu, F. Huang, G. Dong,
568 H. Wei, H. Lin, J. Tang, J. Wang, J. Yang, J. Tu, J. Zhang, J. Ma, J. Yang, J. Xu, J. Zhou, J. Bai,
569 J. He, J. Lin, K. Dang, K. Lu, K. Chen, K. Yang, M. Li, M. Xue, N. Ni, P. Zhang, P. Wang,
570 R. Peng, R. Men, R. Gao, R. Lin, S. Wang, S. Bai, S. Tan, T. Zhu, T. Li, T. Liu, W. Ge, X. Deng,
571 X. Zhou, X. Ren, X. Zhang, X. Wei, X. Ren, X. Liu, Y. Fan, Y. Yao, Y. Zhang, Y. Wan, Y. Chu,
572 Y. Liu, Z. Cui, Z. Zhang, Z. Guo, and Z. Fan. Qwen2 Technical Report. 2024.
- 573 [40] A. Yang, A. Li, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Gao, C. Huang, C. Lv, C. Zheng,
574 D. Liu, F. Zhou, F. Huang, F. Hu, H. Ge, H. Wei, H. Lin, J. Tang, J. Yang, J. Tu, J. Zhang,
575 J. Yang, J. Yang, J. Zhou, J. Zhou, J. Lin, K. Dang, K. Bao, K. Yang, L. Yu, L. Deng, M. Li,
576 M. Xue, M. Li, P. Zhang, P. Wang, Q. Zhu, R. Men, R. Gao, S. Liu, S. Luo, T. Li, T. Tang,
577 W. Yin, X. Ren, X. Wang, X. Zhang, X. Ren, Y. Fan, Y. Su, Y. Zhang, Y. Zhang, Y. Wan, Y. Liu,
578 Z. Wang, Z. Cui, Z. Zhang, Z. Zhou, and Z. Qiu. Qwen3 Technical Report. 2025.
- 579 [41] S. Zhang, L. Diao, C. Wu, Z. Cao, S. Wang, and W. Lin. HAP: SPMD DNN Training on
580 Heterogeneous GPU Clusters with Automated Program Synthesis. In *Proc. EuroSys*, 2024.
- 581 [42] Y. Zhang, D. W. Apley, and W. Chen. Bayesian Optimization for Materials Design with Mixed
582 Quantitative and Qualitative Variables. *Scientific Reports*, (1):4924, 2020.
- 583 [43] Y. Zhao, A. Gu, R. Varma, L. Luo, C.-C. Huang, M. Xu, L. Wright, H. Shojanazeri, M. Ott,
584 S. Shleifer, A. Desmaison, C. Balioglu, P. Damania, B. Nguyen, G. Chauhan, Y. Hao, A. Mathews,
585 and S. Li. PyTorch FSDP: Experiences on Scaling Fully Sharded Data Parallel. In *Proc.*
586 *VLDB*, 2023.
- 587 [44] L. Zheng, Z. Li, H. Zhang, Y. Zhuang, Z. Chen, Y. Huang, Y. Wang, Y. Xu, D. Zhuo, E. P.
588 Xing, J. E. Gonzalez, and I. Stoica. Alpa: Automating Inter- and Intra-Operator Parallelism for
589 Distributed Deep Learning. In *Proc. OSDI*, 2022.

590 A Additional Discussion on Parallelism Types

591 **Data parallelism.** The most basic type of training parallelism is data parallelism (DP) [16] where a
592 batch of training data is split into shards and distributed among each device. These shards are then fed
593 into the local replicas of the models, before the parameter updates from each device are synchronized.
594 While simple and often the fastest, the naive DP approach requires replicating the model on each
595 device, which takes up additional storage on each node. Several methods have since been proposed to
596 perform DP with sharded models, including the Zero Redundancy Optimizer (ZERO) introduced by
597 [26] in the DEEPSPEED package, and Fully-Sharded Data Parallel (FSDP) introduced by [43]. While
598 these frameworks allow for efficient DP implementations, their effectiveness can still heavily depend
599 on the choice of hyperparameters. For example, ZERO involves three different stages of optimization
600 which chooses whether the optimizer states, the model gradients or the model parameters are sharded
601 between each GPU. The choice of sharded items affect the amount of data that has to be stored in
602 each GPU and communicated across GPUs which in turn affects the throughput of the training and
603 the memory usage in each GPU.

604 **Tensor parallelism.** Another method for scaling operation is tensor parallelism (TP) where individual
605 tensors are sharded across multiple devices, so that the matrix multiplication operations are instead
606 done in a distributed manner, allowing these operations to scale to larger sizes than otherwise that
607 would fit on a single GPU. TP initially involved splitting a tensor along a single dimension [31],
608 however has since also incorporated sharding tensors across multiple dimensions as well [2].

609 **Pipeline parallelism.** In pipeline parallelism (PP) [10, 22] we instead partition the model along
610 its execution pipeline with each model partition running synchronously with microbatches of data.
611 The gradients are accumulated for each microbatch and updated at the end of each training step. By
612 sharding the model and training data into smaller chunks, the GPU memory required at any one time
613 becomes lower, allowing for the training of larger models at the cost of more sequential operation
614 rounds and higher cost of communication between each GPU. The tradeoff between training speed
615 and maximum memory usage can be further controlled based on the size of microbatches and the
616 number of model chunks.

617 B Technical Primer on Gaussian Processes and Bayesian Optimization

618 In this section, we provide a technical overview of Gaussian process (GP) regression and on Bayesian
619 optimization (BO). The contents are adapted from [7, 29].

620 A Gaussian process (GP) $\mathcal{GP}(\mu_{\text{prior}}, k)$ with prior mean μ_{prior} and covariance k is a random process
621 where for any subset of input \mathbf{X} , its corresponding output is given by a normal distribution $f(\mathbf{X}) \sim$
622 $\mathcal{N}(\mu_{\text{prior}}(\mathbf{X}), k(\mathbf{X}, \mathbf{X}))$. The prior mean $\mu_{\text{prior}}(x)$ describes the expected value of the random function
623 $f(x)$ at a certain input, while the prior covariance $k(x, x')$ roughly captures that between $f(x)$ and
624 $f(x')$.

625 Suppose that we have an unknown function f drawn from the GP. Given a set $D = (\mathbf{X}, \mathbf{y}) =$
626 $\{(x_1, y_1), \dots, (x_n, y_n)\}$ of inputs and corresponding measurements, where $y_i = f(x_i) + \epsilon_i$ are noisy
627 measurements of the actual function with Gaussian noise $\epsilon_i \sim \mathcal{N}(0, \lambda_i)$. Then, when performing
628 Bayesian inference, we can express the GP posterior mean and covariance as

$$\mu(x) = \mu_{\text{prior}}(x) + k(x, \mathbf{X})(k(\mathbf{X}) + \text{diag}(\lambda))^{-1}(\mathbf{y} - \mu_{\text{prior}}(x)), \quad (4)$$

$$\sigma^2(x) = k(x, x) - k(x, \mathbf{X})(k(\mathbf{X}) + \text{diag}(\lambda))^{-1}k(\mathbf{X}, x). \quad (5)$$

629 In practice, the prior mean and covariance may have hyperparameters θ which specify what functions
630 it is able to model. For example, many kernels representing the prior covariance include lengthscale
631 values, which govern how correlated the function output is when a certain input dimension changes.
632 One method to find the optimal hyperparameters for the kernel is by finding the hyperparameters
633 which maximize the marginal log-likelihood.

634 In Bayesian optimization (BO), the goal is to find the maxima of the unknown function f . This
635 function is black-box with no analytical form. To do so, we can learn more about f by querying it at
636 different inputs, and perform Bayesian inference to update our belief on the unknown function.

637 Given the measurements $\mathcal{D}_t = (\mathbf{X}_t, \mathbf{y}_t) = \{(x_1, y_1), \dots, (x_t, y_t)\}$ in round t of data selection, GP
638 regression can be performed to obtain a posterior mean μ_t and variance σ_t^2 . The next input to query
639 x_{t+1} can be chosen as the input which maximizes some acquisition function. Examples of such
640 acquisition function include the expected improvement [11]

$$\text{EI}_t(x) = \mathbb{E}_{y' \sim \mathcal{N}(\mu_{t-1}(x), \sigma_{t-1}^2(x))} \left[\max(0, y' - \max_{y \in \mathbf{y}_{t-1}} y) \right] \quad (6)$$

641 or the upper confidence bound [33]

$$\text{UCB}_t(x) = \mu_{t-1}(x) + \beta_t \sigma_{t-1}(x) \quad (7)$$

642 where $\beta_t > 0$ is a constant that may vary with t . In all of these acquisition functions, a tradeoff is
643 performed between selecting inputs that the GP is uncertain about (i.e., with high $\sigma_{t-1}^2(x)$) to learn
644 more about those unknown region, and selecting inputs in regions where the function value is known
645 to be higher (i.e., with high $\mu_{t-1}(x)$).

646 C Detailed Discussion on Problem Setting

647 C.1 Hyperparameters Considered

648 In Table 1, we list several hyperparameters which we include in our parallelism configuration and the
649 range of the values. Note that the hyperparameters are constrained to give a valid PC as well; for
650 example, we ensure that $\text{dp} \cdot \text{tp} \cdot \text{pp} = \text{n_gpus}$ to ensure that each dimension do not exceed number
651 of GPUs. Some hyperparameters are also set to their default value when not in use; for example, if
652 $\text{pp} = 1$ (i.e., no PP used) then we restrict $\text{mb} = \text{mc} = 1$ such that PCs are not duplicated. In the code,
653 we generate all possible PCs beforehand so we can ensure that all PCs chosen will be valid according
654 to the constraints.

Table 1: Tunable hyperparameters in a parallelism configuration.

Hyperparameter	Description	Feasible Values	Notes
DP size (dp)	Data parallelism degree	[1, n_gpus]	Requires $\text{dp} \cdot \text{tp} \cdot \text{pp} = \text{n_gpus}$ Only NEMO and must be $\leq \text{tp}$ Only NEMO and only for some experiments Only NEMO and only for MoE
TP size (tp)	Tensor parallelism degree	[1, n_gpus]	
PP size (pp)	Pipeline parallelism degree	[1, n_gpus]	
SP size	Sequence parallelism degree	[1, n_gpus]	
CP size	Context parallelism degree	[1, n_gpus]	
EP size	Expert parallelism degree	[1, n_gpus]	
DP bucket size	Size for gradient reduction buckets (MB)	[1, 4096]	Only COLOSSAL-AI For COLOSSAL-AI, must be ≤ 2 For NEMO, this controls FSDP sharding instead
ZeRO stage	ZeRO stage used	[0, 3]	
ZeRO bucket size	Bucket size for ZeRO communication	[1, 4096]	
Overlap ZeRO communication	Whether to overlap ZeRO communication	True / False	
Overlap DP AllGather	Whether to overlap AllGather	True / False	Only considered for COLOSSAL-AI
# microbatches (mb)	Number of microbatches per forward pass	\leq batch size	
# model chunks (mc)	Number of model chunks for pipelining	\leq # transformer blocks	
Overlap P2P for PP	Overlap PP communication or not	True / False	
Grad. checkpointing	Whether gradient checkpointing is enabled	True / False	

655 C.2 Throughput vs. Time Per Training Step

656 We explain why we choose to maximize throughput instead of minimizing time per training step.

657 As an example, suppose we consider three PCs H, H', H'' where the times per training step are given
658 by $\mathcal{T}(H) = 0.3$, $\mathcal{T}(H') = 0.4$, and $\mathcal{T}(H'') = 0.5$. In this case, H would be the best PC out of the
659 three. We see here the gap of the time per training step between H and H' is $\mathcal{T}(H') - \mathcal{T}(H) = 0.1$,
660 and the same gap size for H' and H'' of $\mathcal{T}(H'') - \mathcal{T}(H') = 0.1$. Meanwhile, the gap between the
661 throughput of the two PCs would be $\mathcal{R}(H) - \mathcal{R}(H') = \mathcal{T}(H)^{-1} - \mathcal{T}(H')^{-1} = 3.\bar{3} - 2.5 = 0.8\bar{3}$
662 and $\mathcal{R}(H') - \mathcal{R}(H'') = \mathcal{T}(H')^{-1} - \mathcal{T}(H'')^{-1} = 2.5 - 2 = 0.5$. We can see that the gap between the
663 best PC becomes enhanced when we consider the throughput, when we compare it with the relative
664 gap size of the training step time.

665 More concretely, if we have a PC which requires time t per training step, then you can reduce it by an
666 amount of Δt , then the throughput would have increased by an amount $(\Delta t)/t^2$. When t becomes
667 smaller, the change in throughput will also increase but at an increasing rate. This therefore means by
668 modeling the throughput, the scores of the good PCs will be more clearly separated.

669 Additionally, we also consider a maximization of throughput since throughput would be bounded by
670 $[0, r_{\max}]$, rather than the time per training step which would be unbounded on one end, i.e., be in the

671 interval $[t_{\min}, \infty)$ which makes suboptimal PCs easier to handle. Also, we frame our problem as a
 672 maximization in order to be consistent with Bayesian optimization works which typically considers a
 673 maximization problem.

674 D Detailed Discussion of Surrogate Models used by OPPA

675 D.1 Domain Knowledge

676 In this section, we elaborate on how the throughput prior mean is constructed in order to obtain the
 677 form for the parallelism-informed prior mean. In summary, we design the prior to incorporate the
 678 following characteristics.

- 679 • *Parallelism coverage.* We model DP/TP traffic via All-Reduce-style collectives and PP via
 680 point-to-point transfers, using a placement-aware split of intra- and inter-node links. For
 681 computation costs, we also explicitly consider the pipeline bubbles.
- 682 • *Topology and protocol awareness.* The communication term uses canonical ring/tree scaling
 683 with hierarchical aggregation (node-local first, then cross-node), while collapsing ZERO
 684 stages into a single All-Reduce operation at the bandwidth level (bytes over the wire are of
 685 the same order) and letting stage-specific latency/overlap differences be absorbed by C (e.g.,
 686 event counts, bucket sizes, communication overhead for each parallelism dimension, etc.).
- 687 • *Hardware agnosticism via learning.* Rather than hard-coding device/network constants, we
 688 expose a small set of effective coefficients that are learned from a few traces. This keeps
 689 the prior portable across models, data types, and interconnects while preserving the correct
 690 asymptotic trends for dp, tp, and pp.

691 To predict the computation time, we assume an idealized machine with infinitely many processors
 692 such that DP and TP can be perfectly parallelized. Meanwhile, PP using an interleaved schedule
 693 incurs additional computation time from the microbatches being ran sequentially, and from pipeline
 694 bubble when the first microbatch is being fed through the pipeline [22]. This additional computation
 695 time from PP, visualized in Fig. 10, is roughly equal to

$$\hat{\mathcal{T}}_{\text{comp}}(H; t_{\text{comp}}) = \frac{t_{\text{comp}}}{n_{\text{gpus}}} \cdot \left(\text{mb} + \frac{\text{pp} - 1}{\text{mc}} \right) \quad (8)$$

696 where mb is the number of microbatches used in PP (set to 1 when PP is not used), mc is the number
 697 of model chunks for PP (also set to 1 when PP is not used), and $t_c = t_f + t_b$ is the total time to
 698 perform the forward and backward passes.

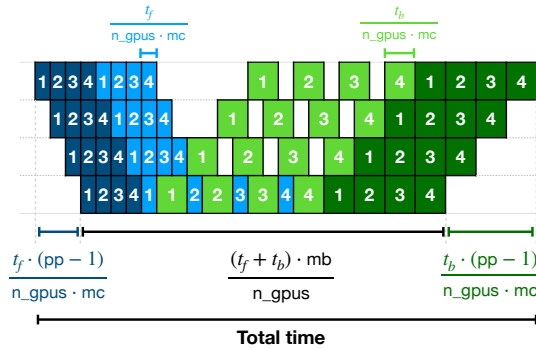


Figure 10: Predicted computation time for PP where t_f and t_b are the time required for the forward and backward stages respectively, for when $\text{pp} = n_{\text{gpus}} = 4$, $\text{mb} = 4$, and $\text{mc} = 2$

699 To predict the communication time, inspired by the model visualized in Fig. 2, we assume that DP
 700 and TP involve All-Reduce communications and PP P2P communications, where all these costs are
 701 modeled separately. We use an extended (α, β, γ) model discussed in [38], and for the intra- and
 702 inter-node communications, we characterize the network performance by the latency α_{intra} and α_{inter} ,

703 the per-byte bandwidth cost β_{intra} and β_{inter} , the incast overhead γ_{intra} and γ_{inter} , and the memory
 704 access overhead δ_{intra} and δ_{inter} . We assume that the inter-node communication costs will be larger
 705 than their intra-node counterparts.

706 For DP and TP, we assume that the Ring All-Reduce implementation is used where the cost to gather
 707 and scatter data of size N across D GPUs is given by

$$C_{\text{AR}}(D, N, \alpha, \beta, \gamma, \delta) = 2(D - 1)\alpha + \frac{D - 1}{D}N(2\beta + \gamma + 3\delta). \quad (9)$$

708 For DP, we consider the gradient synchronization from the All-Reduce procedure. The data size per
 709 GPU for DP All-Reduce, N_{dp} , is

$$N_{\text{dp}} = \lambda_Z \frac{M_{\text{model}}}{\text{tp} \cdot \text{pp}} \quad (10)$$

710 where M_{model} is a learnable total model parameter size, and λ_Z is a tiny fudge for ZERO flavor while
 711 staying in AR-land, and it accounts for param-All-Gather + grad-Reduce-Scatter volume equivalence
 712 with small overhead for ZERO-3.

713 Suppose G_{node} is the number of GPUs per node. The overall communication cost from DP $\widehat{\mathcal{T}}_{\text{comm,dp}}$
 714 would then depend on the configuration of the network as follows:

715 • In a **hierarchical** (multi-node scenario with $\text{dp} > G_{\text{node}}$) system, the cost is given by

$$\widehat{\mathcal{T}}_{\text{comm,dp}}(H; \mathbf{C}) = C_{\text{AR}}(G_{\text{node}}, N_{\text{dp}}, \alpha_{\text{intra}}, \beta_{\text{intra}}, \gamma_{\text{intra}}, \delta_{\text{intra}}) \\ + C_{\text{AR}}(\lceil \text{dp}/G_{\text{node}} \rceil, N_{\text{dp}}/G_{\text{node}}, \alpha_{\text{inter}}, \beta_{\text{inter}}, \gamma_{\text{inter}}, \delta_{\text{inter}}). \quad (11)$$

716 • In a **flat inter-node** (multi-node scenario with $\text{dp} \leq G_{\text{node}}$), the cost is given by

$$\widehat{\mathcal{T}}_{\text{comm,dp}}(H; \mathbf{C}) = C_{\text{AR}}(\text{dp}, N_{\text{dp}}, \alpha_{\text{inter}}, \beta_{\text{inter}}, \gamma_{\text{inter}}, \delta_{\text{inter}}). \quad (12)$$

717 • In a **flat intra-node** (single-node scenario), the cost is given by

$$\widehat{\mathcal{T}}_{\text{comm,dp}}(H; \mathbf{C}) = C_{\text{AR}}(\text{dp}, N_{\text{dp}}, \alpha_{\text{intra}}, \beta_{\text{intra}}, \gamma_{\text{intra}}, \delta_{\text{intra}}). \quad (13)$$

718 Depending on the hardware used, the appropriate cost for the scenario can be selected.

719 For TP, we consider the cost from frequent activation communication (e.g., All-Reduce per layer).
 720 Let $M_{\text{act,tp}}$ be a learnable characteristic data size for one such TP All-Reduce operation. Let $O_{\text{tp}/\text{mb}}$
 721 be the learnable number of these operations per microbatch. The total number of TP communi-
 722 cation operations is $N_{\text{tp,ops}} = O_{\text{tp}/\text{mb}} \cdot \text{mb}$, and the cost of a single TP All-Reduce operation is
 723 $C_{\text{AR}}(\text{tp}, M_{\text{act,tp}}, \alpha_{\text{eff}}, \beta_{\text{eff}}, \gamma_{\text{eff}}, \delta_{\text{eff}})$ where effective parameters $\alpha_{\text{eff}}, \beta_{\text{eff}}, \gamma_{\text{eff}}, \delta_{\text{eff}}$ are chosen as intra-
 724 node or inter-node based on whether the tp group spans multiple nodes (i.e., if $N_H > 1$ and
 725 $\text{tp} > G_{\text{node}}$) or not. Then, the total TP communication cost is the number of communication
 726 operations multiplied by the cost per communication operation, or

$$\widehat{\mathcal{T}}_{\text{comm,tp}}(H; \mathbf{C}) = N_{\text{tp,ops}} \cdot C_{\text{AR}}(\text{tp}, M_{\text{act,tp}}, \alpha_{\text{eff}}, \beta_{\text{eff}}, \gamma_{\text{eff}}, \delta_{\text{eff}}). \quad (14)$$

727 For PP, we consider the point-to-point (P2P) transfers of activations and gradients between pp pipeline
 728 stages where the cost to transfer data of size N is given by

$$C_{\text{P2P}}(N, \alpha, \beta) = \alpha + N \cdot \beta. \quad (15)$$

729 Let $M_{\text{act,pp}}$ be a learnable characteristic data size for one P2P transfer (e.g., the size of an activation
 730 tensor). The number of communication boundaries is $\text{pp} - 1$. Communication occurs for each of
 731 mb microbatch, in both forward (activations) and backward (gradients) directions. The total number
 732 of P2P communications is given by $N_{\text{pp,transfers}} = \max(\text{pp} - 1, 0) \cdot 2\text{mb}$ where a single P2P transfer
 733 uses effective latency α_{eff} and effective per-unit-data cost $\beta_{\text{pp,eff}}$ (derived from base β parameters),
 734 chosen as intra-node or inter-node based on whether communicating stages are on different nodes
 735 (approximated if $N_H > 1$). Given this, the overall cost of all communications related to PP would be
 736 given by

$$\widehat{\mathcal{T}}_{\text{comm,pp}}(H; \mathbf{C}) = N_{\text{pp,transfers}} \cdot C_{\text{P2P}}(M_{\text{act,pp}}, \alpha_{\text{eff}}, \beta_{\text{pp,eff}}). \quad (16)$$

737 Here, if $\text{pp} = 1$, then $\widehat{\mathcal{T}}_{\text{comm,pp}} = 0$.

738 When combining the communication costs for all three types of parallelism, considering in practice
 739 large chunks of communication overlap with compute, we attach non-overlap factors κ and obtain

$$\widehat{\mathcal{T}}_{\text{comm}}(H; \mathbf{C}) = \kappa_{\text{dp}} \cdot \widehat{\mathcal{T}}_{\text{comm,dp}}(H; \mathbf{C}) + \kappa_{\text{tp}} \cdot \widehat{\mathcal{T}}_{\text{comm,tp}}(H; \mathbf{C}) + \kappa_{\text{pp}} \cdot \widehat{\mathcal{T}}_{\text{comm,pp}}(H; \mathbf{C}) \quad (17)$$

740 where \mathbf{C} are constants related to the various costs which are to be inferred. Given (8) and (17), we
 741 can construct the prior mean for throughput to be

$$\widehat{\mathcal{R}}(H; \{t_{\text{comp}}, \mathbf{C}\}) = [\widehat{\mathcal{T}}_{\text{comp}}(H; t_c) + \widehat{\mathcal{T}}_{\text{comm}}(H; \mathbf{C})]^{-1}. \quad (18)$$

742 D.2 Parallelism-Informed Prior Mean for the Maximum Memory Usage

743 We briefly elaborate on the choice of prior mean in (19). As discussed, we only consider the memory
 744 that is required to store the NN parameters, and those to compute the gradient updates.

745 For NN parameters, its sharding can be done on the pipeline or on the layers, allowing us to
 746 approximate the GPU memory required for storing the NN parameters to be inversely proportional to
 747 $\text{pp} \cdot \text{tp}$. Note that assuming the simplest DP implementation, the NN parameters are duplicated and
 748 stored on each DP dimension, and so the maximum memory usage is not affected by the DP.

749 Meanwhile, in the case of backpropagation computation, the maximum memory used will roughly be
 750 proportional to how many model parameters a certain GPU has to perform the forward and backward
 751 passes for, times how many training samples the GPU has to process at any one time. We expect this
 752 quantity to be inversely proportional to the number of total GPUs times the latter to depend on the
 753 number of microbatches used.

754 Combining these two factors, can write the maximum memory usage as

$$\widehat{\mathcal{M}}(H; \theta_{\mathcal{M}}) = \min \{m_1 \cdot (\text{pp} \cdot \text{tp})^{-1} + m_2 \cdot (\text{n_gpus} \cdot \text{mb})^{-1} + m_3, M_0\} \quad (19)$$

755 where, m_1 captures the memory used for storing model parameters, and m_2 captures the memory used
 756 during backpropagation computations, m_3 are any other additional memory overheads unaccounted
 757 for by our model, and $\theta_{\mathcal{M}} = \{m_1, m_2, m_3\}$. Note that since we cannot measure maximum memory
 758 usage above values of M_0 , we apply the min function to clip the GP prior beliefs.

759 D.3 Kernel for Modelling the Unknown Deviations

760 Given the embedding $e(H)$, we use the Matern kernel (29) which is given by

$$k(H, H') = \sigma_k^2 k_{\text{Matern}, \nu}(e(H), e(H'); \ell) = \sigma_k^2 \frac{2^{1-\nu}}{\Gamma(\nu)} (\sqrt{2\nu} d_\ell(H, H'))^\nu K_\nu(\sqrt{2\nu} d_\ell(H, H'))$$

761 where Γ is the Gamma function, K_ν is the modified Bessel function, σ_k is the kernel scaling constant,

$$d_\ell(H, H') = (e(H) - e(H'))^\top \mathbf{L}^{-2} (e(H) - e(H')) \quad (20)$$

762 is the distance between two PC embeddings and $\mathbf{L} = \text{diag}(\ell) = \text{diag}([\ell_1 \cdots \ell_p])$ is the lengthscale.

763 **Justification for kernel based on Euclidean distance.** Note that in our setup, we use a kernel based
 764 on the Euclidean distance despite a discrete search space. This is because we expect each value within
 765 the PCs are typically correlated according to its order (e.g., throughput when $\text{dp} = 1$ will correlate
 766 with when $\text{dp} = 2$, which then correlates with when $\text{dp} = 4$, etc.). In this regard, there exists some
 767 correlation between PCs that have similar values under our embedding (even if the values are all
 768 discrete), justifying our use of GPs.

769 We also observe this later in Table 6 which shows that the lengthscales learned by the GP are larger
 770 in some values even without explicit treatment of the discrete values, suggesting existence of learned
 771 interpolations by OPPA. Additionally, later in Fig. 20, we demonstrate that OPPA which uses a
 772 Matern kernel outperforms OPPA which instead uses a categorical based on the Hamming distance as
 773 in [4], which empirically further justifies our kernel choice.

774 D.4 GP Posterior Beliefs of Throughput and Maximum Memory Usage

775 Suppose we are in the i th round. We let $\mathbf{H}_i = [H_1 \cdots H_i]$ be the list of PCs, $\bar{\mathbf{r}}_i = [\bar{r}_{1, \hat{q}_1} \cdots \bar{r}_{i, \hat{q}_i}]$
 776 and $\sigma_{\bar{\mathbf{r}}_i}^2 = [\sigma_{\bar{r}_{1, \hat{q}_1}}^2 \cdots \sigma_{\bar{r}_{i, \hat{q}_i}}^2]$ be the measured throughput and the corresponding variance, and $\mathbf{m}_i =$
 777 $[m_1 \cdots m_i]$ be the measured maximum memory usage values.

778 Given the data, we first find the optimal hyperparameters $\theta = \{\theta_{\mathcal{R}}, \theta_{\mathcal{M}}, \theta_k\}$. This is done by
 779 maximizing the marginal log-likelihood [29], or

$$\theta = \arg \max_{\theta'} \log p(\bar{\mathbf{r}}_i, \bar{\mathbf{m}}_i | \bar{\mathbf{H}}_i, \theta). \quad (21)$$

780 The GP posterior beliefs of the throughput and maximum memory usage for any PC H' are then
 781 given by normal distributions

$$\mathcal{N}(\mu_{\mathcal{R},i}(H'), \sigma_{\mathcal{R},i}^2(H')) \quad \text{and} \quad \mathcal{N}(\mu_{\mathcal{M},i}(H'), \sigma_{\mathcal{M},i}^2(H')) \quad (22)$$

782 where the mean and variance for the throughput for any PC H' can then be defined as

$$\mu_{\mathcal{R},i}(H') = \widehat{\mathcal{R}}(\mathbf{H}_i; \theta_{\mathcal{R}}) + k(H', \mathbf{H}_i; \theta_k) (k(\mathbf{H}_i, \mathbf{H}_i; \theta_k) + \text{diag}(\sigma_{\bar{\mathbf{r}}_i}^2))^{-1} (\bar{\mathbf{r}}_i - \widehat{\mathcal{R}}(\mathbf{H}_i; \theta_{\mathcal{R}})), \quad (23)$$

$$\sigma_{\mathcal{R},i}^2(H') = k(H', H'; \theta_k) - k(H', \mathbf{H}_i; \theta_k) (k(\mathbf{H}_i, \mathbf{H}_i; \theta_k) + \text{diag}(\sigma_{\bar{\mathbf{r}}_i}^2))^{-1} k(\mathbf{H}_i, H'; \theta_k), \quad (24)$$

783 and the GP posterior mean and variance for the memory usage for any PC H' can then be defined as

$$\mu_{\mathcal{M},i}(H') = \widehat{\mathcal{M}}(\mathbf{H}_i; \theta_{\mathcal{M}}) + k(H', \mathbf{H}_i; \theta_k) (k(\mathbf{H}_i, \mathbf{H}_i; \theta_k) + \lambda I)^{-1} (\mathbf{m}_i - \widehat{\mathcal{M}}(\mathbf{H}_i; \theta_{\mathcal{M}})), \quad (25)$$

$$\sigma_{\mathcal{M},i}^2(H') = k(H', H'; \theta_k) - k(H', \mathbf{H}_i; \theta_k) (k(\mathbf{H}_i, \mathbf{H}_i; \theta_k) + \lambda I)^{-1} k(\mathbf{H}_i, H'; \theta_k) \quad (26)$$

784 where λ is to make the matrix invertible.

785 E Detailed Discussion of PC Selection Method in OPPA

786 E.1 Constrained UCB Criterion

787 The next PC $H_i \in \mathcal{H}$ to trial in round i is chosen to be the PC which maximizes the constrained
 788 upper confidence bound (cUCB) [33, 37], given by

$$\text{cUCB}_i(H) \triangleq \mathbb{E}_{\hat{r}_{H,i-1}, \hat{m}_{H,i-1}} [X_{H,i-1} + \beta_i | X_{H,i-1} - \mathbb{E}[X_{H,i-1}]] \quad (27)$$

789 where $\hat{r}_{H,i-1} \sim \mathcal{N}(\mu_{\mathcal{R},i-1}(H), \sigma_{\mathcal{R},i-1}^2(H))$ and $\hat{m}_{H,i-1} \sim \mathcal{N}(\mu_{\mathcal{M},i-1}(H), \sigma_{\mathcal{M},i-1}^2(H))$ are sam-
 790 pled from their respective GPs as modeled from ①, and $X_{H,i-1} = \hat{r}_{H,i-1} \cdot (1 - \text{sigmoid}(\hat{m}_{H,i-1}))$.
 791 Note that in the case that memory constraint is not violated (i.e., when $\text{sigmoid}(\hat{m}_{H,i-1}) \approx 0$), the
 792 objective in [27] can be reduced to the analytical UCB objective (i.e., $\text{cUCB}_i(H) \approx \mu_{\mathcal{R},i-1}(H) +$
 793 $\beta_i \sigma_{\mathcal{R},i-1}(H)$).

794 E.2 Random Sampling For Additional Exploration

795 In OPPA, we sometimes select PCs at random for additional exploration. There are two scenarios
 796 which triggers a random selection of PC in OPPA.

- 797 1. In the first few chosen PCs. This is because in the beginning there are no PCs which can be
 798 used to infer the hyperparameters for the prior distribution of the GP, therefore a few PCs
 799 are chosen at random to kick-off the modeling process and provide a reasonably diverse set
 800 of samples to infer the hyperparameters well.
- 801 2. When too many out-of-memory errors have been encountered in a row. This is because any
 802 out-of-memory trials will not result in a usable training data for the throughput modeling
 803 and possibly minimal data for the maximum memory GP which does not aid the GP model.
 804 When too many such cases are encountered, we attempt to do random exploration so that
 805 the model can receive some information that can be used to model better with and find new
 806 feasible PCs.

807 For the random selection process, we select a PC using a weighted random strategy, where we
 808 first select a parallelism dimension configuration uniformly at random, then select the remaining
 809 hyperparameters uniformly at random such that they are feasible within the selected parallelism
 810 dimension configuration.

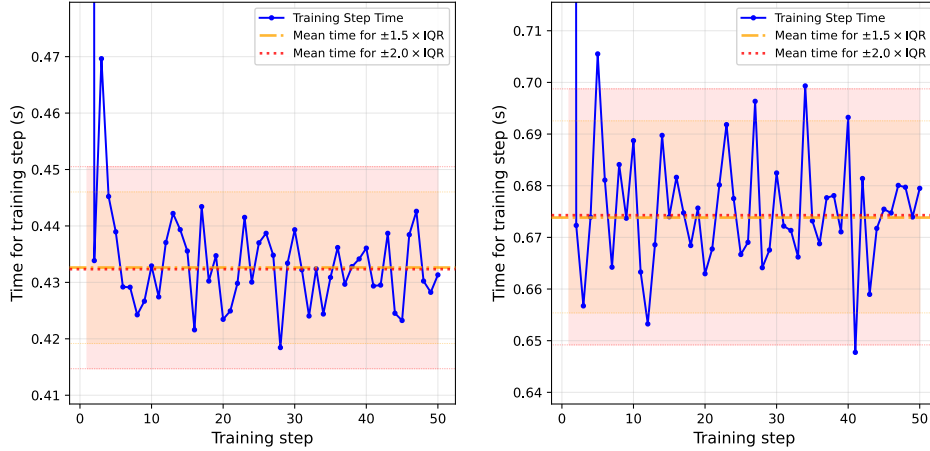


Figure 11: Example of the outlier removal process based on the IQR on different random PC trials. Here, we show the effect of the threshold selected vs. the mean of the measurements, demonstrating the robustness with respect to the different thresholds.

811 F Detailed Discussion of PC Trialing Method

812 F.1 Throughput Estimation

813 Given the time $t_{i,1}, \dots, t_{i,q}$ required for q training steps, we can estimate the throughput and its
 814 predicted variance as

$$\bar{r}_{i,q} = \frac{1}{q} \sum_{j=1}^q \frac{1}{t_{i,j}}, \quad \text{and} \quad \sigma_{\bar{r}_{i,q}}^2 = \frac{1}{q} \sum_{j=1}^q \left(\frac{1}{t_{i,j}} - \bar{r}_{i,q} \right)^2. \quad (28)$$

815 Note that $\bar{r}_{i,q}$ can be viewed as an unbiased estimator of the reciprocal of the time per training
 816 step (i.e., $\mathbb{E}[\bar{r}_{i,q}] = 1/\mathcal{T}(H_i)$ where $\mathcal{T}(H_i)$ is the average time per training step), while $\sigma_{\bar{r}_{i,q}}^2$ is the
 817 variance of the estimator.

818 F.2 Outlier Removal

819 As demonstrated in Fig. 3, not all training time measurements will be representative of the actual
 820 throughput. We therefore perform two actions. First, we remove the first training step $t_{i,1}$ since
 821 it typically corresponds to a warm-up for the training and therefore will usually be an anomaly
 822 measurement. Second, we compute the median and the inter-quartile range (IQR), and remove all
 823 measurements which are away from the median by at least $2 \times IQR$. This is a standard method
 824 for outlier removal in practice. However, we choose a looser threshold for outliers to account for
 825 fluctuating measurements. This generally has little effect on the overall optimization since most
 826 training steps tend to not fluctuate by too much anyway. This is demonstrated by Fig. 11 which shows
 827 that the predicted throughput would have little difference across the threshold of outlier removal,
 828 showing robustness of our method to the outlier removal process. The remaining training points are
 829 then used to compute (28).

830 F.3 Additional Discussion on Early Trial Termination Mechanism

831 In this section, we attempt to prove Thm. 4.1. To do so, we will consider a more general case
 832 for an arbitrary unknown function f . We do so since the problem setting of sequential repeated
 833 measurements is useful in many problem settings beyond optimizing the PC, such as in repeated
 834 scientific experiments where we can decide on-the-fly whether the same experimental setup should
 835 be repeated or not. We will first provide an intuitive justification for why early termination can be
 836 applied, followed by a theoretical proof for Thm. 4.1 and ablation studies on synthetic functions.

837 **F.3.1 Intuition Behind Early Termination Mechanism**

838 Before stating the formal proof, we provide some intuition behind why the early termination mecha-
 839 nism saves on resources while also not affecting the overall performance. In Fig. 12, we present a toy
 840 example of a BO iteration where an input is repeatedly queried.

841 Here, we notice several things. We first see that for a certain input, more repetitions of the query will
 842 result in a lower standard deviation for the measurement. This is simply from the fact that when we
 843 take the mean of several i.i.d. measurements, the mean will tend towards the actual value, and the
 844 variance of the mean predictor will decay at a rate inversely proportional to the number of repeated
 845 measurements. This is reflected in the plots where the measurement at $x = 2$ has smaller and smaller
 846 standard deviation as we have more repeats, and as we will show in Corollary F.2. In these cases, we
 847 see that some repeated queries are therefore necessary to reduce the noise of the measurements.

848 However, as we perform more repeats, there is a diminishing effect on how much more information
 849 is gained from each additional repeats. In the diagram, we see that after 20 repeated measurements
 850 (as in Fig. 12c), the GP would barely change. Similarly, we see that at 20 repeats, we are able to
 851 somewhat conclude that the global optimum is less likely to be in the neighborhood of $x = 2$, as
 852 reflected by the corresponding probability density function, and that this belief remains when we
 853 continue up to 100 repeats (as in Fig. 12d). This shows the additional repeats are unlikely to improve
 854 the predictions by the GP and the confidence of the optimum. We therefore can see that by stopping
 855 after 20 repeats, we would have been able to eventually make the same conclusion about the optimum
 856 in the end while saving on resources that would have been incurred from repeated queries.

857 Finally, we note that additional repeats are more beneficial to recovering the optimum value in the
 858 case that it improves upon the best measurement so far. In Fig. 13, we see that additional repeated
 859 measurements at $x = 2$ eventually results in diminishing returns, as seen in the probability distribution
 860 of the maximum value of $f(x)$ being similar between using 20 repeats and 100 repeats. Meanwhile, in
 861 Fig. 14 where we instead repeatedly query the input $x = 4.2$, we continue to gain information about
 862 the maximum value even when we use 100 repeats, as seen by the sharper probability distribution for
 863 the belief of the maximum $f(x)$. We therefore can see that additional repeats are less informative if
 864 they are for suboptimal inputs. We can therefore use this idea to form a metric to determine if a trial
 865 should be terminated early.

866 **F.3.2 Proof of Thm. 4.1**

867 With this intuition, we will now formally prove the performance of the early termination mechanism.
 868 For completeness, we first state the assumptions for the function and the measurements which follow
 869 from other BO works [12, 21, 33] however with additional assumptions on repeated measurements
 870 from the same input.

871 **Assumption F.1.** Let $f \sim \mathcal{GP}(0, k)$ be an unknown function whose RKHS norm $\|f\|_{\mathcal{H}} \leq B$ is
 872 bounded. The BO algorithm is as noted in Algorithm 3, where in each BO iteration i , an input $x_i \in \mathcal{X}$
 873 is selected, and $\hat{q}_i \leq q_{\max}$ noisy outputs $y_{i,j} = f(x_i) + \varepsilon_{i,j}$ are returned, where $\varepsilon_{i,j} \sim \mathcal{N}(0, s^2)$ are
 874 i.i.d. noise.

875 Note that in this following proof, we consider any function f which has a bounded RKHS norm.
 876 While this property also applies for the unknown throughput \mathcal{R} , it will also be satisfied by other
 877 black-box functions in general as well.

878 We now also show Algorithm 2 which repeats each query q_{\max} times, and our proposed Algorithm 3
 879 which does early termination on some of the rounds. Our theoretical results in this section will
 880 consider Algorithm 3.

881 We first state a result regarding the mean estimator of $f(x_i)$.

882 **Corollary F.2.** Suppose we define

$$\bar{y}_{i,q} = \frac{1}{q} \sum_{j=1}^q y_{i,j}. \quad (29)$$

883 Then, the expected value of \bar{x}_i is $\mathbb{E}[\bar{y}_{i,q}] = f(x_i)$, and its variance bounded by $\mathbb{V}[\bar{y}_{i,q}] = s^2/q$.

884 *Proof.* The results are direct consequences of the summation of expected values and variances of
 885 independent random variables, and the fact that $\mathbb{V}[y_{i,j}] = s^2$ by assumption. \square

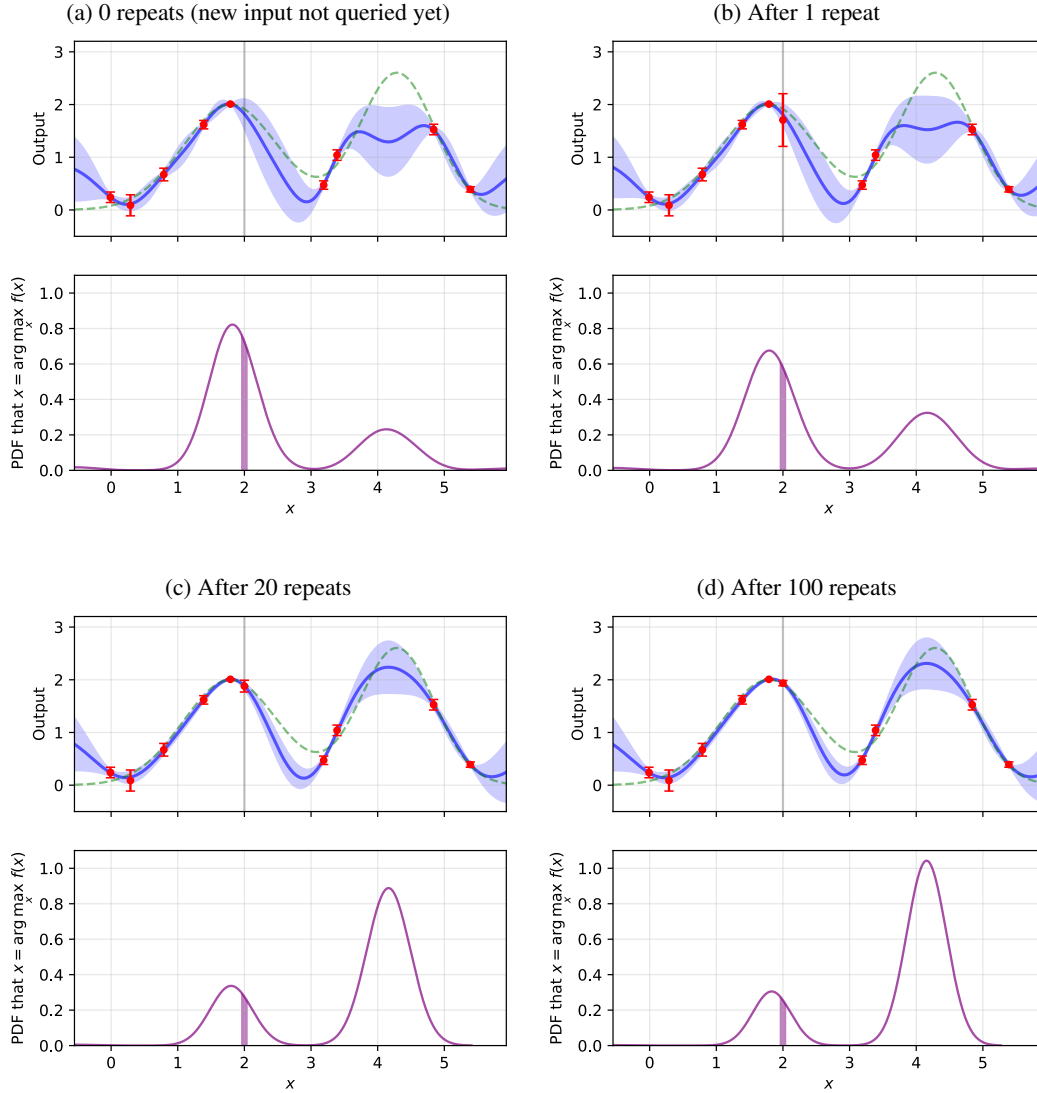


Figure 12: Demonstration of why early termination of certain trials does not affect final result. Figs. 12a to 12d represent the modeled GP and the predicted optimum when a certain input, in this case at $x = \arg \max_{x'} \mu_t(x') + \beta_t \sigma_t(x') = 2$, is repeatedly queried for 0, 1, 20, and 100 times. For each subfigure, the top diagram represents the modeled GP where the dashed green line represents the actual function, the blue line and error band represent the GP mean and standard deviation, and the red points represent the predicted mean and standard deviation (the latter which shrinks with more repeats). In the bottom diagram, we show the probability of each input being the optimal (i.e., probability that $x = \arg \max_{x'} f(x')$), while highlighting the interval around the input $x = 2$.

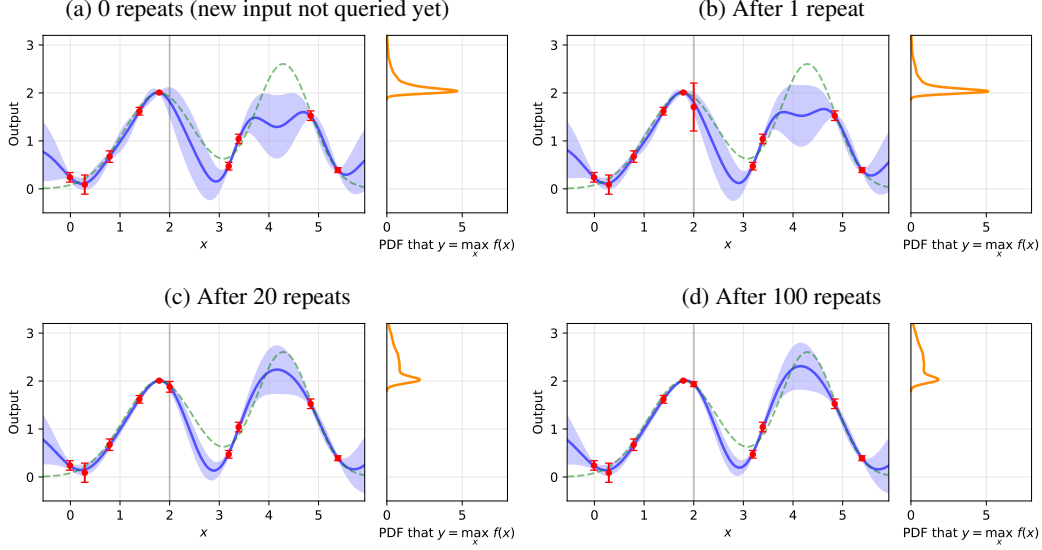


Figure 13: Additional demonstration of why early termination of certain trials does not affect the obtained optimal value. The setup is as done previously in Fig. 12 with queries at $x = 2$. For each subfigure, the left plot represents the modeled GP as shown previously in Fig. 12. In the right diagram, we show the probability density function for the belief of the maximum value $\max_x f(x)$ according to the GP.

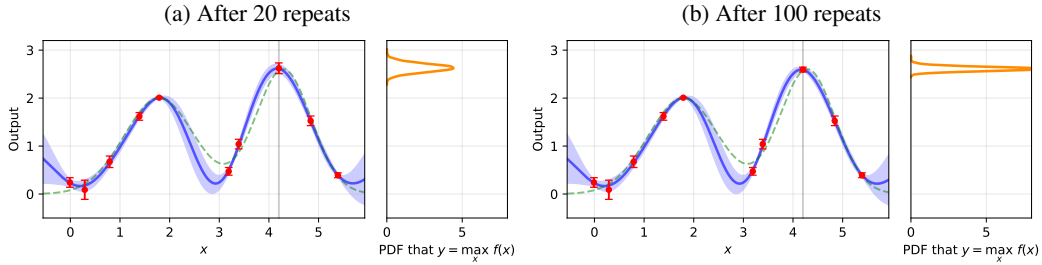


Figure 14: Additional demonstration of why early termination of certain trials does not affect the obtained optimal value. The setup is similar to as done previously in Fig. 12, however instead with queries at $x = 4.2$, and with only 20 repeats and 100 repeats shown (comparable to Figs. 13c and 13d). The information in each subfigure is the same as in Fig. 13.

886 We will now prove the first half of Thm. 4.1. Given the results in Corollary F.2, we next show that we
 887 are able to obtain a GP with good predictive bounds even if some trials are terminated early, i.e., not
 888 ran for up to q_{\max} repeats.

889 **Lemma F.3.** *Suppose we let*

$$\mu_i(x) = k(x, \mathbf{X}_i)(k(\mathbf{X}_i, \mathbf{X}_i) + s^2\mathbf{Q}_i^{-1})^{-1}\mathbf{y}_i, \quad (30)$$

$$\sigma_i^2(x) = k(x, x) - k(x, \mathbf{X}_i)(k(\mathbf{X}_i, \mathbf{X}_i) + s^2\mathbf{Q}_i^{-1})^{-1}k(\mathbf{X}_i, x) \quad (31)$$

890 where $\mathbf{X}_i = [x_1 \cdots x_i]$, $\mathbf{y}_i = [\bar{y}_{1, \hat{q}_1} \cdots \bar{y}_{i, \hat{q}_i}]$, and $\mathbf{Q}_i = \text{diag}([\hat{q}_1 \cdots \hat{q}_i])$. If

$$\beta_i = B + \sqrt{2 \log \frac{\det(k(\mathbf{X}_i, \mathbf{X}_i) + s^2\mathbf{Q}_i^{-1})^{1/2}}{\delta \det(s^2\mathbf{Q}_i^{-1})^{1/2}}} \quad (32)$$

891 then, with probability greater than $1 - \delta$, for all $x \in \mathcal{X}$ and all $i = 1, \dots, N$, we have

$$|f(x) - \mu_{i-1}(x)| \leq \beta_i \sigma_{i-1}(x). \quad (33)$$

Algorithm 2 GP-UCB with Repeated Trials

```
1:  $\mathcal{D}_0 \leftarrow \emptyset$ 
2: for  $i = 1, \dots, N$  do
3:   Fit GP on  $\mathcal{D}_{i-1}$  to obtain mean  $\mu_{i-1}$  and variance  $\sigma_{i-1}^2$ 
4:    $x_i \leftarrow \arg \max_{x \in \mathcal{X}} \mu_{i-1}(x) + \beta_i \sigma_{i-1}(x)$ 
5:   for  $j = 1, \dots, q_{\max}$  do
6:     Sample  $y_{i,j} = f(x_i) + \varepsilon_{i,j}$ 
7:      $\bar{y}_{i,q_{\max}} \leftarrow q_{\max}^{-1} \sum_{j=1}^{q_{\max}} y_{i,j}$ 
8:    $\mathcal{D}_i \leftarrow \mathcal{D}_i \cup \{(x_i, \bar{y}_{i,q_{\max}})\}$ 
```

Algorithm 3 Modified GP-UCB with Early Trial Termination

```
1:  $\mathcal{D}_0 \leftarrow \emptyset$ 
2: for  $i = 1, \dots, N$  do
3:   Fit GP on  $\mathcal{D}_{i-1}$  to obtain mean  $\mu_{i-1}$  and variance  $\sigma_{i-1}^2$ 
4:    $x_i \leftarrow \arg \max_{x \in \mathcal{X}} \mu_{i-1}(x) + \beta_i \sigma_{i-1}(x)$ 
5:   for  $j = 1, \dots, q_{\max}$  do
6:     Sample  $y_{i,j} = f(x_i) + \varepsilon_{i,j}$ 
7:      $\bar{y}_{i,j} \leftarrow j^{-1} \sum_{j'=1}^j y_{i,j'}$ 
8:      $\mathcal{D}'_{i,j'} \leftarrow \mathcal{D}_i \cup \{(x_i, \bar{y}_{i,j'})\}$ 
9:     if  $j > q_{\min}$  and  $\bar{y}_{i,j} < \max_{j' < i} \bar{y}_{j',q_j} + \tau_q$  then
10:      break // Early termination mechanism
11:    $\mathcal{D}_i \leftarrow \mathcal{D}'_{i,j'}$ 
```

892 *Proof.* Given the variance of the mean predictor from Corollary [F.2](#), our scenario can be thought of
893 as having i function measurements with heteroscedastic noise with variances of $s^2/\hat{q}_1, \dots, s^2/\hat{q}_i$.
894 The variance bounds then follow directly from Lemma 7 in [\[12\]](#) where we substitute $\Sigma_i \rightarrow s^2 \mathbf{Q}_i^{-1}$
895 and $\lambda \rightarrow 1$. \square

896 **Corollary F.4.** Given Lemma [F.3](#), for all $x' \in \mathcal{X}$, we have $f(x') - \mu_{i-1}(x_i) \leq \beta_i \sigma_{i-1}(x_i)$.

897 *Proof.* We see that

$$f(x') - \mu_{i-1}(x_i) \leq \mu_{i-1}(x') + \beta_i \sigma_{i-1}(x') - \mu_{i-1}(x_i) \quad \text{by Lemma [F.3](#)} \quad (34)$$

$$\leq \mu_{i-1}(x_i) + \beta_i \sigma_{i-1}(x_i) - \mu_{i-1}(x_i) \quad \text{by how } x_i \text{ is chosen,} \quad (35)$$

$$= \beta_i \sigma_{i-1}(x_i). \quad (36)$$

898 \square

899 We now show the cumulative regret of the problem. Let

$$\mathbb{I}[\mathbf{y}_X; \mathbf{f}_X] = \frac{1}{2} \sum_{x'_i \in X} \log \left(1 + \frac{\sigma_{i-1}^2(x'_i)}{s^2/\hat{q}_i} \right), \quad (37)$$

900 and

$$\gamma_i = \max_{X=[x'_1, \dots, x'_i] \subset \mathcal{X}} \mathbb{I}[\mathbf{y}_X; \mathbf{f}_X] \quad (38)$$

901 be the maximum possible information gain across i rounds. We then prove the following result.

902 **Theorem F.5.** Let $x^* = \arg \max_{x \in \mathcal{X}} f(x)$. With probability at least $1 - \delta$,

$$\sum_{i=1}^N f(x^*) - f(x_i) \leq s \beta_N \sqrt{\frac{8N \gamma_N}{q_{\min}}}. \quad (39)$$

903 *Proof.* This result is similar to previous results for UCB-based methods, e.g., Theorem 3 in [\[33\]](#) or
904 Corollary 9 in [\[12\]](#).

905 With probability at least $1 - \delta$, Lemma F.3 holds. We see that

$$\sum_{i=1}^N (f(x^*) - f(x_i)) = \sum_{i=1}^N \underbrace{(f(x^*) - \mu_{i-1}(x_i))}_{\text{Corollary F.4}} + \underbrace{(\mu_{i-1}(x_i) - f(x_i))}_{\text{Lemma F.3}} \quad (40)$$

$$\leq \sum_{i=1}^N 2\beta_i \sigma_{i-1}(x_i). \quad (41)$$

906 Since

$$\sum_{i=1}^N \sigma_{i-1}^2(x_i) \leq \sum_{i=1}^N \frac{s^2 \sigma_{i-1}^2(x_i)}{\hat{q}_i} \quad (42)$$

$$\leq \frac{s^2}{q_{\min}} \sum_{i=1}^N \log \left(1 + \frac{\sigma_{i-1}^2(x_i)}{s^2/\hat{q}_i} \right) \quad (43)$$

$$\leq \frac{2s^2}{q_{\min}} \gamma_N, \quad (44)$$

907 we can rewrite (41) as

$$\sum_{i=1}^N (f(x^*) - f(x_i)) \leq \sqrt{\sum_{i=1}^N 4\beta_i^2} \sqrt{\sum_{i=1}^N \sigma_{i-1}^2(x_i)} \quad \text{by Cauchy-Schwartz,} \quad (45)$$

$$\leq \beta_N \sqrt{4N} \sqrt{\sum_{i=1}^N \sigma_{i-1}^2(x_i)} \quad \text{since } \beta_i \leq \beta_N, \quad (46)$$

$$\leq s\beta_N \sqrt{\frac{8N\gamma_N}{q_{\min}}} \quad \text{by (44).} \quad (47)$$

908 \square

909 *Remark F.6* (The regret bound in Thm. F.5 is lower when q_{\min} increases). We investigate the upper
910 bound in (47) further to see its dependence on q_{\min} .

911 For simplicity, we will let $\mathbf{K} = k(\mathbf{X}_N, \mathbf{X}_N)$. First, we see that

$$\log \det (\mathbf{K} + s^2 \mathbf{Q}_N^{-1}) = \log \det s^2 \mathbf{Q}_N^{-1} + \log \det (I + s^{-2} \mathbf{Q}_N \mathbf{K}) \quad (48)$$

$$\leq \log \det s^2 \mathbf{Q}_N^{-1} + \log \det (I + s^{-2} q_{\max} \mathbf{K}), \quad (49)$$

912 which means that

$$\beta_N = B + \sqrt{2 \log \frac{\det (\mathbf{K} + s^2 \mathbf{Q}_N^{-1})^{1/2}}{\delta \det (s^2 \mathbf{Q}_N^{-1})^{1/2}}} \quad (50)$$

$$= B + \sqrt{2 \log \frac{1}{\delta} + \log \det (\mathbf{K} + s^2 \mathbf{Q}_N^{-1}) - \log \det s^2 \mathbf{Q}_N^{-1}} \quad (51)$$

$$\leq B + \sqrt{2 \log \frac{1}{\delta} + \log \det (I + s^{-2} q_{\max} \mathbf{K})}. \quad (52)$$

913 Furthermore, we see that

$$\gamma_N = \max_{X=[x'_1, \dots, x'_N] \subset \mathcal{X}} \frac{1}{2} \sum_{x'_i \in X} \log \left(1 + \frac{\sigma_{i-1}^2(x'_i)}{s^2/\hat{q}_i} \right) \quad (53)$$

$$\leq \max_{X=[x'_1, \dots, x'_N] \subset \mathcal{X}} \frac{1}{2} \sum_{x'_i \in X} \log \left(1 + \frac{\sigma_{i-1}^2(x'_i)}{s^2/q_{\max}} \right). \quad (54)$$

914 Therefore we see that both β_N and γ_N are upper bounded by terms which are independent of q_{\min} ,
915 and so the constants in the upper bound provided in Thm. F.5 do not hide any additional dependencies
916 with respect to q_{\min} . This shows that the upper bound of cumulative regret from (47) decays at a rate
917 of $1/\sqrt{q_{\min}}$.

918 We will now prove the second half of Thm. 4.1. We first prove the following result.

919 **Lemma F.7.** Define $c_1 \triangleq \sqrt{2 \log(2Nq_{\max}/\delta)}$. With probability at least $1 - \delta$, for all $i = 1, \dots, N$
 920 and all $q = 1, \dots, q_{\max}$, we have

$$|\bar{y}_{i,q} - f(x_i)| \leq \frac{c_1 s}{\sqrt{q}}. \quad (55)$$

921 *Proof.* From Corollary F.2, we know that the $\bar{y}_{i,q}$ is normally distributed with mean $f(x_i)$ and
 922 standard deviation s/\sqrt{q} . By Chernoff bounds, for each $i = 1, \dots, N$ and each $q = 1, \dots, q_{\max}$, we
 923 would have $\bar{y}_{i,q} - f(x_i) > c_1 s/\sqrt{q}$, and $f(x_i) - \bar{y}_{i,q} > c_1 s/\sqrt{q}$ where either event happens with
 924 probability no greater than $\delta N/2q_{\max}$. This means that with probability no greater than δ/Nq_{\max} ,
 925 we have $|\bar{y}_{i,q} - f(x_i)| > c_1 s/\sqrt{q}$. Therefore, by union bound, we have for all $i = 1, \dots, N$
 926 and each $q = 1, \dots, q_{\max}$, we would have $|\bar{y}_{i,q} - f(x_i)| \leq c_1 s/\sqrt{q}$ with probability greater than
 927 $1 - (\delta/Nq_{\max})(Nq_{\max}) = 1 - \delta$. \square

928 **Theorem F.8.** Define

$$\tau_q = c_1 s \left(\frac{1}{\sqrt{q_{\min}}} + \frac{1}{\sqrt{q}} \right). \quad (56)$$

929 Then, with probability at least $1 - \delta$, for all $i \leq N$, if we have $f(x_i) < \max_{j < i} f(x_j)$, then $\hat{q}_i < q_{\max}$.

930 *Proof.* With probability at least $1 - \delta$, Lemma F.7 applies.

931 To prove the statement above, we show its contrapositive. Suppose we have $\hat{q}_i = q_{\max}$, or that the trial
 932 for x_i does not terminate early. This implies that $\bar{y}_{i,q} \geq \max_{j < i} \bar{y}_{j,\hat{q}_j} + \tau_q$ for all $q = q_{\min}, \dots, q_{\max}$.
 933 For any q in this range, we would then have

$$f(x_i) \geq \bar{y}_{i,q} - \frac{c_1 s}{\sqrt{q}} \quad (57)$$

$$\geq \max_{j < i} \bar{y}_{j,\hat{q}_j} + \tau_q - \frac{c_1 s}{\sqrt{q}} \quad (58)$$

$$\geq \max_{j < i} f(x_j) - \frac{c_1 s}{\sqrt{q_{\min}}} + \tau_q - \frac{c_1 s}{\sqrt{q}} \quad (59)$$

$$= \max_{j < i} f(x_j). \quad (60)$$

934 This proves the contrapositive which in turn proves the original statement. \square

935 Finally, Thms. F.5 and F.8 can be combined with appropriate union bounds to achieve a more formal
 936 version of Thm. 4.1.

937 F.3.3 Ablation Studies on Toy Examples

938 In Fig. 15, we provide a brief empirical demonstration of efficiency gains due to early termination.
 939 We see that when the noise variance is too high, querying the function once per input would give
 940 measurements which are too noisy to give good information. Meanwhile, by repeating each query a
 941 maximum number of times, we can obtain a good estimate of the actual function and allow the BO
 942 algorithm to arrive at the optimal using few queries. However, we see that when early termination
 943 is allowed, we can still arrive at the optimal input as before, while not requiring all queries to be
 944 repeated the maximum number of times. This shows that early termination allows for efficiency gains
 945 while minimally sacrificing on the actual optimization process.

946 In practice, since it is difficult to determine γ_N , β_N , s , and c_1 exactly, we instead fix β_i and τ_q
 947 to some constant. In OPPA, we choose $\beta_i = 1$ and $\tau_q = 10^{-3}$. We find that these values work well for
 948 our methods. Furthermore, in our proofs we do not consider the constrained BO setting. Despite
 949 this, the early termination can still be used in practice to achieve good results which we show in the
 950 experiments.

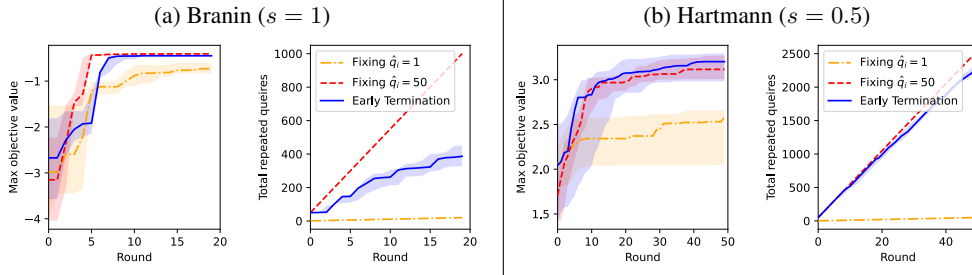


Figure 15: Results of BO with early termination (with $q_{\min} = 1$) on different synthetic functions where each query has different noise levels. For each case, we show the best queried objective value (left plot), and the cumulative number of total repeated queries made with each BO iteration (right plot).

951 G Pseudocode of OPPA

952 We present the pseudocode for OPPA in Algorithm 4.

953 In practice, we find that fixed hyperparameter values (i.e., setting $\beta_t = 1$, $\tau_q = 10^{-3}$, and $q_{\min} = 5$
 954 in all iterations as done in our main experiments) gives consistent results since the parameters are
 955 generally robust to different choices of hyperparameter values anyway.

956 In the experiments, the value of q_{\max} are set in accordance to Table 2 depending on the benchmark;
 957 however in practice can be set as large as the practitioner wants depending on how accurately they
 958 would like to measure the actual throughput of the PC. As plotted in Fig. 16, we find that after
 959 a certain point, a larger number of repeats will typically not have large effects on the measured
 960 throughput of a PC, showing that OPPA is quite robust to different choices of q_{\max} . For OPPA, larger
 961 q_{\max} is more tolerable since the early termination mechanism will not exhaust such budget anyway.

962 We also consider removing outliers which are beyond two times the IQR from the mean; this is be
 963 more lenient towards fluctuating training effects, however can also be set to 1.5 times the IQR as
 964 typically done as well with little consequences, as described in App. F.2.

Algorithm 4 OPTIMIZER FOR PARALLELISM CONFIGURATIONS (OPPA)

- 1: Generate all valid PCs \mathcal{H}
 - 2: $i \leftarrow 1$
 - 3: **while** time budget not exhausted **do**
 - 4: **if** $i < N_{\text{random}}$ **then**
 - 5: Select H_i randomly
 - 6: **else**
 - 7: // Step ① – Modeling throughput and memory usage
 - 8: Construct $\mu_{\mathcal{R},i-1}$ and $\sigma_{\mathcal{R},i-1}^2$ according to (23) and (24) respectively
 - 9: Construct $\mu_{\mathcal{M},i-1}$ and $\sigma_{\mathcal{M},i-1}^2$ according to (25) and (26) respectively
 - 10: // Step ② – Selecting the next PC to trial
 - 11: $H_i \leftarrow \arg \max_{H \in \mathcal{H} \setminus \{H_1, \dots, H_{i-1}\}} \text{cUCB}_i(H)$ where cUCB_{i-1} is defined in (27)
 - 12: // Step ③ – Trialing selected PC
 - 13: **for** $q = 1, \dots, q_{\max}$ **do**
 - 14: Measure training step time as $t_{i,j}$
 - 15: Compute $\bar{r}_{i,q}$ and $\sigma_{\bar{r}_{i,q}}^2$ according to (28)
 - 16: **if** $I_{i,q} = 0$ **then**
 - 17: **break** // Early termination condition from (3)
 - 18: Measure maximum memory usage as m_i
 - 19: $\hat{q}_i \leftarrow q$ to track the number of training steps ran in round i
 - 20: $i \leftarrow i + 1$
 - 21: **return** H_{i^*} where $i^* = \arg \max_i : (m_i < M_0) \wedge (\hat{q}_i = q_{\max}) \bar{r}_{i, \hat{q}_i}$
-

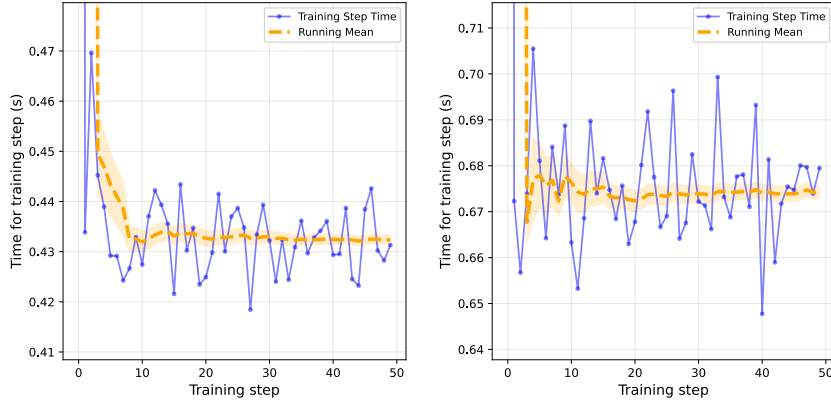


Figure 16: The effect of number of training steps repeatedly measured on the predicted throughput. The blue lines represent a training trial and the measured time per training step. The orange line represents the average training time (i.e., reciprocal of throughput) predicted up to that certain training step.

965 H Additional Information on Experimental Setup

966 H.1 Training and Hardware Configurations

967 We list the models used in our experiments in Table 2, along with the allotted search time and how
 968 many trials we repeat on them. All models used are based on the transformer architecture, and were
 969 retrieved from Huggingface.

970 Note that in all of our plots, we plot the median value (with a line) and also the lower and upper
 971 quartiles (with a fainter band over and under the line). We do so since we find that the values are
 972 often asymmetrically skewed, and therefore opted to show the quartile values to more accurately
 973 represent the distribution of these values. Also note that the repeated trials were reduced for larger
 974 models due to restrictions in compute budget.

Table 2: Details of models used in our experiments and the corresponding training scenario

Case	Model	# Params	Batch size	Max. seq. length	Search Time	q_{\max}	Repeats
BERT	BERT Base Uncased	110M	256	256	20 mins	50	10
Qwen2-1.5b	Qwen-2	1.5B	64	1024	20 mins	30	10
Llama3-1b	Llama 3	1B	64	1024	60 mins	20	5
Llama2-7b	Llama 2	7B	256	1024	60 mins	30	5
Qwen3-32b	Qwen-3	32B	256	8192	80 mins	30	5

975 In Table 3, we list the hardware configuration used in our experiments. The hardware configurations
 976 used are based on the resources that are available to the authors.

Table 3: Configurations of tested hardware.

Config. Name	GPU Model (Memory)	GPUs per Node	# Node	Multi-Node Characteristic
8 GPU _s	NVIDIA RTX A5000 (24GB)	8	1	-
16 GPU _s	NVIDIA RTX 3080 (10GB)	8	2	Docker Overlay Network
32 GPU _s	NVIDIA A100 (40GB)	4	8	High-Performance Compute
128 GPU _s	NVIDIA GH200 (96GB)	4	32	High-Performance Compute

977 H.2 Other Methods Ran For Comparison

978 We list the other methods we have ran along with their implementation details here.

- 979 • RANDOM. This involves randomly selecting a PC from \mathcal{H} to trial in each round until the
 980 time budget is exhausted.

- 981 • XGBOOST [3]. This is the method DEEPSPEED [28] uses to configure the PC, and is
982 adapted to work with the hyperparameters in our PC. The method involves training an
983 XGBOOST model based on the measured throughput, then selecting the next PC to trial
984 as the one whose predicted throughput is the highest. This is repeated until time budget is
985 exhausted.
- 986 • COST-MODEL. This involves using the cost model as described in Apps. D.1 and D.2,
987 learned based on several randomly selected PCs, to obtain a prediction of throughput and
988 maximum memory usage, and perform a one-shot selection of the best PC according to the
989 predictions.
- 990 • VANILLA-BO. This performs BO whose GP has a constant mean and a Matern kernel with
991 $\nu = 5/2$. The cUCB criterion is used for PC selection. The BO loop is implemented using
992 BOTORCH [1].
- 993 • OPPA. This is the method proposed in Sec. 4 which involves modifying BO to include
994 parallelism-informed GP prior beliefs and early trial termination. The hyperparameters used
995 are as stated in App. G.

996 We note that due to the search space employed for \mathcal{H} , we do not consider methods which performs
997 non-adaptive optimization with a cost model. This is because those methods optimize with respect to
998 the computation graph rather than the hyperparameters which we discussed in App. C.1 and since
999 they do not use the same information as OPPA to perform optimization, making it futile to compare
1000 between the two since they focus on optimizing different aspects of training parallelism. Nonetheless,
1001 we provide some comparisons with selected methods of such nature, namely, [15] and [19] in Fig. 9.

1002 I Additional Results

1003 I.1 Plots of Best Achieved Throughput vs. Other Quantities

1004 In Fig. 17, we plot the best achieved throughput vs. the number of training trials that have been ran.
1005 We see that in this view, OPPA still outperforms other methods. While the margin may be smaller in
1006 some instances, we see that OPPA is able to achieve the good results more consistently as seen by
1007 the error bars when compared to some of the other methods. This also demonstrates that the prior
1008 belief used in OPPA alone would have helped in achieving a better performance regardless of the
1009 early termination mechanism in OPPA.

1010 In Fig. 18, we show the achieved throughput is plotted against the number of training steps run in the
1011 training trials, showing that the difference in efficiency of OPPA becomes more pronounced. From
1012 these results see that OPPA is both more time-efficient and query-efficient which can be useful when
1013 the overhead to run one trial may become higher, for example when the framework is adapted to run
1014 on a cluster with a job scheduler.

1015 I.2 Optimal PCs recovered by OPPA

1016 In Tables 4 and 5, we show the PCs that were recovered by OPPA for training scenarios ran on the
1017 COLOSSAL-AI and NEMO frameworks respectively. Note that the optimal PCs chosen match quite
1018 well with intuition where for smaller models DP tends to be prioritized. Meanwhile for larger models
1019 and training scenarios which are done across multiple nodes, PP and potentially TP is prioritized. In
1020 all of these cases, there is a range for the hyperparameters within the potential PC, which suggests
1021 that multiple choices of PC can achieve similar throughput.

1022 I.3 Flexibility Across Different Parallel Training Frameworks

1023 Another advantage of adaptively selecting actual training trials to run can be seen when we change the
1024 parallel training framework used to train the NN. In Fig. 19, we run similar experiments as in Fig. 4,
1025 however instead performing training on NEMO [13]. In these experiments, we consider additional
1026 hyperparameters within our PC, such as the size of the sequence parallelism dimension. We see that
1027 compared to Fig. 4, the obtained maxima are different because the training processes are implemented
1028 differently across the two frameworks, and the fact that the adjustable hyperparameters in the PC are
1029 also different. Regardless, we see that even though most methods can eventually recover the optimal

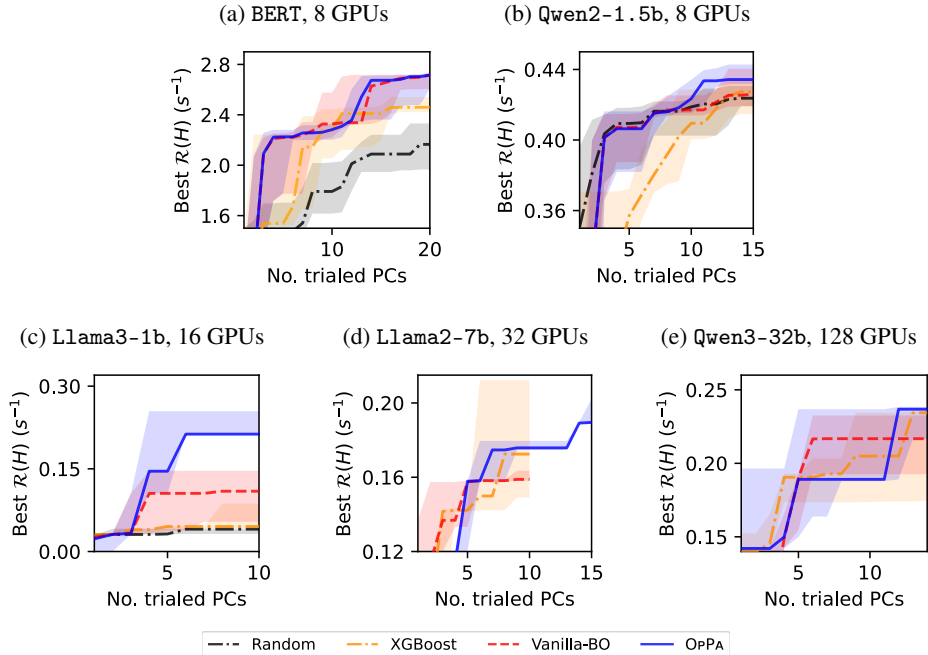


Figure 17: Results of the best obtained throughput (higher is better) plotted against the number of PCs that have been trialed. The lines represent the median value of the best obtained throughput across five trials, while the error bar represent the quartile values. Note that for the experiments on 32 GPUs, we are unable to run the optimization beyond the allotted time due to resource constraints and therefore are only able to plot some methods for a fewer number of queried PCs.

Table 4: Example of optimal PCs selected by OPBA in different training scenarios when training on COLOSSAL-AI framework. Values are based on the found optimal PCs across multiple repeated trials. A range of value shows that a certain parameter shows a spread across multiple trials (i.e., no strong preference towards one value), while a dash shows that the parallelism dimension associated with that hyperparameter is not used. Note that these may not be the best PCs for the training scenarios, but are instead examples of what are recovered by OPBA.

Hyperparameter	BERT, 8 GPUs (Fig. 4a)	Qwen2-1.5b, 8 GPUs (Fig. 4b)	ViT, 8 GPUs (Fig. 21)	Llama3-1b, 16 GPUs (Fig. 8a)	Llama2-7b, 32 GPUs (Fig. 8b)
DP size (dp)	8	4	8	1	2
TP size (tp)	1	1	1	1	1
PP size (pp)	1	2	1	16	16
DP bucket size	1 - 64	1 - 4096	1 - 64	-	1 - 4096
ZeRO stage	0	1	0	-	1
ZeRO bucket size	-	1 - 64	-	-	64 - 4096
Overlap ZeRO communication	-	True/False	-	-	False
Overlap DP AllGather	-	True/False	-	-	False
# microbatches (mb)	-	8	-	64	32
# model chunks (mc)	-	1 - 2	-	1	1
Overlap P2P for PP	-	True/False	-	False	True/False
Grad. checkpointing	False	False	True	False	False

1030 PC, OPBA is able to do so much more rapidly and, more importantly, much more consistently. The
 1031 latter fact can be seen in the error bars of the other methods which are much larger than that of OPBA.
 1032 This demonstrates that OPBA does not overfit to any one specific training framework, but instead can
 1033 simultaneously cater to different training frameworks.

1034 I.4 PC Optimization for Other Models

1035 In Fig. 21, we presented the results for tuning the PC for ViT model [6]. We see that the results here
 1036 show that OPBA is able to select better PCs compared to the other methods, consistent with other
 1037 results in the paper.

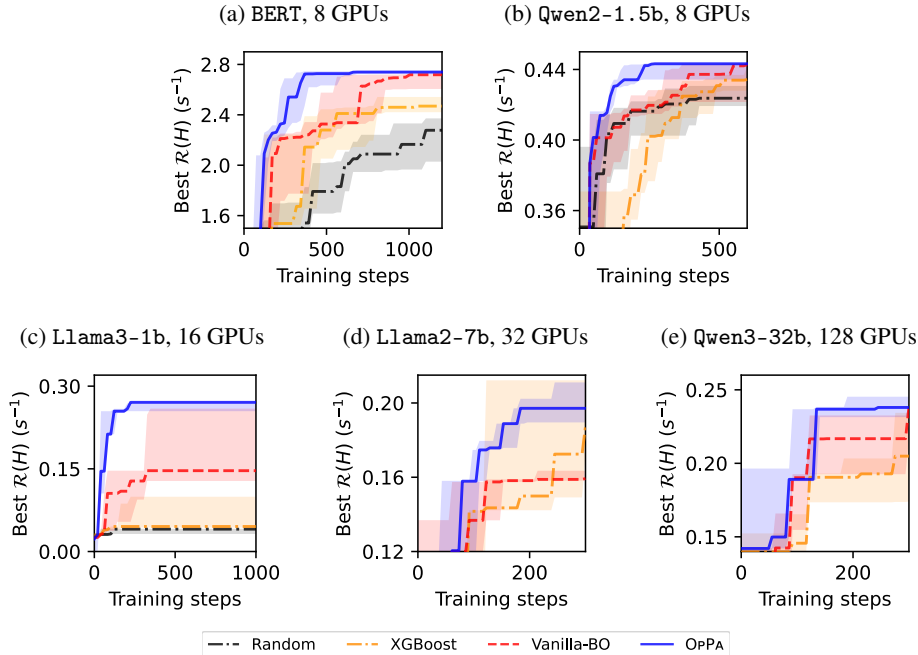


Figure 18: Results of the best obtained throughput (higher is better) plotted against the total number of training steps ran during in the training trials. The lines represent the median value of the best obtained throughput across five trials, while the error bar represent the quartile values.

Table 5: Example of optimal PCs selected by OPPA in different training scenarios when training on NEMO framework. The interpretation is as in Table 4.

Hyperparameter	BERT, 8 GPUs (Fig. 19a)	Qwen2-1.5b, 8 GPUs (Fig. 19b)	MoE, 8 GPUs (Fig. 22)	Qwen3-32b, 128 GPUs (Fig. 8c)
DP size (dp)	8	4	8	16-32
TP size (tp)	1	1	1	4
PP size (pp)	1	2	1	1-2
SP size	1	1	1	1-4
EP size	-	-	2	-
CP size	-	-	-	1
DP bucket size	1-64	1-4096	1-4096	1-4096
ZeRO stage	0	0	0	0
Overlap ZeRO communication	-	-	-	-
Overlap DP AllGather	False	False	True	True
# model chunks (mc)	-	1	-	-/4
Overlap P2P for PP	-	-/True/False	-	True/False
Grad. checkpointing	False	False	False	False

1038 Additionally, in Fig. 22, we presented the results for tuning the PC for transformer-based Mixture-
 1039 of-Experts (MoE) model [30] with the additional tuning of the expert parallelism (EP) degree being
 1040 performed by the same framework in addition to the existing DP, TP and PP. For the experiments we
 1041 consider parallel training of a model with 16 transformer units, 32 attention heads and 32 experts.
 1042 Each trial is ran for up to $q_{\max} = 50$ steps, and in OPBA, we set $q_{\min} = 10$ to account for higher
 1043 variation in training step times. In the results, we again see consistently better performances compared
 1044 to the other competing methods. Specifically, even if OPBA may eventually lead to similar results to
 1045 vanilla BO at the end of the optimization process, OPBA is able to do so more consistently (lower
 1046 variance across random runs), and is able to arrive at the optimal PC much more rapidly and efficiently.
 1047 These results further demonstrate that OPBA is able to also generalize to other models as well.

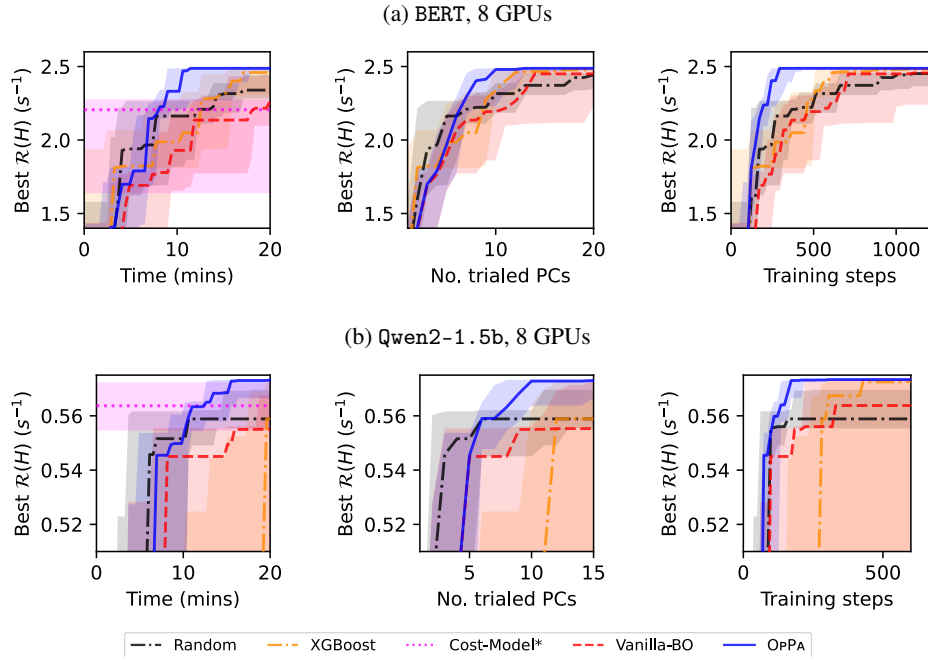


Figure 19: The best obtained throughput (higher is better) vs. the duration each method has been run for, when training on a system with a single node and using the NEMO framework [13]. Both cases train NN with identical architectures as the experiments in Figs. 4a and 4b respectively, and are ran on a single node.

1048 I.5 Predicted Throughput and Memory by Parallelism-Informed GP Prior Beliefs

1049 In Figs. 23 to 25, we compare the modeled throughput with the actual values in different training
 1050 scenarios. We see that in this case, the parallelism-informed GP prior beliefs allow for the values to be
 1051 modeled adequately well, but more importantly, allow for the PC which achieves the best throughput
 1052 to also have the highest values, and therefore be identified correctly. We find that for BO, surrogate
 1053 models only need to model the good PCs well in order to select a good PC in the end. Meanwhile,
 1054 when not using the parallelism-informed prior beliefs, the GP posterior beliefs learn the patterns
 1055 much less efficiently or do not learn them at all. This correlates well with the results in the main text
 1056 where standard BO selects a worse PC compared to OPPA which uses a better GP prior belief.

1057 In Fig. 26, we compare the modeled maximum memory with the actual measured value. In both
 1058 cases, a GP has been used, however, with and without a parallelism-informed GP prior belief since
 1059 only VANILLA-BO and OPPA are the only methods we tested which explicitly models the memory
 1060 usage. Here, we see that OPPA is able to better model the memory usage due to its use of the GP
 1061 prior belief. This is reflected in the confusion matrices which shows that after training, OPPA is able
 1062 to more accurately detect when a certain PC will result in out-of-memory errors.

1063 To additionally demonstrate the interpretability of the GP for the surrogate model, in Table 6, we show
 1064 the results for the lengthscales learned by the kernel (as given in (20)). We see that hyperparameters
 1065 that have larger effects on the resulting throughput or are less well-modeled by our GP prior belief
 1066 will typically correspond to the shorter lengthscales. For example, for the training of BERT, we
 1067 see that the parameters for TP dimension size and for the number of microbatches (for PP) have
 1068 shorter lengthscales. This matches our intuition where the throughput would be more sensitive to the
 1069 increased TP or PP being used (for the worse).

1070 I.6 Effects of Prior Misspecification

1071 To investigate the robustness of our method with respect to a misspecified prior, we perform experi-
 1072 ments to see how OPPA behaves as our cost-model prior becomes increasingly inaccurate. To do so,
 1073 we add a perturbation term into our cost function where we increase the magnitude of the perturbation

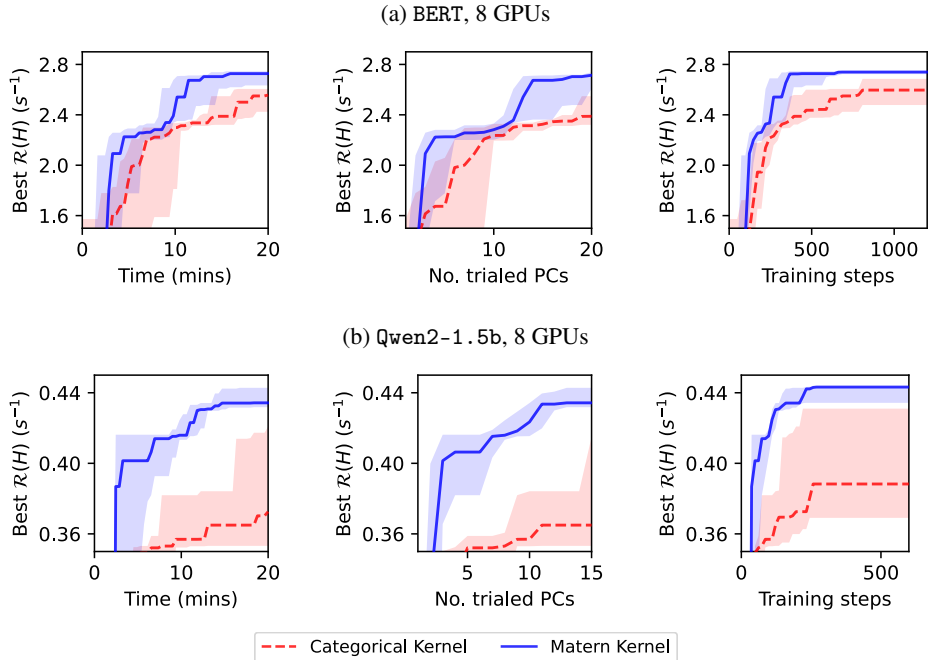


Figure 20: The best obtained throughput (higher is better) vs. the duration OP-PA has been run for, when changing from the Matern kernel to the categorical kernel [4]. Both cases train NN with identical architectures as the experiments in Figs. 4a and 4b respectively (with the results for the Matern kernel being exactly the same as in those figures), and are ran on a single node.

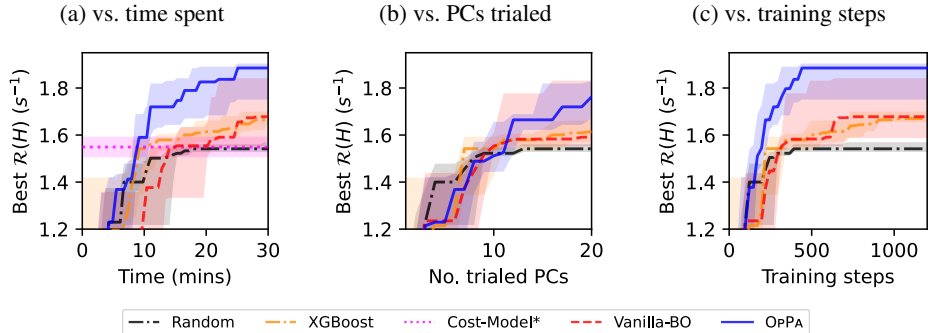


Figure 21: Results for training ViT model [6] with batch size of 256.

1074 term up to about 25% and 50% that of the maximum throughput obtained. We present the median
 1075 obtained throughput across 5 random trials in Table 7. We see that even when the cost-function prior
 1076 is adversarially constructed (by knowingly adding an incorrect term into the cost function), we are
 1077 still able to obtain good performances to the unperturbed cost-model prior even if the convergence
 1078 is slightly slower. This suggests that even in this extreme case, the GP is able to correct for the
 1079 inaccuracies in the cost-prior effectively.

1080 In practice, prior misspecification typically will due to the cost function not being sufficiently complex
 1081 to match the actual parallel training dynamics, because of incomplete or inaccurate domain knowledge
 1082 about the actual system rather than due to an adversarial construction of the cost function. This is the
 1083 case in the cost function we have chosen in our paper where there is a discrepancy between the cost
 1084 function alone and the actual throughput as demonstrated in Fig. 6c, resulting from our cost function
 1085 not modeling the effects of all hyperparameters in the PC. Under practical scenarios, we therefore
 1086 would not expect the results to be as extreme as what we have seen in the presented results, and that a
 1087 GP should be able to effectively model the throughputs.

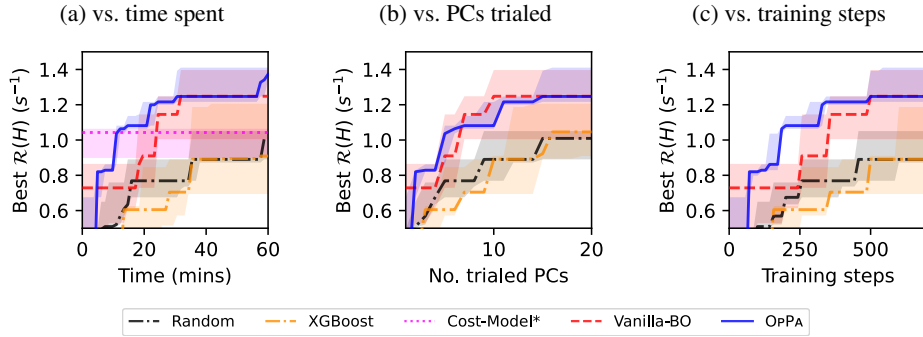


Figure 22: Results for training Mixture-of-Experts [30] models.

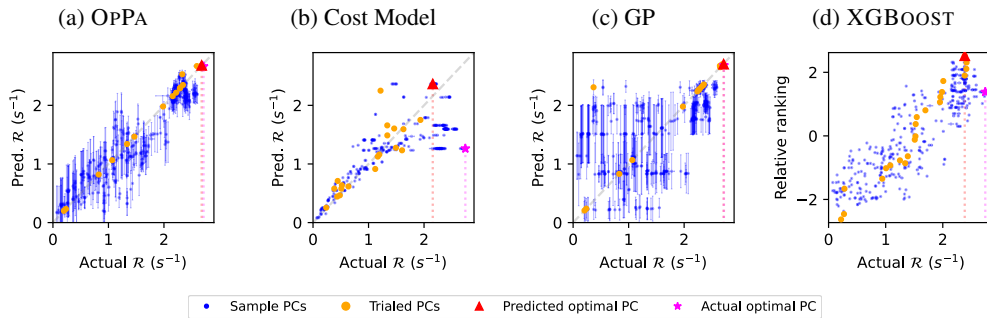


Figure 23: Comparison of modeled throughput vs. the actual throughput for training of BERT model on 8 GPUs after 20 training trials.

1088 **I.7 Additional Results for Ablation Studies of Components in OPPA**

1089 In Fig. 27, we demonstrate how early termination and parallelism-informed GP prior beliefs affect
 1090 the overall achieved throughput. First, we see that when parallelism-informed GP prior beliefs are
 1091 used, the performance is no worse than when no GP prior belief is used, although this benefit is more
 1092 pronounced for the Qwen2-1.5b example, possibly due to the increased complexity in the training
 1093 setup. Meanwhile, with early termination, we see that the performance is better in terms of time and
 1094 number of training steps needed, while not sacrificing the performances when considering the number
 1095 of PCs that are trialed to achieve a certain performance. This shows that the benefit gain comes
 1096 from being able to shorten the duration of the training trials while not sacrificing the throughput
 1097 predictions.

1098 **I.8 Effect of q_{\min} on OPPA performance**

1099 In Fig. 28, we see how the choice of q_{\min} affects the achieved throughput. For the BERT case, we
 1100 see that early termination clearly improves the time required for optimization, as seen where when
 1101 $q_{\min} = q_{\max}$ the obtained PC is the worst when all methods are allotted the same amount of time. As
 1102 q_{\min} decreases, we see that there is less drop in the amount of time required since each training step
 1103 is dominated by the time to setup each training trial. However, we still see that when we plot the
 1104 number of training steps for the optimization, we see that a smaller q_{\min} will require fewer training
 1105 steps to arrive at the same optimal PC. However, this trend breaks down when q_{\min} is too small, likely
 1106 since the value obtained is too noisy to give good information. Nonetheless, even in this case, we still
 1107 obtain good results. For the Qwen2-1.5b case, similar trends can be seen where reducing q_{\min} is able
 1108 to reduce the time and the number of training steps required to find the optimal PC up to a certain
 1109 point.

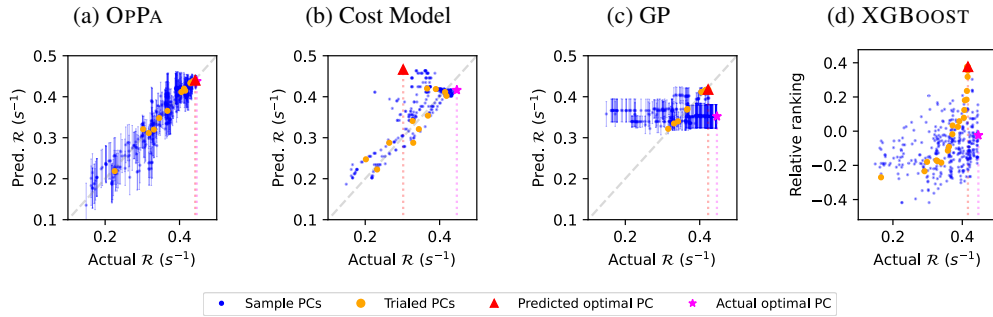


Figure 24: Comparison of modeled throughput vs. the actual throughput for training of Qwen2-1.5b model on 8 GPUs after 20 training trials.

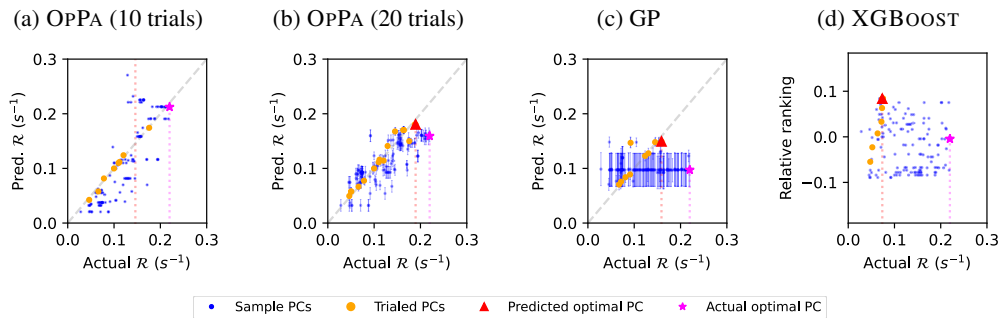


Figure 25: Comparison of modeled throughput vs. the actual throughput for training of Llama-2 model on 32 GPUs for, in order, OPPA after 10 trials, OPPA after 20 trials, GP after 10 trials and XGBOOST after 10 trials. Note that among the three methods, only OPPA ran for up to 20 trials given the time constraint.

1110 J Limitations and Broader Impacts

1111 While there are many other factors of parallel training that we could optimize, we have mainly
 1112 considered hyperparameters commonly found in existing parallel NN training frameworks which
 1113 would generally be tuned manually by practitioners. We believe that OPPA can also be extended to
 1114 optimize in other more sophisticated PC search spaces that may arise in practice as well. Maximizing
 1115 the throughput during NN training would allow the same amount of computation to possibly be done
 1116 in a more efficient manner, both in terms of time and compute resources. While this may allow faster
 1117 development of NNs for both good and bad use cases, overall it would still have a positive impact
 1118 since it allows for higher efficiency which reduces waste in computation time and other feasible
 1119 resources that come with it.

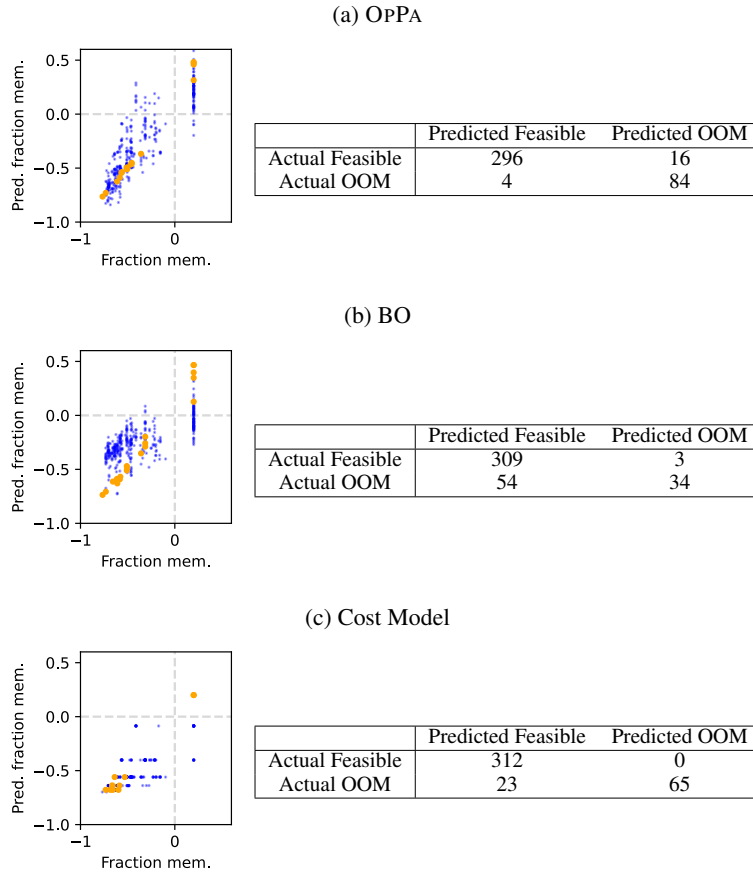


Figure 26: Comparison of modeled maximum memory usage vs. the actual maximum memory usage for training of Qwen2-1.5b model on 8 GPUs after 20 training trials. In each row, the left hand graph shows the predicted value vs. the actual value (with the predictive variance are omitted for clarity), while the right hand table is the confusion matrix.

Table 6: Example of the log lengthscales learned by the kernel of the GP for the throughput. The bolded value are to highlight hyperparameters with particularly shorter lengthscales.

Hyperparameter	BERT, 8 GPUs	Qwen2-1.5b, 8 GPUs
DP size (dp)	1.877	5.233
TP size (tp)	-1.071	6.024
PP size (pp)	2.293	-1.849
DP bucket size	3.484	5.994
ZeRO stage	0.434	5.279
ZeRO bucket size	3.402	6.516
Overlap ZeRO communication	3.949	7.581
Overlap DP AllGather	3.822	7.589
# microbatches (mb)	-1.322	-1.741
# model chunks (mc)	-0.293	0.124
Overlap P2P for PP	3.854	6.974
Grad. checkpointing	-0.334	6.586

Table 7: Effect of perturbing the GP prior belief on the resulting optimization process. The values reported are the median throughput obtained for the PC found after certain number of minutes of running OPFA with the perturbed prior belief.

Percent perturbation magnitude	10 mins. of search	20 mins. of search
0 (original prior)	0.425	0.446
About 25	0.424	0.447
About 50	0.415	0.447

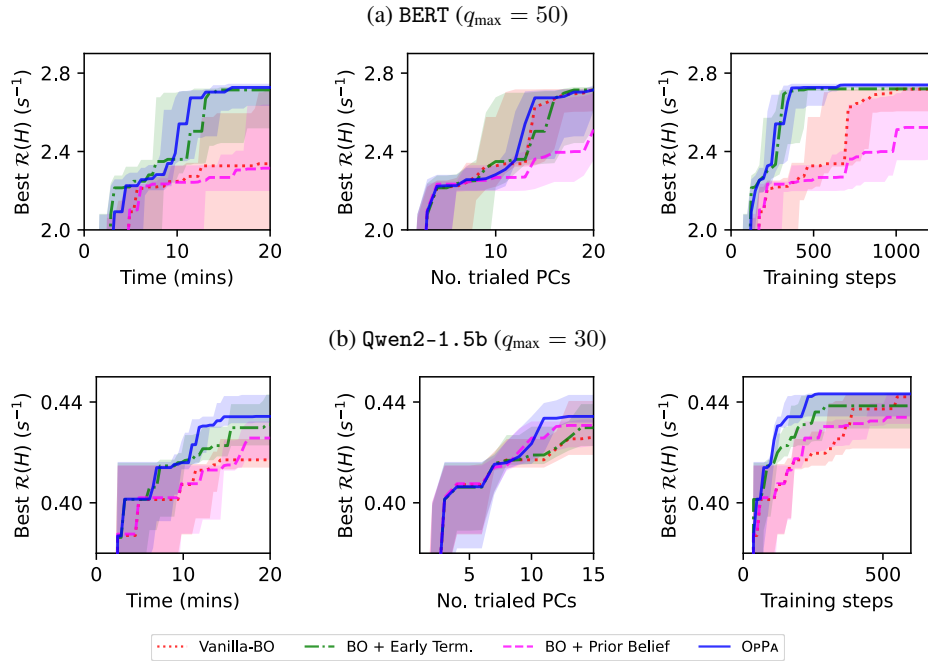


Figure 27: The effects of different components of OPBA on the optimal PC found. In each row, we present the results in a certain training setup as noted in the subheader. We present the obtained throughput versus, from left to right, the time the method has been ran for, the number of unique PCs trialed, and the number of training steps that has been ran across all of the trials.

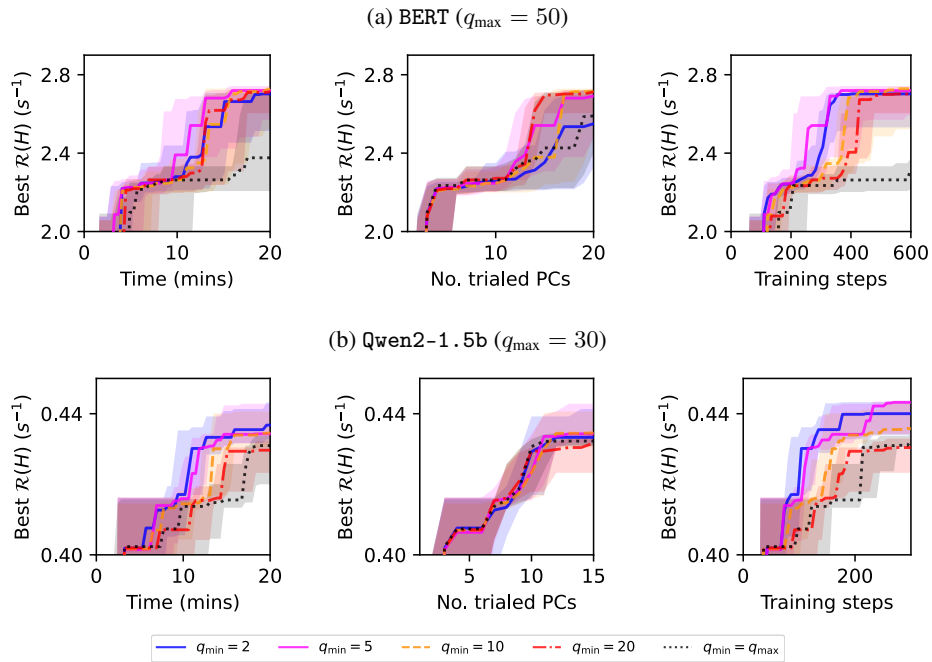


Figure 28: The effects of q_{\min} on the optimal PC found. In each row, we present the results in a certain training setup, where we present the obtained throughput versus, from left to right, the time the method has been ran for, the number of unique PCs trialed, and the number of training steps that has been ran across all of the trials.