
A Deep Learning Dataloader with Shared Data Preparation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Parallely executing multiple training jobs on overlapped datasets is a common
2 practice in developing deep learning models. By default, each of the parallel jobs
3 prepares (i.e., loads and preprocesses) the data independently, causing redundant
4 consumption of I/O and computations. Although a centralized cache component
5 can reduce the redundancies by reusing the data preparation work, each job's
6 random data shuffling results in a low sampling locality causing heavy cache
7 thrashing. Prior work tries to improve the sampling locality by enforcing all the
8 training jobs loading the same dataset in the same order and pace. However, such
9 a solution is only efficient under strong constraints: all jobs are trained on the
10 same dataset with the same starting moment and training speed. In this paper, we
11 propose a new data loading method for efficiently training parallel DNNs with much
12 flexible constraints. Our method is still highly efficient when different training
13 jobs use different but overlapped datasets and have different starting moments
14 and training speeds. To achieve this, we propose a dependent sampling algorithm
15 (DSA) and a domain-specific cache policy. Moreover, a novel tree data structure is
16 designed to efficiently implement DSA. Based on the proposed technologies, we
17 implemented a prototype, named JOADER, which can share data preparation work
18 as long as the datasets are overlapped for different training jobs. We evaluate the
19 proposed JOADER, showing a greater versatility and superiority of training speed
20 improvement (up to 200% in ResNet18) without affecting the accuracy.

21 1 Introduction

22 The rapid development of deep learning frameworks and tools [3, 22] has enabled and facilitated
23 researchers and practitioners in a multitude of disciplines (e.g., physics [28], chemistry [18],
24 biology [12], and Earth science [23] in addition to computer science) to start developing their own
25 DNN models solving various problems at hand. In such cases, they usually need to simultaneously
26 run multiple training jobs on the same dataset or overlapped datasets¹, due to the purpose of model
27 selection [15], hyper-parameter tuning [6], network architecture searching [25], etc.

28 Although widely used, the training of DNNs is usually time-consuming and tricky, which may affect
29 the development efficiency, especially for general users from fields other than artificial intelligence.
30 In literature, various approaches have been proposed to reduce the DNN training time such as
31 data preparation optimization [4], communication overhead reduction [10, 9, 29], GPU memory
32 optimization [8, 16, 26], and compiler-based operator optimization [7, 13, 17].

¹This could happen when we train models on different datasets sharing a particular portion of data, e.g., running the full experiment on ImageNet and further tuning parameters on TinyImageNet [11]).

33 Different from the above existing work, our focus is on the data preparation aspect, or more
34 specifically, the data loading aspect, in the context of parallel DNN training. On the one hand,
35 Mohan et al. [19] have observed that data preparation, which consists of data loading and the
36 subsequent data preprocessing, occupies a significant portion of parallel DNN training time. This is
37 due to the fact that multiple training jobs on the same data copy essentially create redundant data
38 preparation efforts. Ideally, to reduce the load of hardware and boost the training efficiency, the
39 prepared data should be cached in memory and serve for multiple jobs. However, the cache may
40 thrash with a high probability as each job adopts a fully random shuffling strategy that independently
41 shuffles the dataset in each epoch. On the other hand, it is until recently have some studies investigated
42 the shuffling process itself [20, 19, 21]. These studies investigate to what extent the fully random
43 shuffling is required and have observed that various shuffling strategies breaking the full randomness
44 still yield competitive results. We name such randomness as *correlated randomness* as opposed to
45 the full randomness. For example, prior work [19, 20] proposed to solve the cache thrashing issue by
46 enforcing that all jobs iterate over the dataset in the same order and pace, and observed no accuracy
47 loss. However, such a solution maximizes its performance under some strong constraints, i.e., all
48 training jobs have to start simultaneously, use the same dataset, and have similar training speeds.

49 This paper proposes a new data loading method for training multiple parallel jobs. Specifically, we
50 aim to relax the constraints of prior work in three aspects: 1) all jobs can start and end freely, 2)
51 different datasets can be used as long as they are partially overlapped, and 3) jobs' training speeds
52 may vary widely. For the above purposes, we first propose the *dependent sampling algorithm* (DSA)
53 to schedule the sampling for each job while ensuring their correlated randomness. DSA is inspired
54 by the correlated sampling theory [5] and it mainly involves three steps: 1) divide the dataset into
55 two parts of intersection set and difference set, 2) correlate the selections between intersection and
56 difference, and 3) share some operations for the intersection set to improve locality. Furthermore, we
57 design a domain-specific cache policy *RefCnt* that is tailored for parallel DNN training. It is aware of
58 the fact that all jobs iterate over the dataset once in an epoch, and thus it can evict data that is least
59 likely to be used in the near future. Finally, we design a novel data structure *dependent sampling tree*
60 to efficiently implement the DSA by elaborately organizing the intersection sets and difference sets.
61 Based on the above technologies, we implement a prototype system named JOADER and integrate it
62 into PyTorch for evaluation. The evaluation results show that JOADER can boost training efficiency in
63 more flexible cases. For example, it achieves up to 200% training speedup when models of different
64 sizes are parallelly trained without affecting the accuracy.

65 The main contributions of this paper include:

- 66 • A new dataloader for parallel DNN training on overlapped datasets.
- 67 • A sampling algorithm to increase the sampling locality with guaranteed randomness. The algorithm
68 reaches the global optima when there are two training jobs.
- 69 • A domain-specific cache policy and a novel tree-based data structure for efficient implementation.
- 70 • Experimental evaluations demonstrating the performance improvements.

71 The rest of this paper is organized as follows. Section 2 introduces the background information.
72 Section 3 describes the dependent sampling algorithm and our cache policy, followed by the
73 description of the dependent sampling tree in Section 4. Section 5 presents the evaluation of
74 the implemented dataloader JOADER. The paper is concluded in section 6.

75 **2 Background**

76 **2.1 Problem Statement**

77 In each epoch, a DNN training job sends data requests to the storage device and passes through
78 the dataset in random order. For the brevity of the presentation, we assume that each data request
79 involves only one data point/element, and we name the data access order over all the data elements
80 as the *access path* for each epoch. Consider an example in Figure 1, where we assume two training

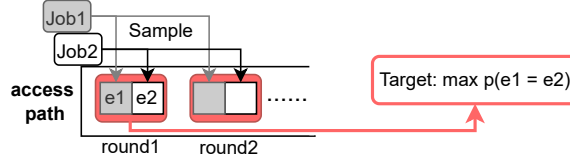


Figure 1: The access path of two training jobs. The target is to maximize the probability of sampling the same element in each round for all the jobs.

81 jobs are running, and the two sampled data elements ('e1' and 'e2') requested by the two jobs are
 82 packaged and prepared together in one *round*.

83 Cache can hold the prepared data to serve the future duplicated requests. To fully utilize the cache
 84 when multiple DNN models are being trained simultaneously in a server, one could generate a
 85 cache-friendly access path by encouraging all jobs to request the same data within a short period, i.e.,
 86 improving the sampling locality in each round. Meanwhile, we still need to ensure the correlated
 87 randomness for each job. To this end, we formulate the problem as follows:

88 **Problem 1** Assume there are n training jobs on overlapped datasets, where the i -th job J_i trains a
 89 model on dataset D_i with cardinality $|D_i|$. In each round, $p(e_j^i) = \frac{1}{|D_i|}$ indicates the probability of
 90 data element e_j being picked via sampling uniformly at random by J_i . For a given round, the goal is
 91 to maximize $p(e^1 = e^2 = \dots = e^n)$, while maintaining $p(e^i) = \frac{1}{|D_i|}$.

92 2.2 Related Work

93 **Dataloader for multiple training jobs.** Cerebro [20] is proposed to avoid redundant data preparation
 94 work for the task of model selection. It partitions the dataset across the servers into clusters and hopes
 95 the models iterate over data from one server to another instead of shuffling data. Although it can help
 96 avoid redundant data preparation work across servers, the work is still repeated on a single server.
 97 CoordDL [19] provides a more general solution for hyper-parameter optimization. In one turn, it
 98 distributes one batch to each job and caches the batch into the buffer for future usage. The batches are
 99 evicted until all jobs have consumed them in an epoch. However, CoordDL still has some limitations:
 100 1) all jobs must be set up with the same batch size; 2) when the jobs have different training speeds,
 101 the faster jobs must wait for the slower ones; 3) CoordDL cannot handle situations where the jobs
 102 arrive at different moments or are trained on different datasets.

103 **Domain-specific caching.** Crafting a cache for a specific domain is not a new idea [24]. For example,
 104 Quiver [14] uses local solid-state drive (SSD) caches to eliminate the impact of slow reads from the
 105 remote storage. However, it is too expensive to prepare enough memory against cache thrashing for
 106 DNN training, especially considering the fact we need to cache not only the previous reading but also
 107 the results of data preprocessing. In this paper, we propose a well-designed sampling algorithm to
 108 improve the sampling locality with a small cache.

109 **Correlated sampling.** As a theoretical problem, correlated sampling [5] aims to minimize the
 110 probability of two sampling results being equal, given that the two players are sampling uniformly at
 111 random from two subsets of the same set. Some sampling strategies [27] give the optimal solution for
 112 two players. In this paper, we practically solve a related problem by designing an algorithm and a
 113 novel data structure for efficient implementation while enjoying their theoretical results.

114 3 Dependent Sampling Algorithm

115 3.1 Algorithm Design

116 To solve Problem 1, we propose a dependent sampling algorithm (DSA) to maximize the sampling
 117 locality while ensuring correlated randomness. In the following, we first present the algorithm in
 118 two-job case and then extend it to the n -job case.

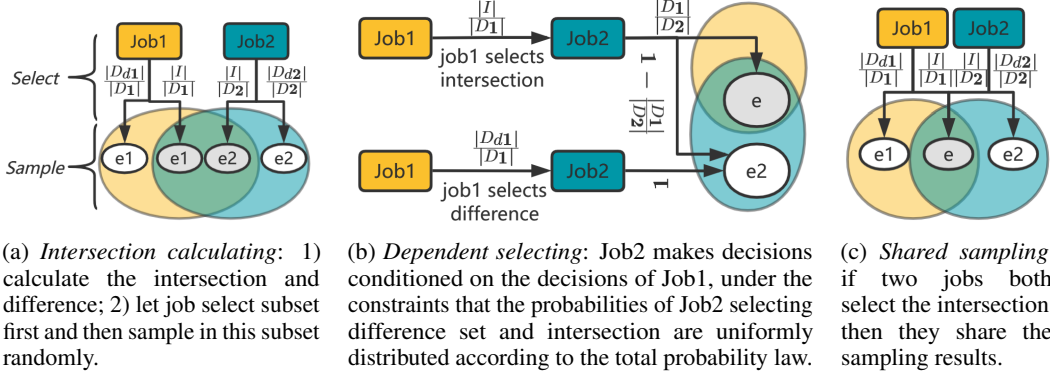


Figure 2: Dependent sampling algorithm in the two-job case.

119 **Two-job Case.** The proposed *dependent sampling algorithm* (DSA) involves three steps: 1)
 120 intersection calculating, 2) dependent selecting, and 3) shared sampling. For a better illustration, we
 121 assume two jobs J_1 and J_2 are being trained on two overlapped datasets D_1 (the yellow circle) and
 122 D_2 (the green circle), respectively, as shown in Figure 2.

123 First, we divide the datasets into three subsets: *intersection* set $I = D_1 \cap D_2$ and two *difference*
 124 sets $D_{d1} = D_1 \setminus I$, $D_{d2} = D_2 \setminus I$. For the default *independent sampling algorithm* (ISA), which is
 125 currently used by PyTorch, it then processes the data sampling uniformly at random in the dataset
 126 for each job independently. As shown in Figure 2a, the ISA algorithm can be divided into two steps:
 127 *selecting* and *sampling* procedures. In the selecting procedure, each job selects between its difference
 128 and intersection sets for sampling. In the sampling procedure, each job randomly picks an element e
 129 from the subset chosen in the previous step. This two-step process maintains the uniform distribution
 130 of sampling for each job.

131 Then, the proposed DSA takes a careful modification in the sampling procedure. As shown in
 132 Figure 2c, if the two jobs both select the intersection set, one shares the sampling result to another.
 133 We call this modification *shared sampling*. Based on shared sampling, the probability that two jobs
 134 selecting the same element is equal to that both jobs selecting the intersection set.

135 To further improve locality, we need to improve the probability of both two jobs selecting the
 136 intersection, which can be achieved according to the conditional probability, as shown in Figure 2b.
 137 The independent events in the selecting procedure can be transformed into the dependent ones,
 138 making the probability of the second job conditional. That is, J_2 makes decisions conditioned on the
 139 decision of J_1 .

140 To maximize the sampling locality, we formulate the conditional probabilities of dependent selecting
 141 in case of $|D_1| < |D_2|$. For J_1 , the dependent selecting procedure is the same as that in ISA: selecting
 142 I with the probability of $\frac{|I|}{|D_1|}$, and D_{d1} with the probability of $\frac{|D_{d1}|}{|D_1|}$. The selecting procedure of
 143 J_2 makes the decision according to J_1 . If J_1 selects D_{d1} , then J_2 must select another difference
 144 set D_{d2} . If J_1 select I , then J_2 selects I with the probability of $\frac{|D_1|}{|D_2|}$ and D_{d2} with probability
 145 $\frac{|D_2| - |D_1|}{|D_2|}$. Figure 2b shows the above selecting procedure. For the case of $|D_2| < |D_1|$, we just need
 146 to exchange the order of J_1 and J_2 to apply the above rules.

147 According to above conditional probabilities, the probability of $p(e_j^1 = e_i^2)$ is $\frac{|I|}{\max(|D_2|, |D_1|)}$, which
 148 is also the theoretical upper bound as discussed in Theorem 1. The proof of this theorem can be found
 149 in Proof 1.

150 **Theorem 1** Assume there are 2 jobs J_1, J_2 sampling uniformly at random from two datasets D_1, D_2
 151 respectively, and the elements they sampled are e_j^1, e_i^2 . The probability of these two elements being
 152 equal cannot be larger than $\frac{|I|}{\max(|D_2|, |D_1|)}$ where $I = D_1 \cap D_2$.

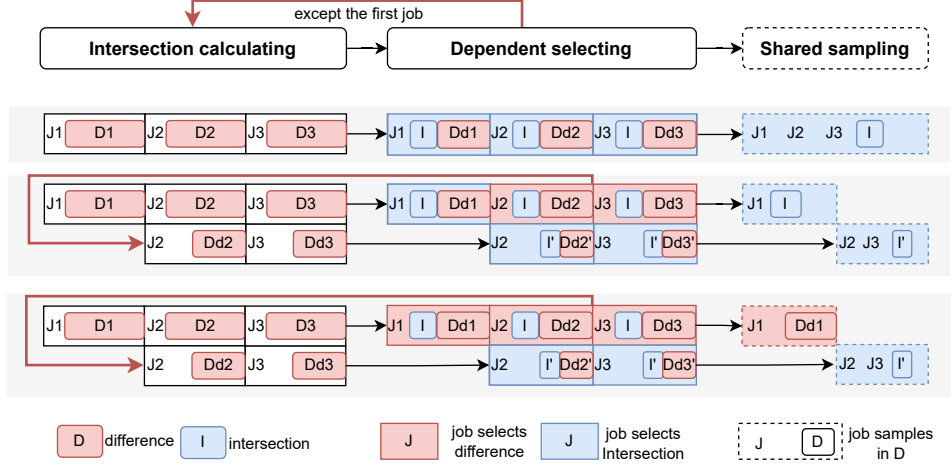


Figure 3: An illustration for the n -job case of DSA. We assume three jobs $\{J_1, J_2, J_3\}$ with datasets $\{D_1, D_2, D_3\}$, and there are three possible sampling cases: 1) all jobs select intersection, 2) some jobs select intersection and some select difference, and 3) all jobs select difference. Note that if there is one job selecting difference, all following jobs must select difference.

153 **N-job Case.** DSA of n -jobs is similar to two-jobs at the beginning. The difference happens when
 154 some jobs select difference sets. We need to go through the above procedure recursively for these
 155 jobs rather than sampling in the current difference set directly, as shown in Figure 3. The reason for
 156 recursive execution is that some data may only be shared by part of the jobs. In the extreme case, if
 157 all datasets are mutually disjoint, then the intersection in the first recursion is empty and the DSA is
 158 useless. Therefore, we need to further check whether there are shared data for part of the jobs.

159 The conditional probability for n -jobs that are sorted by the cardinalities of their datasets is formulated
 160 as follows. 1) The first job J_1 selects the intersection with the probability of $\frac{|I|}{|D_1|}$ while selecting the
 161 difference with the probability of $\frac{|D_{d1}|}{|D_1|}$. 2) If the job J_{i-1} selects the difference set, then job J_i must
 162 select the difference set. If job J_{i-1} selects the intersection set, then job J_i selects the intersection
 163 with the probability of $\frac{|D_{i-1}|}{|D_i|}$ while selecting the difference set with the probability of $\frac{|D_i| - |D_{i-1}|}{|D_i|}$.
 164 For the above algorithm, we can observe that if one job selects the difference, then all subsequent
 165 jobs will choose the difference, which is an important feature for the data structure in Section 4.

166 The pseudo code of DSA are shown in Appendix B for completeness. Here, we provide the theoretical
 167 result as follows, whose proofs can be found in Appendix D.2.

168 **Theorem 2** Assume there are multiples jobs with the i -th job J_i trained on dataset D_i . The
 169 probability of any element e being picked for J_i via DSA is $p(e) = \frac{1}{|D_i|}$.

170 3.2 RefCnt Cache Policy

171 If we allow jobs to vary greatly in speed, then some data may be evicted before they are consumed
 172 by all jobs. Therefore, we need a cache policy that can take out the data that is least likely to be
 173 used (i.e., the data that least jobs will request) in the near future. Classical cache policies (e.g., LRU
 174 and LFU) are not suitable for deep learning scenarios, where each training job reads each data only
 175 once in an epoch. So the total number of request to each data equals to the number of jobs using the
 176 data for training. To this end, if we could record the count of data request, the number of request to
 177 this data in near future could be derived. The fewer jobs that access the data, the higher the eviction
 178 priority of that data.

179 Based on the above understanding, we design the caching policy *RefCnt*. We maintain a reference
 180 count for each data element. If a job will request the element in the future, we increase its reference

181 count; if a job reads the data, we decrease its reference count. When the cache is full, the data with
 182 the lowest reference count will be evicted.

183 4 Dependent Sampling Tree

184 The time complexity of DSA is $O(nm)$, where n is the number of jobs and m is the maximum
 185 cardinality of the datasets. The worst case happens if all jobs select the difference set in each recursion,
 186 which means we need n recursions for n jobs. Each recursion has the time complexity of $O(m)$.²
 187 To reduce the time complexity of sampling, we can maintain the intersections for DSA to avoid the
 188 intersection calculation in each recursion, which can make each recursion cost constant time. In the
 189 worst case of DSA, the complexity of DSA is reduced from $O(nm)$ to $O(n)$.

190 We propose the dependent sampling tree for organizing the intersections. Although the number of
 191 intersections is $2^n - 1$ for n sets, only $n - 1$ intersections are valuable. In DSA, we only need to
 192 calculate the intersection of datasets for those jobs that select the difference, as shown in Figure 3.
 193 When the jobs are sorted by the cardinality of their datasets, the job selecting the difference will
 194 make all following jobs select the difference sets. Therefore, the possible combinations for jobs
 195 that select the difference must be the suffixes of the sorted jobs array. For example, if there are
 196 three sorted jobs $[J_1, J_2, J_3]$, the groups of jobs choosing the difference set can only be one of the
 197 following three sets, including $[J_1, J_2, J_3]$, $[J_2, J_3]$, and $[J_3]$. Thus, the sets to be calculated are
 198 $\{D_1 \cap D_2 \cap D_3, D_2 \cap D_3\}$.

199 4.1 Tree Definition

200 For n datasets $\{D_1, D_2, \dots, D_n\}$ in ascending order w.r.t. $|D_1| < |D_2| < \dots < |D_n|$, we need
 201 to calculate $n - 1$ intersections $\{D_1 \cap \dots \cap D_n, D_2 \cap \dots \cap D_n, \dots, D_{n-1} \cap D_n\}$. We build a
 202 dependent sampling tree, where each intersection is an internal vertex, and each difference set is
 203 a leaf. Meanwhile, all vertices in a path from the root to a leaf represent a complete dataset. The
 204 definition of dependent sampling tree is as follows

205 **Definition 1** *Dependent sampling tree is a binary tree that stores a collection of datasets*
 206 *$\{D_1, D_2, \dots, D_n\}$, with $|D_1| < |D_2| < \dots < |D_n|$. In the tree, the root is the intersection of*
 207 *all datasets $I = D_1 \cap D_2 \dots \cap D_n$. The right child is the difference of the smallest dataset $D_1 \setminus I$.*
 208 *The left child is a new dependent sampling tree of $\{D_2 \setminus I, D_3 \setminus I, \dots, D_n \setminus I\}$. When there is only*
 209 *one dataset, the tree has one vertex that contains the dataset.*

210 If there are n datasets, the height of the dependent sampling tree is $n - 1$. Figure 4a shows an
 211 example with three sets $\{D_1, D_2, D_3\}$ with $|D_1| < |D_2| < |D_3|$. In this tree, vertex A contains
 212 $D_A = D_1 \cap D_2 \cap D_3$, and vertex B contains $D_B = (D_2 \setminus D_A) \cap (D_3 \setminus D_A)$. For the leaf vertices
 213 C, E, and F, they maintain the other differences.

214 4.2 Tree Operations

215 When executing the DSA algorithm, dependent sampling tree should maintain the datasets into
 216 intersection and difference sets dynamically. Furthermore, dependent sampling tree should also
 217 support the situations when a job finishes its execution or a new job starts. In this part, we discuss the
 218 operations provided by dependent sampling tree, which are *sampling*, *deletion*, and *insertion*.

219 **Sampling.** The intersections of the flow in Figure 3 can be represented by one or multiple vertexes of
 220 the rightmost path in the dependent sampling tree. Thus, the process of DSA is to traverse the path
 221 from the root to the most-right leaf in the tree. When traversing to an internal vertex, the jobs will

²1) The intersection calculation needs to traverse the dataset so its complexity is linear to the cardinality of the dataset. 2) In the worst case, only the first job needs to make the decision, and others follow it. Therefore, dependent selecting needs constant time. 3) Sampling only needs to pick a random number between 0 and the size of the dataset, which needs constant time.

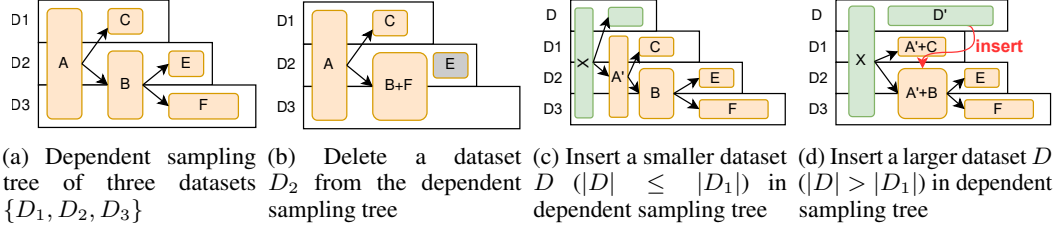


Figure 4: Illustration for the proposed dependent sampling tree.

222 decide whether to select this vertex according the conditional probability. For those jobs that do not
 223 select this vertex, they will make decisions in the right child tree.

224 Sometimes the intersection can be composed of multiple vertexes. For example, if only the job of D_3
 225 selects difference at A in Figure 4a, then it can sample both in B and F . In these cases, the algorithm
 226 needs to first select a vertex that belongs to the intersection and then sample.

227 **Dataset Deletion.** When a job finishes training or is killed, we should delete its dataset from the tree.
 228 The *deletion* involves two steps: 1) deleting the corresponding leaf vertex, and 2) merging its parent
 229 vertex with its sibling vertex. For example, if we want to delete D_2 in Figure 4a, we need to delete
 230 the leaf vertex E and merge vertex B with F , as shown in Figure 4b. The time complexity is $O(1)$
 231 since deleting and merging both cost constant time.

232 **Dataset Insertion.** The *insertion* operation depends on the cardinality of the new dataset D . When D
 233 is smaller than the smallest set in the sampling tree, the procedure involves three steps: 1) calculating
 234 the intersection of it and the root of the tree, 2) setting the new intersection as the new root of the
 235 dependent sampling tree, and 3) setting the difference set of the old root as the right child and the
 236 difference of D as the left child for the new root. For example, Figure 4c shows the new tree after
 237 inserting a smaller dataset D . The new root X is the intersection of D and the set in vertex A . The
 238 old tree becomes the right sub-tree, and the difference of the dataset D becomes the left leaf.

239 When the new dataset D is larger than the smallest dataset in the sampling tree, we need to insert it
 240 recursively: 1) calculating the intersection of D and root A and set it as new root X , 2) adding the
 241 difference of A to its children, 3) deleting the old root A and adding their children to the new root, 4)
 242 recursively executing step 1 on X 's right child with the difference set of D until the difference set is
 243 the smallest in the current sub-tree. An example of this procedure is shown in Figure 4d. In the worst
 244 case, the time complexity of *insertion* is $O(n|D|)$ if the new dataset D is the largest³.

245 5 Evaluation

246 In this section, we evaluate JOADER on ImageNet with the family of ResNet models. We denote the
 247 default dataloader strategy in PyTorch as the 'Baseline' method, and further compare JOADER with
 248 the state-of-the-art method CoordDL [19]. Due to the space limit, we mainly show two experiments in
 249 the paper and more experimental results can be found in the appendix. First, we profile the server
 250 when there are multiple training jobs to show the bottlenecks and compare JOADER with CoordDL
 251 and Baseline in synchronous case. Second, we evaluate JOADER in asynchronous cases: 1) multiple
 252 models trained at different speeds, 2) multiple jobs arrived at different moments, and 3) multiple
 253 datasets only partially overlapped. For the results in the appendix, we summarize them as follows: 1)
 254 ablation studies show that both DSA and RefCnt improve the training efficiency, 2) the I/O speed is
 255 not the bottleneck for data loading, and 3) we train ResNet18 with PyTorch and JOADER separately
 256 to show JOADER does not affect the accuracy.

257 The evaluated models are the basic models with their default settings in torchvision [1], and trained
 258 on top of the PyTorch 1.6.0 DL framework. The experiments were conducted on a GPU server with

³In practice, the dataset is a bitmap set of data indices in the dataset, e.g., a unique integer id that identifies data. Therefore, the insert operation is very fast, which only takes a little time compared with training time.

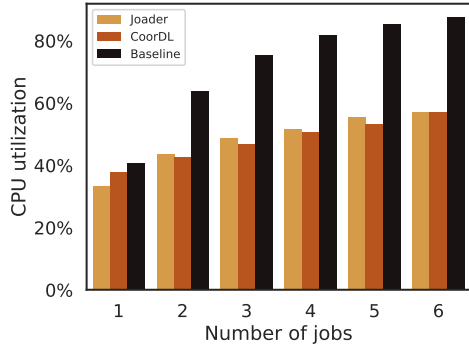


Figure 5: The CPU utilization when training multiple ResNet18 models simultaneously. CPU quickly becomes the bottleneck for the Baseline dataloader of PyTorch, while both CoordDL and JOADER can greatly reduce CPU utilization.

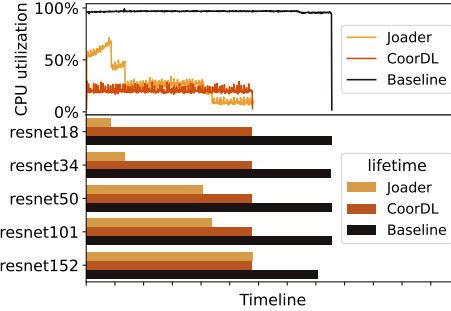
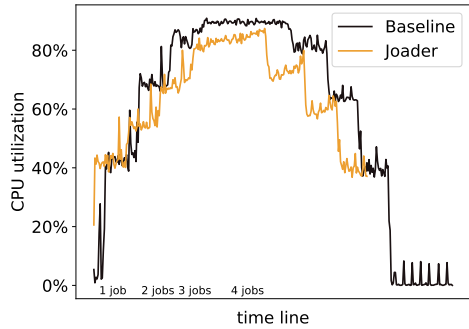
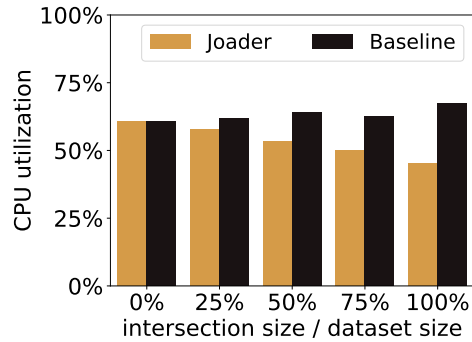


Figure 6: The CPU utilization and lifetime when training 5 models with different speeds. JOADER can reduce the CPU utilization and keep each job’s own speed. In contrast, fast job must wait for the slow job for CoordDL (e.g., ResNet18 costs the same time with ResNet152).



(a) The CPU utilization and lifetime when four jobs start at different moments. JOADER is faster and causes less CPU load because it can reuse the data preparation work from previous jobs.



(b) The CPU utilization when datasets are partially overlapped. JOADER can reduce 10% CPU utilization when half of datasets are in the intersection set for 2 jobs.

Figure 7: Training jobs arrive at different moments and have different datasets

259 two Intel Xeon Gold 5118 CPUs @ 2.30GHz (24 physical cores and 48 threads), 500GB RAM, and
 260 6 TITAN RTX GPUs. The server ran Ubuntu 18.04 with GNU/Linux kernel 4.15.0. The disk is
 261 Symbios Logic MegaRAID SAS-3 3316 of 1GB/s read speed.

262 5.1 Performance in Synchronous Cases

263 In this experiment, we start training multiple ResNet18 models with different hyper-parameters. **We
 264 try to make fair comparison by forcing the jobs training on the GPUs of the same version to get the
 265 theoretically same training speeds. All jobs train the same model on the same dataset and start at the
 266 same time.**

267 We start from training 1 model to training 6 models on ImageNet, and profile the server for different
 268 dataloaders. The experiments show that the CPU utilization tends to be the bottleneck, as shown in
 269 Figure 5.

270 The CoordDL and JOADER both can reduce the CPU utilization greatly by sharing the data preparation.
 271 Additionally, it shows that JOADER can achieve close performance to CoordDL when all training jobs
 272 start at the same time and have a similar speed. Note that CoordDL is designed for such synchronous
 273 cases while JOADER can also handle asynchronous cases as we will later show. The time cost of each

274 epoch is shown in Table 1. With 6 training jobs, JOADER saves 44.8% training time with over 40%
 275 less CPU utilization, compared the default Baseline in PyTorch.

Table 1: The time cost of different dataloaders when training multiple ResNet18 models in an epoch.

Number of jobs	1	2	3	4	5	6
CoorDL (min)	32.82	31.73	32.73	34.76	35.28	37.98
JOADER (min)	32.70	32.49	33.72	36.32	37.44	37.55
Baseline (min)	31.74	35.01	40.04	52.11	57.67	67.85

276 5.2 Performance in Asynchronous Cases

277 Next, we consider three asynchronous cases when multiple models differ in training speeds, arriving
 278 moments, and datasets.

279 **Multiple models with different training speeds.** We first simultaneously train five models
 280 (ResNet18, ResNet34, ResNet50, ResNet101, and ResNet152) with different training speeds, and
 281 show the lifetime and CPU utilization results in Figure 6. With the Baseline dataloader provided
 282 by PyTorch, each job preprocesses the data individually so that the CPU load is extremely heavy.
 283 CoorDL enforces multiple jobs requesting the data in lockstep to share the data preparation work, and
 284 it achieves less CPU load and higher training speed than Baseline. However, the faster jobs (smaller
 285 models) must wait for the slower jobs, making CPU idle during the training period. JOADER can both
 286 reduce the CPU load by sharing the preprocessing work and keep these jobs’ own training speeds.
 287 For example, ResNet18 runs $4\times$ faster using JOADER compared with CoorDL. Furthermore, this
 288 experiment suggests that our system has less profitable for larger models because GPU tends to be
 289 the bottleneck for those models.

290 **Multiple models arriving at different moments.** Next, we consider the situation that multiple DNN
 291 training jobs start at different moments. Since CoorDL must wait for all the jobs to be ready for
 292 training, we only compare JOADER with Baseline in this experiment. Figure 7a shows the CPU
 293 utilization when four ResNet18 training jobs arrive at different moments. JOADER takes up less CPU
 294 load than the Baseline. This is due to the fact Baseline repeats the data preparation work for new
 295 jobs, whereas JOADER reuses the prepared data from previous jobs. As a result, JOADER saves 17.1%
 296 training time and 8.7% CPU utilization.

297 **Multiple models with partially overlapping datasets.** Finally, we show the CPU utilization when
 298 the two datasets only partially overlap in Figure 7b. The datasets are re-sampled from ImageNet.
 299 We use the ratio between intersection size and dataset size to measure the overlapping degree. The
 300 larger the intersection, the more data to share. We still compare with the Baseline method only since
 301 CoorDL requires all jobs having the same access path over datasets. From the figure, we can observe
 302 that the CPU utilization of the same dataset is only 60% of that of different datasets. Additionally,
 303 the CPU utilization always decrease as the intersection size increases, showing that JOADER is valid
 304 with some overlapping elements. The reason is that the DSA algorithm can get the same sampling
 305 results with the max probability and the cache policy evicts the useless data with higher priority.

306 6 Conclusion

307 Accelerating DNN training is a fundamental challenge in deep learning. This paper presents an
 308 approach to boosting training efficiency by caching data prep work and reusing them to parallel
 309 training jobs. To prevent cache thrashing, we propose DSA to improve sampling locality with ensured
 310 randomness. We also propose a novel tree-based data structure to efficiently implement the proposed
 311 algorithm. Furthermore, a domain-specific cache policy is proposed for evicting data that is least
 312 used in the near future. At last, we evaluate the proposed approach and the results show significant
 313 improvement in asynchronous cases while achieving close performance in synchronous cases to
 314 CoorDL which is specially designed for such cases.

315 References

- 316 [1] The package consists of popular datasets, model architectures, and common image
317 transformations for computer vision. [https://pytorch.org/vision/stable/index.](https://pytorch.org/vision/stable/index.html)
318 [html](https://pytorch.org/vision/stable/index.html), 2021.
- 319 [2] Rust program language. <https://github.com/rust-lang/rust>, 2021.
- 320 [3] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu
321 Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A system for
322 large-scale machine learning. In *12th {USENIX} symposium on operating systems design and*
323 *implementation ({OSDI} 16)*, pages 265–283, 2016.
- 324 [4] Jonghyun Bae, Jongsung Lee, Yunho Jin, Sam Son, Shine Kim, Hakbeom Jang, Tae Jun Ham,
325 and Jae W Lee. {FlashNeuron}::{SSD-Enabled} {Large-Batch} training of very deep neural
326 networks. In *19th USENIX Conference on File and Storage Technologies (FAST 21)*, pages
327 387–401, 2021.
- 328 [5] Mohammad Bavarian, Badih Ghazi, Elad Haramaty, Pritish Kamath, Ronald L Rivest, and
329 Madhu Sudan. Optimality of correlated sampling strategies. *Theory of Computing*, 16(1), 2020.
- 330 [6] James Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. *Journal*
331 *of machine learning research*, 13(2), 2012.
- 332 [7] Tianqi Chen, Thierry Moreau, Ziheng Jiang, Lianmin Zheng, Eddie Yan, Haichen Shen, Meghan
333 Cowan, Leyuan Wang, Yuwei Hu, Luis Ceze, et al. {TVM}: An automated end-to-end
334 optimizing compiler for deep learning. In *13th {USENIX} Symposium on Operating Systems*
335 *Design and Implementation ({OSDI} 18)*, pages 578–594, 2018.
- 336 [8] Tianqi Chen, Bing Xu, Chiyuan Zhang, and Carlos Guestrin. Training deep nets with sublinear
337 memory cost. *arXiv preprint arXiv:1604.06174*, 2016.
- 338 [9] Jiawei Fei, Chen-Yu Ho, Atal N Sahu, Marco Canini, and Amedeo Sapia. Efficient sparse
339 collective communication and its application to accelerate distributed deep learning. In
340 *Proceedings of the 2021 ACM SIGCOMM 2021 Conference*, pages 676–691, 2021.
- 341 [10] Sayed Hadi Hashemi, Sangeetha Abdu Jyothi, and Roy H Campbell. Tictac: Accelerating
342 distributed deep learning with communication scheduling. *SysML 2019*, 2019.
- 343 [11] Guoxi Xu Jiayu Wu, Qixiang Zhang. Tiny imagenet challenge.
- 344 [12] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf
345 Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Židek, Anna Potapenko, et al.
346 Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- 347 [13] Taebum Kim, Eunji Jeong, Geon-Woo Kim, Yunmo Koo, Sehoon Kim, Gyeong-In Yu, and
348 Byung-Gon Chun. Terra: Imperative-symbolic co-execution of imperative deep learning
349 programs. *Advances in Neural Information Processing Systems*, 34, 2021.
- 350 [14] Abhishek Vijaya Kumar and Muthian Sivathanu. Quiver: An informed storage cache for deep
351 learning. In *18th {USENIX} Conference on File and Storage Technologies ({FAST} 20)*, pages
352 283–296, 2020.
- 353 [15] Arun Kumar, Robert McCann, Jeffrey Naughton, and Jignesh M Patel. Model selection
354 management systems: The next frontier of advanced analytics. *ACM SIGMOD Record*,
355 44(4):17–22, 2016.
- 356 [16] Gangmuk Lim, Jeongseob Ahn, Wencong Xiao, Youngjin Kwon, and Myeongjae Jeon. Zico:
357 Efficient {GPU} memory sharing for concurrent {DNN} training. In *2021 USENIX Annual*
358 *Technical Conference (USENIX ATC 21)*, pages 161–175, 2021.

- 359 [17] Zhongyi Lin, Evangelos Georganas, and John D Owens. Towards flexible and compiler-friendly
360 layer fusion for cnns on multicore cpus. In *European Conference on Parallel Processing*, pages
361 232–248. Springer, 2021.
- 362 [18] Adam C Mater and Michelle L Coote. Deep learning in chemistry. *Journal of chemical*
363 *information and modeling*, 59(6):2545–2559, 2019.
- 364 [19] Jayashree Mohan, Amar Phanishayee, Ashish Raniwala, and Vijay Chidambaram. Analyzing
365 and mitigating data stalls in dnn training. *Proceedings of the VLDB Endowment*, 14(5):771–784,
366 2021.
- 367 [20] Supun Nakandala, Yuhao Zhang, and Arun Kumar. Cerebro: A data system for optimized deep
368 learning model selection. *Proceedings of the VLDB Endowment*, 13(12):2159–2173, 2020.
- 369 [21] Thao Truong Nguyen, François Trahay, Jens Domke, Aleksandr Drozd, Emil Vatai, Jianwei
370 Liao, Mohamed Wahib, and Balazs Gerofi. Why globally re-shuffle? revisiting data shuffling in
371 large scale deep learning. In *IEEE International Parallel & Distributed Processing Symposium*,
372 2022.
- 373 [22] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
374 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative
375 style, high-performance deep learning library. *Advances in neural information processing*
376 *systems*, 32, 2019.
- 377 [23] Markus Reichstein, Gustau Camps-Valls, Bjorn Stevens, Martin Jung, Joachim Denzler, Nuno
378 Carvalhais, et al. Deep learning and process understanding for data-driven earth system science.
379 *Nature*, 566(7743):195–204, 2019.
- 380 [24] Michael Stonebraker. Operating system support for database management. *Communications of*
381 *the ACM*, 24(7):412–418, 1981.
- 382 [25] Haotian Tang, Zhijian Liu, Shengyu Zhao, Yujun Lin, Ji Lin, Hanrui Wang, and Song
383 Han. Searching efficient 3d architectures with sparse point-voxel convolution. In *European*
384 *Conference on Computer Vision*, pages 685–702. Springer, 2020.
- 385 [26] Yuke Wang, Boyuan Feng, Gushu Li, Shuangchen Li, Lei Deng, Yuan Xie, and Yufei Ding.
386 {GNNAdvisor}: An adaptive and efficient runtime system for {GNN} acceleration on {GPUs}.
387 In *15th USENIX Symposium on Operating Systems Design and Implementation (OSDI 21)*,
388 pages 515–531, 2021.
- 389 [27] Yun-Fu Wu. Correlated sampling techniques used in monte carlo simulation for risk assessment.
390 *International Journal of Pressure Vessels and Piping*, 85(9):662–669, 2008. Special Issue:
391 Advances in Structural Integrity of Nuclear Components in Asian Power Plants.
- 392 [28] Lenka Zdeborová. Understanding deep learning is also a job for physicists. *Nature Physics*,
393 16(6):602–604, 2020.
- 394 [29] Zhenwei Zhang, Qiang Qi, Ruitao Shang, Li Chen, and Fei Xu. Prophet: Speeding up distributed
395 dnn training with predictable communication scheduling. In *50th International Conference on*
396 *Parallel Processing*, pages 1–11, 2021.

397 **Checklist**

- 398 1. For all authors...
- 399 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
400 contributions and scope? [Yes]
- 401 (b) Did you describe the limitations of your work? [Yes]
- 402 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 403 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
404 them? [Yes]
- 405 2. If you are including theoretical results...
- 406 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 407 (b) Did you include complete proofs of all theoretical results? [Yes]
- 408 3. If you ran experiments...
- 409 (a) Did you include the code, data, and instructions needed to reproduce the main
410 experimental results (either in the supplemental material or as a URL)? [Yes] See
411 <https://anonymous.4open.science/r/Joader-6F4D>
- 412 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
413 were chosen)? [Yes]
- 414 (c) Did you report error bars (e.g., with respect to the random seed after running
415 experiments multiple times)? [No]
- 416 (d) Did you include the total amount of compute and the type of resources used (e.g., type
417 of GPUs, internal cluster, or cloud provider)? [Yes]
- 418 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 419 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 420 (b) Did you mention the license of the assets? [N/A]
- 421 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 422 (d) Did you discuss whether and how consent was obtained from people whose data you're
423 using/curating? [N/A]
- 424 (e) Did you discuss whether the data you are using/curating contains personally identifiable
425 information or offensive content? [N/A]
- 426 5. If you used crowdsourcing or conducted research with human subjects...
- 427 (a) Did you include the full text of instructions given to participants and screenshots, if
428 applicable? [N/A]
- 429 (b) Did you describe any potential participant risks, with links to Institutional Review
430 Board (IRB) approvals, if applicable? [N/A]
- 431 (c) Did you include the estimated hourly wage paid to participants and the total amount
432 spent on participant compensation? [N/A]

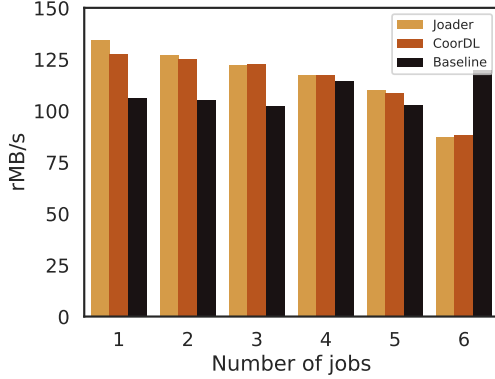


Figure 8: When training multiple ResNet18 models simultaneously, the I/O speed is shown above.

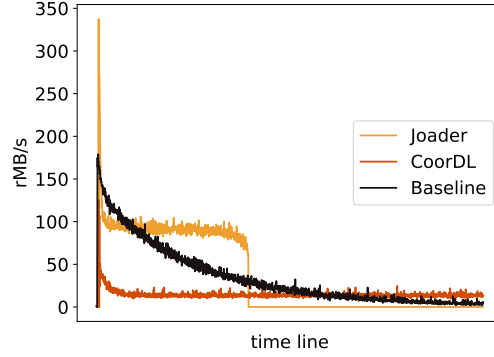


Figure 9: When training 5 models with different speeds at the same time, the I/O speeds are shown above.

433 A Evaluation

434 In this section, we first present the I/O speed in the synchronous cases and asynchronous cases to
 435 show JOADER can reduce the redundant I/O work. Then, we evaluate the algorithm DSA and RefCnt
 436 separately for the ablation test. Finally, we show the loss and accuracy trace of ResNet18 in 40
 437 epochs to demonstrate JOADER does not affect convergence speed. We denote the default dataloader
 438 strategy in PyTorch as the ‘Baseline’ method, and further compare JOADER with the state-of-the-art
 439 method CoordL. **At last, we show the cost of the operation in DSA.**

440 A.1 I/O speed

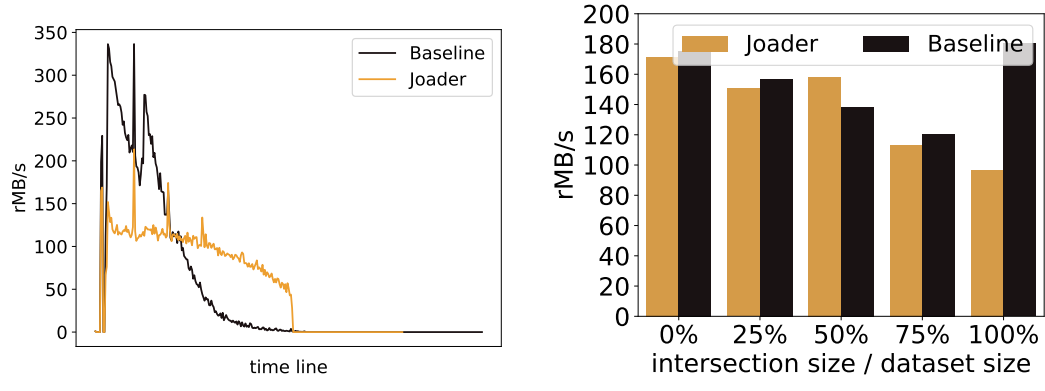
441 In this part, we show the I/O speed in the synchronous and asynchronous cases. Note that I/O is
 442 not the bottleneck in our server (Maximum I/O speed of the server is about 1GB/s). Figure 8 shows
 443 the I/O from training one job to training six jobs. The I/O speed of Baseline tends to increase while
 444 that of CoordL and JOADER tends to decrease, because CoordL and JOADER can reuse the data
 445 in memory. When we train 5 jobs (ResNet18, ResNet34, ResNet50, ResNet101, ResNet152) with
 446 different speed, the I/O speed is shown in Figure 9. Although the I/O speed of CoordL is far less
 447 than JOADER, the fast job must wait for a slow job causing inefficiency. JOADER can make jobs run
 448 at their own speeds. Figure 10a show the I/O speed for four jobs that start at different moments. Due
 449 to the redundant I/O work, the baseline has a more considerable I/O speed. For JOADER, the more
 450 data is shared, the smaller the I/O speed, as shown in Figure 10b.

451 A.2 Algorithm evaluation

452 In this section, we first compare ISA, the default sampling algorithm in PyTorch, with DSA to show
 453 the performance in two cases: 1) datasets overlap partially, and 2) datasets vary in size. Then we
 454 further compare the RefCnt with the generic cache policy in the above cases. We evaluate these
 455 algorithms in different cache sizes and the evaluation metric is the count of *cache misses*, indicating
 456 the number of elements not read from the cache. For n datasets D_1, D_2, \dots, D_n , the maximum count
 457 of cache misses is $|D_1| + |D_2| + \dots + |D_n|$ that all the data is read from the storage, while the
 458 minimum is $|D_1 \cup D_2 \cup \dots \cup D_n|$ that all the repeated data is hit in cache.

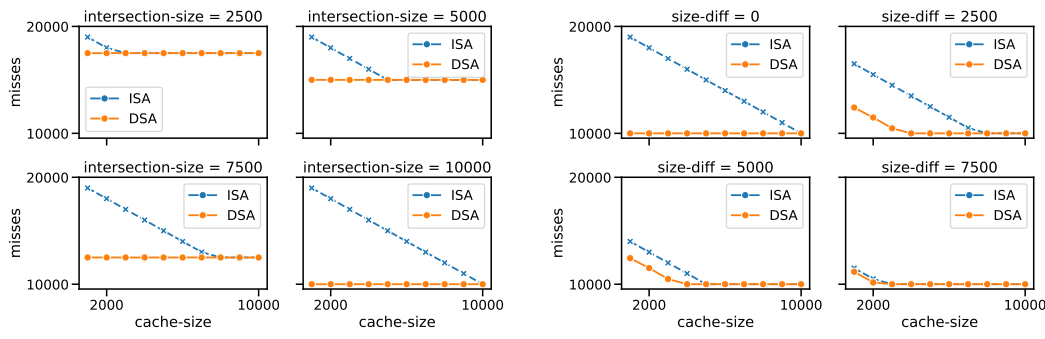
459 We construct the following cases for evaluation⁴ in Table 2:

⁴1) $D = [0, 1e4]$ means dataset D contains all number from 0 to 10000. 2) $D = \text{sample}([0, 13333], 10000)$ means sample a subset D with 10000 of size from $[0, 13333]$ uniformly at random



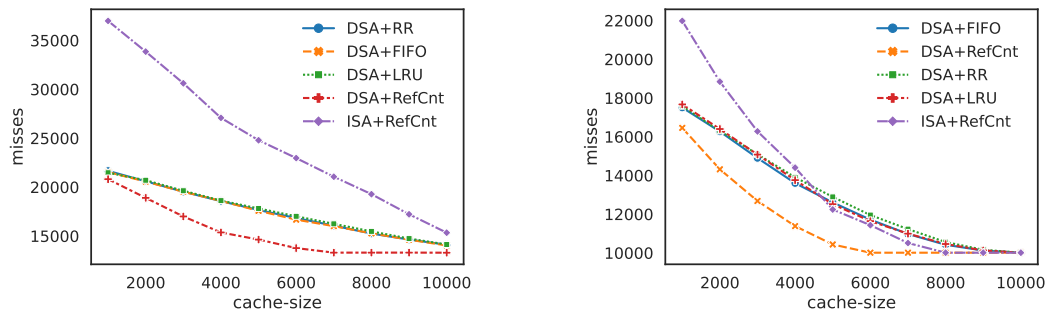
(a) When four jobs start at different moments, the I/O speeds are shown above. (b) I/O speed are shown above when intersection sizes are different.

Figure 10: Training jobs arrive at different time and have different datasets



(a) Misses count of DSA and ISA for different intersections with a cache of size from 1 to 10000. DSA can always get the minimum misses. (b) Misses count of DSA and ISA for different dataset sizes with a cache of size from 1 to 10000. DSA can reduce lots of misses.

Figure 11: Misses count of DSA and ISA in various scenes for 2 jobs



(a) Misses count for 4 jobs in an utterly random scene. RefCnt can reduce 10% misses for the same cache size. (b) Misses count for 4 jobs with datasets of different sizes. RefCnt can reduce 10% misses for the same cache size.

Figure 12: The number of misses with different cache policies applied with dependent sampling algorithm.

Table 2: Configuration of the number of jobs and datasets

Algorithm	Case	Num	way of construct dataset
Sampling Algorithm	overlapping partially	2	$D_1 = [0, 1e4]; D_2 = [0, 1e4]$
		2	$D_1 = [0, 1e4]; D_2 = [7500, 17500]$
		2	$D_1 = [0, 1e4]; D_2 = [5000, 15000]$
		2	$D_1 = [0, 1e4]; D_2 = [2500, 12500]$
		4	$D_1, D_2, D_3, D_4 = sample([0, 13333], 10000)$
	varying in size	2	$D_1 = [0, 1e4]; D_2 = [0, 1e4]$
		2	$D_1 = [0, 1e4]; D_2 = [0, 7500]$
		2	$D_1 = [0, 1e4]; D_2 = [0, 5000]$
		2	$D_1 = [0, 1e4]; D_2 = [0, 2500]$
		4	$D_1 = [0, 1e4]; D_2 = [0, 7500], D_3 = [0, 5000], D_4 = [0, 2500]$
Cache Policy	overlapping partially	4	$D_1, D_2, D_3, D_4 = sample([0, 13333], 10000)$
	varying in size	4	$D_1 = [0, 1e4]; D_2 = [0, 7500], D_3 = [0, 5000], D_4 = [0, 2500]$

460 **Dependent sampling algorithm.** In this part, we show DSA can significantly reduce the cache
461 misses in various of scenes. To simulate various scenes, we generate multiple datasets with different
462 sizes and different intersections. We use DSA to access them in shuffling order from the storage with
463 the cache of different sizes (e.g., from one-slot to holding all datasets) to evaluate the number of
464 misses, compared with the experimental results with ISA (shuffling independently).

465 We start the the evaluation in two jobs training on *dataset*₁ and *dataset*₂, respectively. We set up
466 eight different configurations for the two datasets, which can be divided into two types: the different
467 sizes of intersections and different sizes of the datasets, as shown in Table 2.

468 For different intersection sizes, Figure 11a plots the relations between cache size and the count of
469 cache misses. For ISA, the misses count is 2x higher than the size of dataset in the case of one-slot
470 cache, shows that all elements are missed and it decreases linearly as the cache size increase. Unless
471 the cache can hold all the data of intersection, the misses count could be minimized. However, DSA
472 can always get the minimum cache misses regardless of the cache size. The reason is that in the
473 selecting procedure, all jobs must select the same subset when the datasets' sizes are equal, due to the
474 conditional probability of 1. Thus, their sampling results would be the same, meaning that all jobs
475 will access the same element simultaneously when the datasets with the same size.

476 Figure 11b plots the evaluation results on for the datasets with different sizes. For one-slot cache,
477 DSA can reduce 50% cache misses compared with ISA when the dataset sizes differ by 25%. The
478 cache held 30% dataset could get the best performance for DSA while ISA needs to cache 75%
479 dataset to reach optimal. DSA needs to cache 25% data because the size difference brings some
480 uncertainty in the selecting procedure.

481 Then, we evaluate the DSA algorithm on four DNN training jobs to show the outstanding performance
482 of the DSA algorithm on multiple datasets. Multiple datasets bring more complexity to the intersection
483 as there can be $2^n - 1$ intersection sets for n sets. To simulate an utterly random scene, we
484 random sample 10,000 elements from 13,333 elements four times, making four different datasets for
485 evaluation. Compared ISA in the random scene, DSA can reduce 50% misses with a one-slot cache
486 (20000 misses out of 40000 misses), as shown in Figure 13a.

487 DSA also provides good performance for multiple jobs training on datasets with different sizes.
488 Compared to ISA in Figure 13b, one-slot cache can reduce 30% cache misses (16000 misses out of
489 22000 misses) with DSA when four jobs are training on four datasets, with a 25% different size.

490 **RefCnt cache policy.** In this part, we compare cache policy RefCnt with the generic cache policies
491 LRU, FIFO, and random replacement (RR) to show 10% misses reduction with the same cache size.
492 We evaluate these policies in above two situations of four jobs: utterly random scene and different
493 dataset sizes. Figure 12 shows the results, which can be summarized as two improvements: 1) RefCnt

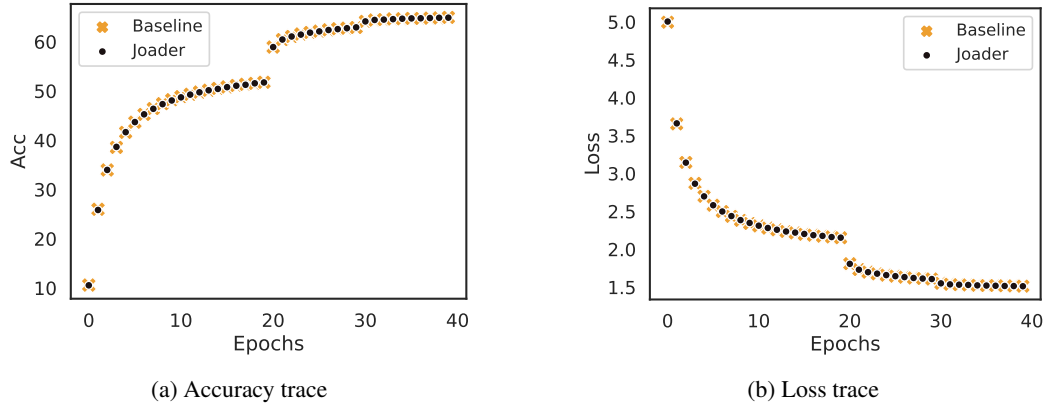


Figure 13: Trace of training ResNet18 in 40 epochs

494 cache policy can reduce 10% misses in the same cache size compared with general cache policies,
 495 and 2) RefCnt cache policy can get the best performance with only caching the 60% dataset while the
 496 generic cache policies require holding the whole dataset.

497 A.3 Correctness experiments

498 To demonstrate JOADER does not affect convergence speed w.r.t. accuracy and loss, we train ResNet18
 499 on ImageNet-1K with JOADER and PyTorch⁵. Figure 13 shows that the loss descent trajectory and
 500 accuracy trajectory are almost the same.

501 A.4 Cost of operation in sampling tree

502 **We have conducted an additional experiment by randomly inserting 128 datasets into our system and**
 503 **then evaluating the time cost of dataset operations. The size of each dataset is between 1,000,000 and**
 504 **2,000,000 elements (note that ImageNet contains 1,400,000 images). The average cost of inserting each**
 505 **dataset is 0.57 seconds. And the time cost of the operations remains nearly constant and does not change**
 506 **as the number of data increases.**

507 **Deletion and sampling are also very efficient in our algorithm. Specifically, sampling takes an average of**
 508 **0.000054 seconds for each element, and deletion takes an average of 0.00000105 seconds. We will add the**
 509 **corresponding description and experiment in our rebuttal revision.**

510 B Dependent Sampling Algorithm

511 In this section, we show the pseudo-code of the DSA, which involves two functions: selecting and
 512 sampling. The selecting algorithm is shown in Algorithm 1, and the sampling algorithm is shown in
 513 Algorithm 2.

514 C Sampling in Dependent Sampling Tree

515 The *sampling* operation involves two steps: 1) *Deciding*: choose a vertex for each job, and 2)
 516 *Sampling*: select an element in the corresponding vertex for each job randomly.

517 In *deciding*, we start the procedure from the root of the tree. In the beginning, all jobs are marked as
 518 *undecided*. In each vertex, all undecided jobs need to decide whether to select the current vertex, i.e.,
 519 the intersection I in the DSA algorithm. If the job chooses the current vertex, we will remove them
 520 from the undecided job set. Then, the remaining undecided jobs are put down to the child and make
 521 decisions recursively until the undecided job set is empty, like Algorithm 1.

⁵To activate DSA, we run 2 jobs at different moments

Algorithm 1: Selecting Procedure

Data: n jobs J_1, \dots, J_n and n datasets D_1, \dots, D_n **Result:** The subset that each job selects

```
1  $JOB \leftarrow [J_1, \dots, J_n]$ ; // Array of all jobs
2  $DS \leftarrow [D_1, \dots, D_n]$ ; // Array of all datasets sorted by their cardinality
3  $S \leftarrow \emptyset$ ; // Set of data that are not selected
4 while  $JOB \neq \emptyset$  do
5    $I \leftarrow \bigcap D_i - S, \forall D_i \in DS$ ;
6    $c \leftarrow |I|$ ;
7   for  $D, J$  in  $DS, JOB$  do
8     //  $\frac{c}{|D|-|S|}$  is the conditional probability  $\frac{|D_{i-1}|}{|D_i|}$ 
9     if  $\text{rand}(0, 1) \leq \frac{c}{|D|-|S|}$  then
10       $J$  selects  $I$ ;
11       $DS \leftarrow DS - D$ ;
12       $JOB \leftarrow JOB - J$ ;
13       $c \leftarrow |D|$ ;
14    else
15      break;
16  if There is no job selecting  $I$  then
17     $JOB[0]$  selects  $D[0] - I$ ;
18     $DS \leftarrow DS - D[0]$ ;
19     $JOB \leftarrow JOB - JOB[0]$ ;
20   $S \leftarrow S \cup I$ ;
```

Algorithm 2: Sampling Procedure

Data: n jobs J_1, \dots, J_n and k subsets S_1, \dots, S_k **Result:** The sampling results of each job

```
1  $JOB \leftarrow [\{J_i, \dots, J_j\}, \dots]$ ; // The job set where job selects the correspond subset
2  $SUBSET \leftarrow [S_1, \dots, S_k]$ ; // The subset that the correspond job set selects
3 for  $i$  in  $0..k$  do
4    $subset \leftarrow SUBSET[i]$ ;
5    $job\_set \leftarrow JOB[i]$ ;
6   Sample  $e$  uniformly at random in  $subset$ ;
7   for  $job$  in  $job\_set$  do
8      $job$  picks  $e$ ;
```

522 However, in *deciding*, there may be an intersection that contains multiple vertices of the sampling
523 tree. As shown in Figure 4c, if job J_2 with D_2 and job J_3 with D_3 do not choose the root X in the
524 first recursion, then they should make decision in the intersection $(D_2 - D_X) \cap (D_3 - D_X)$, which
525 contains the datasets in vertices A and B . In this case, we should combine these vertices into a big
526 set first and start the deciding procedure once. If someone chooses the big set, it should choose a
527 vertex from the big set randomly again.

528 After the deciding procedure, each job is assigned to a vertex randomly. Then, we need to select an
529 element in the corresponding vertex randomly. Based on the *united sampling* discussed in Section 3,
530 those jobs that are assigned to the same vertex should share the same sampling result.

531 Each data item in the dataset should be read exactly once for a job in an epoch. Intuitively, the data
532 item should be deleted from the vertex to avoid being read by the job repeatedly. This process raises
533 the issue that other jobs may need the deleted element that some have read. For example, suppose in
534 the Figure 4a, job J_1 fetches the element e_1 from vertex A while job J_2 and job J_3 fetch the element

Table 3: Notation in sampling algorithm

J_i	The i -th job
D_i	Dataset that J_i trains upon
I	Intersection of all datasets
D_{di}	$D_i - I$
$ D $	Cardinality of set D
J_i^I	Event that J_i chooses the intersection
J_i^D	Event that J_i do not choose the intersection
e_j^i	Element e_j picked by J_i in one round
S	Set of elements that are not selected

535 from other vertices. Thus, for J_1 , we should remove e_1 from vertex A , but J_2 and J_3 still need the
 536 element e_1 .

537 The compensation procedure is introduced in sampling to solve this issue, which adds the deleted
 538 element to the child vertex of the job that still needs it. To decide which vertex needs compensation,
 539 we construct a job set for each vertex that contains all the jobs that will sample in this vertex. For
 540 example, the job set of vertex A is $\{J_1, J_2, J_3\}$ and the set of vertex B is $\{J_2, J_3\}$ in Figure 4a. And
 541 after deleting some elements, we should construct a compensation set containing all the jobs that
 542 still need the deleted elements. If the compensation set includes the job set, the deleted elements are
 543 added to the corresponding vertex.

544 **Partially Sampling.** When some jobs are much faster than other jobs, we only need the sampling
 545 results for these fast jobs. The sampling procedures are the same as the above. However, the tree
 546 cannot be ordered after sampling because these datasets decrease fast and can be less than the above
 547 datasets in the tree soon. Therefore, we need to reorder the tree, just like the intersection procedure.

548 D Proof of Algorithm

549 In this section, we present the derivation and the proof for the proposed algorithm. Proposition 1
 550 shows the first principle for randomness. Table 3 presents some notations for the following proof.

551 **Proposition 1** For any job J_i , the probability of choosing any element from dataset D_i is $\frac{1}{|D_i|}$.

552 D.1 Derivation in Two-job Case

553 There are 2 jobs J_1 and J_2 which training upon datasets D_1 and D_2 , while the intersection is
 554 $I = D_1 \cap D_2$ and two difference sets are $D_{d1} = D_1 - I, D_{d2} = D_2 - I$. Their cardinalities are
 555 $|D_1|, |D_2|, |I|, |D_{d1}|$, and $|D_{d2}|$. The event that J_1 selects the intersection is J_1^I , while the event that
 556 J_1 selects the difference set is J_1^D .

557 Due to the randomness, the constraints that we need to maintain are

$$\begin{cases} p(J_1^I) = \frac{|I|}{|D_1|}, p(J_1^D) = \frac{|D_{d1}|}{|D_1|} \\ p(J_2^I) = \frac{|I|}{|D_2|}, p(J_2^D) = \frac{|D_{d2}|}{|D_2|}. \end{cases} \quad (1)$$

558 Since J_2 makes the decision conditioned on J_1 , we can the formulate the equations according to the
 559 law of total probability in below

$$\begin{cases} p(J_2^I) = p(J_2^I | J_1^I) * p(J_1^I) + p(J_2^I | J_1^D) * p(J_1^D) \\ p(J_2^D) = p(J_2^D | J_1^I) * p(J_1^I) + p(J_2^D | J_1^D) * p(J_1^D). \end{cases} \quad (2)$$

560 After *shared sampling*, our target is to maximize the probability of J_1 and J_2 both choose the
 561 intersection, which is positive correlated with $p(J_2^I|J_1^I)$. The $p(J_2^I|J_1^I)$ is

$$\begin{aligned} p(J_2^I|J_1^I) &= \frac{p(J_2^I) - p(J_2^I|J_1^D) * p(J_1^D)}{p(J_1^I)} \\ &\leq \min(1, \frac{p(J_2^I)}{p(J_1^I)}) \\ &\leq \min(1, \frac{|D_1|}{|D_2|}). \end{aligned} \quad (3)$$

562 Therefore, we can get the conditional probability in two cases

$$p(J_2^I|J_1^I) = \begin{cases} 1, & \text{if } |D_1| > |D_2| \\ \frac{|D_1|}{|D_2|}, & \text{if } |D_1| \leq |D_2|. \end{cases} \quad (4)$$

563 When $|D_1| \leq |D_2|$, the other conditional probabilities of J_2 conditioned on J_1 can be derived as
 564 follows

$$\begin{cases} p(J_2^I|J_1^D) = 0 \\ p(J_2^D|J_1^D) = 1 \\ p(J_2^D|J_1^I) = 1 - \frac{|D_1|}{|D_2|}. \end{cases} \quad (5)$$

565 In *dependent selecting* procedure, the probability of both choosing intersection set is $p(J_1^I) * p(J_2^I|J_1^I)$.
 566 With *data partition*, we can get

$$p(e_j^1 = e_i^2) = p(J_1^I) * p(J_2^I|J_1^I) \leq \frac{|I|}{\max(|D_2|, |D_1|)}, \quad (6)$$

567 while e_j^1, e_i^2 is the sampling results of J_1, J_2 in one round.

568 We then prove the the proposed dependent algorithm gives the optimal solution for the scenario of
 569 two-jobs.

570 **Proof 1** Assume that the distribution of the event exists, and the probability of $p(e_j^1 = e_i^2)$ is unknown.
 571 When e_j^1 and e_i^2 are equal, they must belong to the intersection set D_i . Therefore, $p(e_j^1 = e_i^2)$ is equal
 572 to $p(e_j^1 = e_i^2, e_j^1 \in D_i, e_i^2 \in D_i)$. Due to the definition of joint probability, we can get

$$p(e_j^1 = e_i^2) \leq p(e_j^1 \in D_i), \quad \text{and} \quad p(e_j^1 = e_i^2) \leq p(e_i^2 \in D_i).$$

573 Due to the constraint of randomness, e_j^1 and e_i^2 are both sampled under the uniform distribution, so
 574 $p(e_j^1 \in D_i) = \frac{|D_i|}{|D_1|}$ and $p(e_i^2 \in D_i) = \frac{|D_i|}{|D_2|}$. Then $p(e_j^1 = e_i^2) \leq \frac{|D_i|}{|D_1|}$ and $p(e_j^1 = e_i^2) \leq \frac{|D_i|}{|D_2|}$ both
 575 hold. Thus, for any distribution of the event, it holds that

$$p(e_j^1 = e_i^2) \leq \frac{|D_i|}{\max(|D_2|, |D_1|)}.$$

576 D.2 Derivation in N-job Case

577 Assume there are n jobs $\{J_1, \dots, J_n\}$ and n datasets $\{D_1, \dots, D_n\}$, while the intersection $I =$
 578 $\bigcap_{i=1}^n D_i$. For the k -th job where $k > 1$, the probability of the job J_k chooses the intersection is

$$p(J_k^I) = p(J_k^I|J_1^I, J_2^I, \dots, J_{k-1}^I) * p(J_1^I, J_2^I, \dots, J_{k-1}^I) + x, \quad (7)$$

579 where x is the sum of probabilities of J_k chooses intersection in other conditions. With the similar
 580 process in Proof 1, we can get the result for k -jobs, where the maximum probability of they all choose

581 intersection is also $p(J_1^I \dots J_k^I) \leq \frac{|I|}{\max(D_1, \dots, D_k)}$. Then, we can get

$$\begin{aligned}
p(J_k^I | J_1^I, J_2^I, \dots, J_{k-1}^I) &\leq \min\left(\frac{p(J_k^I)}{p(J_1^I, J_2^I, \dots, J_{k-1}^I)}, 1\right) \\
&\leq \min\left(\frac{\frac{|I|}{|D_k|}}{\frac{|I|}{\max(D_1, \dots, D_{k-1})}}, 1\right) \\
&\leq \min\left(\frac{\max(D_1, \dots, D_{k-1})}{|D_k|}, 1\right).
\end{aligned} \tag{8}$$

582 Assume that the datasets are sorted in ascending order w.r.t. their cardinalities, then the formula can
583 be written as for the maximum probability

$$p(J_k^I | J_1^I, J_2^I, \dots, J_{k-1}^I) = \frac{|D_{k-1}|}{|D_k|}. \tag{9}$$

584 The probability of n jobs all choose intersection I is

$$p(J_1^I, J_2^I, \dots, J_n^I) = \frac{|I|}{|D_1|} * \frac{|D_1|}{|D_2|} * \dots * \frac{|D_{n-1}|}{|D_n|} \leq \frac{|I|}{|D_n|}, \tag{10}$$

585 which is the theoretical maximum according to a similar Proof 1.

586 D.3 Proof of Randomness

587 We prove the Proposition 1 for Algorithm 1. Although multiple jobs share the sampling results in the
588 sampling procedure, the random sampling is not changed. Therefore, to prove Proposition 1, we only
589 need to prove the subsets are randomly selected for each job. The Lemma is stated as follow

590 **Lemma 1** *In each loop of selecting procedure, the intersection I is selected with the probability of*
591 $\frac{|I|}{|D_i|}$ *for every job J_i .*

592 In the line 8-12 of Algorithm 1, the job J_i chooses the intersection only if the previous job has chosen
593 the intersection I , that the probability is

$$\begin{aligned}
p(J_i^I) &= p(J_i^I | J_{i-1}^I, \dots) * p(J_{i-1}^I, \dots) \\
&= \frac{|I|}{|D_i| - |S|}.
\end{aligned} \tag{11}$$

594 Therefore, the Lemma 1 is satisfied only if the probability of entering this loop is $\frac{|D_i| - |S|}{|D_i|}$, that is

595 **Lemma 2** *For each loop, the set of elements that are not selected is S , and the probability of entering*
596 *this loop is $\frac{|D_i| - |S|}{|D_i|}$ for job J_i .*

597 **Proof 2** *We prove the Lemma 2 by induction on loop index k .*

598 **Base case.** *Show Lemma 2 holds for the first loop.*

599 *The first loop ($k = 1$) must be entered, while S is an empty set and $\frac{|D_1| - |S|}{|D_1|} = 1$. Thus Lemma 2 is*
600 *satisfied.*

601 **Induction step.** *Show that for any $k \geq 1$, if the k -th recursion holds Lemma 2, then $(k + 1)$ -th loop*
602 *also holds.*

603 *The probability of entering the k -th loop is $\frac{|D_k| - |S|}{|D_k|}$, and if job J_i does not select the intersection in*

604 the k -th loop, then it entering the $(k + 1)$ -th loop. The probability of J_i does not select intersection is

$$\begin{aligned}
p(J_i^D) &= p(J_k^D) + p(J_{k+1}^D) + \dots + p(J_i^D) \\
&= p(J_k^D) + p(J_{k+1}^D | J_k^I) * p(J_k^I) + \dots + p(J_i^D | J_{i-1}^I) * p(J_{i-1}^I) \\
&= p(J_k^D) + (1 - p(J_{k+1}^I | J_k^I)) * p(J_k^I) + \dots + (1 - p(J_i^I | J_{i-1}^I)) * p(J_{i-1}^I) \\
&= p(J_k^D) + p(J_k^I) - p(J_{k+1}^I) + p(J_{k+1}^I) - \dots - p(J_i^I) \\
&= 1 - p(J_i^I) \\
&= 1 - \frac{|I|}{|D_k| - |S|} * \frac{|D_k| - |S|}{|D_{k+1}| - |S|} * \frac{|D_{i-1}| - |S|}{|D_i| - |S|} \\
&= \frac{|D_i| - |S| - |I|}{|D_i| - |S|}
\end{aligned} \tag{12}$$

605 where J_k is the first job in job List JOB, then we can get the probability of entering the next loop is
606 $\frac{|D_i| - |S| - |I|}{|D_i|}$ while new S is $S \cup I$, which also satisfies the Lemma 2.

607 **Conclusion.** Since both the base case and the inductive step have been proved as true, the Lemma 2
608 holds for every loop by mathematical induction.

609 E System Implementation

610 In this section, we describe the implementation of JOADER, which is a data loading management
611 system for multiple DNN training jobs. The proposed dependent sampling algorithm applied with the
612 dependent sampling tree is the core in JOADER. In details, we describe the system overview, three
613 main components, and the adaption for distributed DNN training.

614 E.1 Overview

615 **API.** JOADER API allows users to do two things: 1) creating and altering datasets, and 2) registering
616 DNN training jobs. Before training, the user needs to create the dataset with a name in JOADER first.
617 Then, to execute the training job on the dataset, the job should register itself to JOADER with the
618 dataset's name. If multiple jobs are training on that dataset, JOADER will attach the dataset to the
619 jobs for sharing reading and preprocessing. Each job in an epoch will be assigned to a unique id in
620 JOADER, and JOADER will dispatch data according to the job id.

621 **Architecture.** JOADER includes a frontend and a backend. The backend is implemented in Rust [2]
622 that accesses and processes data efficiently, while the frontend is implemented in Python to provide a
623 user-friendly API. They communicate with each other in RPC. Figure 14 shows the architecture of
624 JOADER, which consists of three components: Sampler, Loader, and Cache.

625 E.2 Sampler and Loader

626 The Sampler is responsible for data sampling in JOADER. An instance of the Sampler consists of a
627 sampling tree for a dataset. The datasets are organized in tables like the relational database. Each
628 data tuple has a unique integer id (i.e., primary key) to identify itself, consisting of multiple elements,
629 e.g., image, label, and the bounding boxes in ImageNet.

630 Each job registered to the sampler needs to specify the job's name, dataset's name, names of the
631 needed columns in the table, and a predicate for filtering. Then, the sampler collects the ids of data
632 tuples that meet the requirements and inserts the ids into the sampling tree corresponding to the
633 dataset. Each sampling tree should pick the element uniformly at random in the DS algorithm. These
634 elements (id of data tuple) will be transformed into the data requests according to the columns of
635 each job needed. If the elements and the columns are the same, these data requests will be merged
636 into one request.

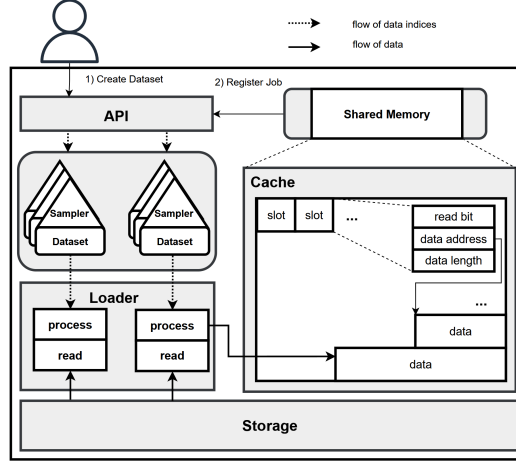


Figure 14: The architecture of JOADER.

637 The Loader is responsible for reading and processing data w.r.t. the data request from Sampler. In
 638 JOADER, each loader corresponds to one type of storage, e.g., POSIX file system, key-value database,
 639 or distributed storage. Loaders encapsulate different storages and provide a unified interface to
 640 Sampler to present good scalability and compatibility.

641 E.3 Specific Cache Implementation

642 The Cache is used to store and share data between JOADER and the DNN training jobs, which is
 643 implemented by shared memory to avoid the cost of memory copying, serialization, and network. To
 644 manage data and reduce data race, the layout of cache is slotted, making multiple slots in the head of
 645 the cache, as shown in Figure 14. Each element in the cache is managed by a slot containing the start
 646 address, length, and read bit of the data. The read bit is used to determine whether the data has been
 647 consumed.

648 When Loader tries to load data, it needs to require slots and a contiguous block memory for storing
 649 this data. Then, the Loader dispatches the slot ids to training jobs that need the data. These jobs will
 650 access the data soon and set the read bit. When there is not enough memory, the reclamation program
 651 should be triggered. Data can be classified as consumed and unconsumed in the cache according to
 652 the read bit. Because those unconsumed data were sampled earlier than the data that Loader is trying
 653 to require memory, they will be accessed earlier than the new data and cannot be evicted.

654 We can only reclaim those consumed elements, which are prioritized in different levels. For example,
 655 there are two training jobs J_1 and J_2 all want to read four data $\{d_1, d_2, d_3, d_4\}$. In the first round
 656 of sampling, d_1 is consumed by both jobs. In the second round, job J_1 consumed d_2 while job J_2
 657 consumed d_3 . Then the memory reclamation program is triggered. Considering the constraint that
 658 each job should traverse the dataset once. In this time, data d_1, d_2, d_3 are in cache while data d_4 is no
 659 longer needed anymore while d_2, d_3 are still needed. Therefore, d_1 will be prioritized for eviction.

660 In JOADER, we reclaim data according to the number of reference of data, as shown in section 3.2.

661 E.4

662 In our Joader, the time complexity of dataset operations, i.e., sampling, deletion, and insertion, are
 663 $O(1)$, $O(1)$, and $O(|D|)$, respectively, for each dataset. Moreover, the dependent sampling tree
 664 manages the **index** of each input rather than the input itself, and thus executing the operations is
 665 extremely fast. In our implementation, we use a bitmap to represent the dataset in our dependent
 666 sampling tree. The bitmap is very efficient for these operations. In practice, the operations only count
 667 up to a minor fraction of time consumption, which takes a tiny proportion of a total training epoch.

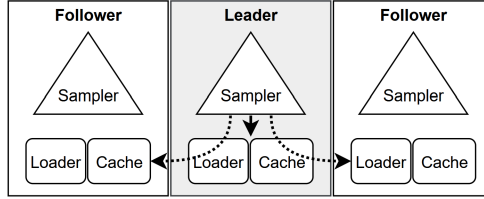


Figure 15: JOADER for distributed DNN training, where the leader JOADER is responsible for the sampling globally.

668 To further answer the question, we have conducted an additional experiment by randomly inserting
 669 128 datasets into our system and then evaluating the time cost of dataset operations. The size of each
 670 dataset is between 1,000,000 and 2,000,000 elements (note that ImageNet contains 1,400,000 images).
 671 The average cost of inserting each dataset is 0.57 seconds. And the time cost of the operations remains
 672 nearly constant and does not change as the number of data increases.

673 Deletion and sampling are also very efficient in our algorithm. Specifically, sampling takes an average
 674 of 0.000054 seconds for each element, and deletion takes an average of 0.00000105 seconds.

675 E.5 Towards Distributed Training

676 In model parallelism, the model is segmented into different parts that can run concurrently in different
 677 nodes. Only one node is responsible for data preparation. Therefore, JOADER need not change
 678 anything to apply to the distributed training in model parallelism.

679 In data parallelism, the dataset is divided into several partitions, where the number of partitions is
 680 equal to the total number of available nodes in the cluster. The model is replicated to the worker
 681 nodes. Each worker operates on its subset of the dataset to train the model locally. For each worker
 682 node, we need to set up a JOADER for sharing data preparation to multiple DNN training jobs in this
 683 node. However, the sampling work should be done globally to avoid redundant sampling.

684 To solve the above issue, we set multiple JOADER for multiple nodes, while only one leader JOADER
 685 can sample. The leader needs to dispatch the sampling results to the followers, as shown in Figure 15.
 686 During training, each sub-process of the distributed training job in data parallelism needs to register
 687 itself to the local JOADER in the same worker node with the same name and the assigned id.
 688 Meanwhile, the local JOADER should register a sub-sampler in leader JOADER. The local JOADER
 689 should keep fetching sampling results from the leader JOADER and load them to the local cache.
 690 By doing this, only leader JOADER is responsible for sampling and dispatching sampling results to
 691 followers. Notice that the cache is also distributed in different worker nodes, storing more data than a
 692 single cache.

693 The same data is sent to the same worker node when the leader dispatches sampling results. PyTorch
 694 uses the hash partition algorithm to dispatch data. For example, suppose a worker node array $nodes =$
 695 $[node_1, node_2, \dots, node_n]$. The data with integer id is sent to $nodes[id\%n]$. However, things will
 696 be more complicated when dealing with multiple training jobs. For example, two training jobs are
 697 on the worker node arrays $[node_1, node_2, \dots, node_{n-1}]$ and $[node_2, node_3, \dots, node_n]$, respectively.
 698 Due to the offset between the two arrays, the same data cannot be sent to the same worker node for
 699 the two jobs by one hash partition.

700 In JOADER, we solve it by hashing the data id twice. In the first hash, we try to dispatch the data to
 701 all the worker node in cluster. Therefore, the same data can be sent to the same node. However, some
 702 training jobs are not training on some worker nodes and the data can not be sent to these nodes for
 703 these jobs. Therefore, in the second hash, for the data dispatched wrongly, we will try to dispatch it
 704 to the nodes locally w.r.t. the configuration of each job.