

# Ferret: An Efficient Online Continual Learning Framework under Varying Memory Constraints

Anonymous CVPR submission

Paper ID 715

## Abstract

001 *In the realm of high-frequency data streams, achieving*  
002 *real-time learning within varying memory constraints is*  
003 *paramount. This paper presents Ferret, a comprehensive*  
004 *framework designed to enhance online accuracy of Online*  
005 *Continual Learning (OCL) algorithms while dynamically*  
006 *adapting to varying memory budgets. Ferret employs a fine-*  
007 *grained pipeline parallelism strategy combined with an it-*  
008 *erative gradient compensation algorithm, ensuring seam-*  
009 *less handling of high-frequency data with minimal latency,*  
010 *and effectively counteracting the challenge of stale gradi-*  
011 *ents in parallel training. To adapt to varying memory bud-*  
012 *gets, its automated model partitioning and pipeline plan-*  
013 *ning optimizes performance regardless of memory limita-*  
014 *tions. Extensive experiments across 20 benchmarks and 5*  
015 *integrated OCL algorithms show Ferret’s remarkable ef-*  
016 *iciency, achieving up to  $3.7\times$  lower memory overhead*  
017 *to reach the same online accuracy compared to compet-*  
018 *ing methods. Furthermore, Ferret consistently outperforms*  
019 *these methods across diverse memory budgets, underscor-*  
020 *ing its superior adaptability. These findings position Ferret*  
021 *as a premier solution for efficient and adaptive OCL frame-*  
022 *work in real-time environments.*

## 023 1. Introduction

024 Data is crucial for Machine Learning (ML), forming the ba-  
025 sis for algorithms and models [29, 52, 74]. In real-world  
026 applications [23, 70], data arrives in high-frequency streams  
027 with varying distributions [51, 77]. This makes data time-  
028 sensitive and short-lived [11, 47, 75], rendering offline-  
029 trained models based on historical data ineffective for future  
030 data of unknown distribution [60]. Thus, the significance of  
031 *Online Continual Learning* (OCL) is growing [6, 39, 76],  
032 as it enables learning over data streams to adapt to dynamic  
033 data distributions in real-time.

034 In the literature, OCL tackles two main challenges: 1)  
035 mitigating catastrophic forgetting [50], where the model re-

tains previously learned knowledge while acquiring new in- 036  
formation (e.g., regularization-based [2, 9, 10, 19], replay- 037  
based [12, 38, 67], sampling-based [3, 4, 82], others [25, 038  
64], etc.), and 2) enhancing rapid adaptation [11, 47], which 039  
involves swiftly adjusting to new data or tasks (e.g., latency- 040  
oriented [28, 69], buffering [53, 81], others [26, 68], etc.). 041  
In general, the increasing demand for resource-limited sys- 042  
tems that can seamlessly integrate new information with 043  
minimal latency has driven the popularity of OCL [31] 044  
Therefore, this paper explores the challenge of rapid adap- 045  
tation under varying memory constraints in OCL. 046

To effectively address the above OCL challenge, it is 047  
essential to explore solutions beyond mentioned algorithmic 048  
improvements by also optimizing the underlying frame- 049  
work. An efficient OCL framework must prioritize both 050  
processing speed and memory management under the lim- 051  
ited memory capacity so that it can efficiently handle un- 052  
limited data streams with dynamic data distributions for 053  
increased *online accuracy* [11] (i.e., a metric measuring 054  
real-time accuracy for continuous new data predictions). 055  
Specifically, the framework should quickly process incom- 056  
ing data to extract valuable insights and make informed de- 057  
cisions [55, 73] by minimizing both the latency from data 058  
receipt to its initial processing and the time taken for the 059  
learning process itself. Additionally, the framework is not 060  
only expected to operate within a predetermined memory 061  
allotment but also to demonstrate scalability across diverse 062  
memory capacities [18, 56]. This duality ensures that the 063  
framework remains efficient regardless of the memory re- 064  
sources available, thereby maintaining consistent perfor- 065  
mance in dynamic environments. 066

Numerous ML frameworks have been proposed that offer 067  
innovative approaches to scalable and flexible ML de- 068  
velopment [5, 17, 30, 35, 36, 41, 54, 83, 85]. For instance, 069  
Ray [54] facilitates distributed computing on any scale, 070  
while Pytorch [5] excels in dynamic computation graphs for 071  
model training. Despite their advancements, these frame- 072  
works often do not specifically address the unique require- 073  
ments of learning over streaming data [28], which is a 074  
key focus of OCL. Recently, there are some frameworks 075

dedicated to OCL by prioritizing real-time data processing [43, 45, 79]. Nevertheless, they either lack general applicability or fail to balance processing speed with consumed memory, leading to reduced online accuracy and low memory scalability, underscoring the need for innovative solutions in this domain.

In this work, we propose an OCL framework named Ferret, designed to achieve **eFficiEnt** pipeline **leaRning** over **fRequEnt** data **sTreams** for enhanced online accuracy across memory constraints. Ferret comprises a fine-grained pipeline parallelism component with an iterative gradient compensation algorithm and a model partitioning and pipeline planning component. Firstly, to facilitate rapid adaptation over frequent streaming data for higher online accuracy, Ferret employs a fine-grained pipeline parallel strategy, allowing precise control over each pipeline stage for seamless data management. Additionally, to mitigate the impact of stale gradients in parallel processing, Ferret integrates a novel iterative gradient compensation algorithm. Secondly, to guide the selection of optimal model partition schemes and pipeline configurations under given memory budgets, Ferret solves the involved multivariate optimization problem through a bi-level optimization algorithm.

Our contributions can be outlined as follows:

- We propose a framework named Ferret for boosting the online accuracy of OCL algorithms under memory constraints. To the best of our knowledge, this is the first work focusing on enhancing OCL by employing pipeline parallelism and scheduling.
- To process high-frequent data streams without delay, Ferret employs a fine-grained pipeline parallelism strategy, enabling interleaved processing of incoming streaming data. Furthermore, Ferret utilizes an iterative gradient compensation algorithm to efficiently mitigate the effects of stale gradients across different pipeline stages, preventing performance degradation.
- We derive the optimal parameters for automatic model partitioning and pipeline planning by mapping the involved multi-variable optimization problem into a bi-level optimization problem.
- Extensive experiments on 20 benchmarks demonstrate that our proposed framework consistently enables more efficient OCL within given memory budgets. The code is open-sourced for reproduction.

## 2. Related Work

The current OCL research focuses on two areas: mitigating catastrophic forgetting and enhancing rapid adaptation.

**Mitigating catastrophic forgetting:** Catastrophic forgetting, often quantified by the test accuracy [28, 45, 47], poses a significant barrier to the efficacy of OCL in dynamic environments, where the ability to preserve historical information is crucial. Multiple directions have emerged to

reduce catastrophic forgetting, including: 1) regularization-based techniques [2, 9, 10, 19] impose constraints on weight updates to preserve important parameters that are crucial for past tasks. 2) replay-based techniques [12, 38, 67] help the model to rehearse old knowledge alongside new information by maintaining a memory of previous data. 3) sampling-based techniques [3, 4, 82] enhance the efficiency of replay mechanisms by selectively choosing the most relevant data samples for rehearsal. 4) other techniques [25, 64] focus on various novel approaches, such as modular networks and dynamically allocated resources, to protect previously learned information from being overwritten.

**Enhancing rapid adaptation:** Rapid adaptation is in scenarios where immediate processing of incoming data is required [11, 47], which is often quantified by the online accuracy [11] defined as  $oacc_{\mathcal{A}}(t) = \sum_{i=1}^t acc(y^i, \hat{y}^i)/t$ . Strategies developed to enhance rapid adaptation include: 1) latency-oriented techniques [28, 69] iteratively generate predictions and update model parameters immediately upon the arrival of streaming data by discarding data that cannot be processed in time. 2) buffering-oriented techniques [53, 81] buffers and samples incoming data streams and apply periodic batch-training [61]. 3) other techniques [26, 68] introduce novel methods like adapting model structures in response to new tasks and learning how to learn efficiently, to adapt to dynamic data distributions rapidly.

## 3. Motivation

To effectively navigate the challenges posed by OCL, it is crucial to expand our approach beyond merely refining mentioned OCL algorithms, by also enhancing the underlying ML framework to adaptively balance processing speed with efficient memory management under memory constraints. Particularly, boosting processing speed is essentially reducing data processing time, which can be represented as  $t^l + F/R_h P_h$ , where  $t^l$  denotes the latency from data arrival to processing,  $F$  denotes the required floating point operations (FLOPS) by the underlying OCL algorithm,  $R_h$  and  $P_h$  denote the hardware utilization rate and the theoretical floating point operations per second (FLOPs) of the hardware, respectively. Clearly, only  $t^l$  and  $R_h$  are optimizable by the framework.

Existing ML frameworks mainly focus on: 1) distributed and parallel computing [41, 54, 85], 2) Optimized model training and deployment [5, 35, 83], and 3) others including security [17, 36] and debugging [30]. These frameworks facilitate scalable and flexible ML, yet they rarely tackle the challenges of managing streaming data. Regrettably, the few ML frameworks designed for OCL [43, 45, 79] either lack general applicability or fail to concurrently optimize  $t^l$  and  $R_h$  within memory limitations. For instance, Kraken [79] is tailored for recommendation systems. Conversely, while Camel [45] and LifeLearner [43] boost  $R_h$

via buffering and sampling, they also raise  $t^l$  and memory usage, reducing online accuracy.

Pipeline parallelism can naturally process sequential streaming data while utilizing batch training. This motivates us to incorporate pipeline parallelism into OCL to simultaneously minimize  $t^l$  and maximize  $R_h$  under a given memory budget, thereby boosting online accuracy. We achieve this balance through refined scheduling strategies and better hardware integration, ensuring optimal resource utilization within the constraints of memory budgets.

## 4. Problem Formulation

In this section, we define the problem we aim to address. Note that the notations used throughout this paper are defined in Sec. 9 in the appendix.

Consider a general learning problem defined over a feature space  $\mathcal{X}$  and a label space  $\mathcal{Y}$  that aims to minimize a loss function  $\mathcal{L}(D^t; \theta)$  where data  $D^t = (\mathbf{x}^t, \mathbf{y}^t) \in \mathcal{X} \times \mathcal{Y}$  arrives at timestamp  $t$ . Our objective is to rapidly derive an updated model  $\theta^t$  with  $D^t$  and  $\theta^{t-1}$  under a given memory constraint  $M$ , so that the online accuracy of  $\theta^t$  is high. Unlike updating a model offline with a pre-collected dataset,  $D^t$  will be discarded after updating  $\theta^{t-1}$  in OCL.

Directly optimizing online accuracy in our objective during runtime is hard, as the online accuracy is only calculable after obtaining labels of incoming data. Instead, we measure the volume of data values learned by the model as a proxy to estimate and optimize online accuracy. Formally, assuming  $D^t$  has an initial data value of  $V_{D^t}$  and its data value declines as a time-dependent exponential decay function [75], we define the Adaptation Rate as follows.

**Definition 4.1** (Adaptation Rate of A OCL framework). Consider a OCL framework  $\mathcal{A}$  receives a data  $D^t = (\mathbf{x}^t, \mathbf{y}^t)$  at timestamp  $t$  that has an initial data value of  $V_{D^t}$ , and updates a model  $\theta^{t-1}$  in the hypothesis space  $\Theta$  at timestamp  $t + r_{\mathcal{A}}^t$  ( $r_{\mathcal{A}}^t = +\infty$  if  $D^t$  is discarded). Let the data value of  $D^t$  decline as a time-dependent exponential decay function, and new data  $D^t$  constantly arrives until  $t = T$ . The Adaptation Rate of  $\mathcal{A}$  is defined as

$$R_{\mathcal{A}}^T = \frac{\sum_{t=0}^T e^{-cr_{\mathcal{A}}^t} V_{D^t}}{T}, \quad (1)$$

where the constant  $c$  describes the reduction rate of  $V_{D^t}$ .

With Def. 4.1, our objective can be formulated as

$$\max_{\mathcal{A}} R_{\mathcal{A}}^T \text{ s.t. } \mathcal{M}_{\mathcal{A}} \leq M, \quad (2)$$

where  $\mathcal{M}_{\mathcal{A}}$  is the memory footprint of  $\mathcal{A}$  during training.

## 5. Methodology

The workflow of Ferret is shown in Fig. 1, comprising a fine-grained pipeline parallelism component (A), followed by a model partitioning and pipeline planning component

(B). In A, the model is trained using Ferret’s fine-grained pipeline parallelism to manage high-frequent data streams with minimal latency. Given the high degree of parallelism within the system, gradient staleness can become significant and variable, potentially causing severe model degradation. To mitigate this issue, an iterative gradient compensation algorithm is applied prior to model updating. In B, the model is profiled to optimize Eq. 2, determining the optimal model partition scheme and pipeline configuration.

### 5.1. Fine-grained Pipeline Parallelism

#### 5.1.1 Architectural design

Ferret utilizes an asynchronous pipeline parallelism strategy with 1F1B scheduling to process streaming data immediately upon arrival. To efficiently handle high-frequency data streams without delay, it is imperative that  $t^f + t^b$  is minimized. However,  $t^f + t^b$  is inherently lower bounded by  $(\max_i \hat{t}_i^f + \max_i \hat{t}_i^b)$ , indicating that some of the data must be discarded if  $t^d$  is less than this lower bound. To prevent the loss of data, Ferret enhances system throughput by deploying  $N \leq \lceil (t^f + t^b)/t^d \rceil$  workers, each performing pipeline parallelism concurrently over interleaved data streams. Specifically, the  $i$ -th data is processed by the  $n$ -th worker if and only if  $i \equiv c_n^d \pmod{\lceil (t^f + t^b)/t^d \rceil}$ . This strategy, while effective in reducing latency, significantly increases memory usage. Therefore, Ferret balances the trade-offs between  $\mathcal{R}$  and  $\mathcal{M}$  by collectively employing four techniques: activation recomputation [13], gradient accumulation [58], back-propagation omission and worker removal, allowing precise control over each pipeline stage for seamless data management.

**T1. Activation Recomputation:** Activation recomputation exchanges additional computational overhead for reduced memory usage, as Fig. 1a illustrated. In Ferret, a binary indicator  $c_n^r$  within configuration  $C$  denotes whether activation recomputation is enabled for the  $n$ -th worker. When activation recomputation is enabled (*i.e.*,  $c_n^r = 1$ ), an additional forward pass is executed prior to the backward pass, effectively managing memory consumption at the expense of increased computational load.

**T2. Gradient Accumulation:** Gradient Accumulation allows multiple forward and backward passes to accumulate gradients before model updating, thereby decreasing the frequency of parameter updates, as Fig. 1b depicted. In Ferret, the parameter  $c_{n,j}^a$  in configuration  $C$  defaults to 1, indicating the number of gradient accumulation steps before model updating for the  $j$ -th stage in the  $n$ -th worker. By utilizing gradient accumulation, the  $j$ -th stage in the  $n$ -th worker only stores  $(1 + \lceil (P - j - 1)/c_{n,j}^a \rceil)$ , instead of  $(P - j)$ , models, thereby optimizing memory usage.

**T3. Back-propagation Omission:** To further reduce memory usage, back-propagation omission skips all backward passes that depend on previous model parameters, as

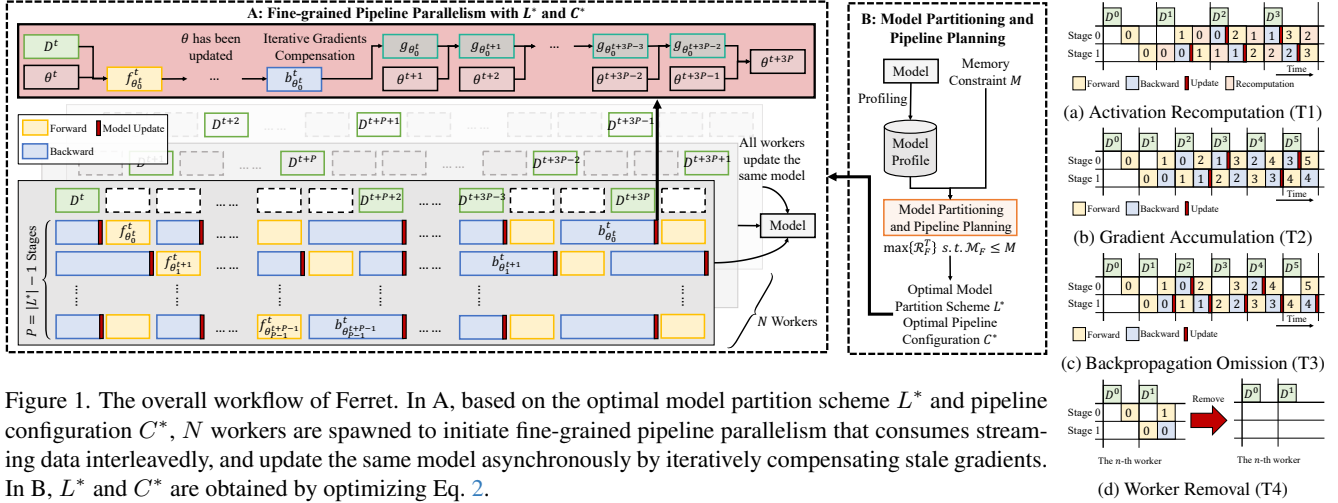


Figure 1. The overall workflow of Ferret. In A, based on the optimal model partition scheme  $L^*$  and pipeline configuration  $C^*$ ,  $N$  workers are spawned to initiate fine-grained pipeline parallelism that consumes streaming data interleavedly, and update the same model asynchronously by iteratively compensating stale gradients. In B,  $L^*$  and  $C^*$  are obtained by optimizing Eq. 2.

279 Fig. 1c illustrated. In Ferret, the parameter  $c_{n,j}^o$  in con-  
 280 figuration  $C$  defaults to 0, indicating the number of back-  
 281 propagation omission steps for the  $j$ -th stage in the  $n$ -th  
 282 worker. This approach reduces memory overhead by eliminat-  
 283 ing the need to store multiple versions of models.

284 **T4. Worker Removal:** Spawning  $N$  workers increases  
 285 the system throughput but also linearly increases the mem-  
 286 ory footprint. When resources are highly constrained, the  
 287  $n$ -th worker can be shut down and removed to reduce the  
 288 memory overhead by setting  $c_n^d = -1$  in configuration  $C$ .

289 Finally, assume the initial data value of any data is  $V_D$ ,  
 290 and the retained value of the data when updating a sub-  
 291 set of model parameters is proportional to the size of the  
 292 subset model parameters, given  $L$  and  $C$ ,  $\mathcal{R}$  and  $\mathcal{M}$  of  
 293 the fine-grained pipeline parallelism strategy can be respec-  
 294 tively formulated as

$$295 \mathcal{R}_F^T = \sum_{n=1}^{N-1} \sum_{i=0}^{P-1} \frac{|w_i|}{\sum_{j=0}^{P-1} (|w_j|) c_{n,i}^a} \frac{1}{c_{n,i}^a} \sum_{j=0}^{c_{n,i}^a-1} A_{i,j}, \text{ where } A_{i,j} =$$

$$296 \frac{e^{-c((P+j)t^f + (P-i+j)t^b + c_n^r(P-i+j)t^f)} V_D}{LCM(\{c_{n,k}^o + 1 | k \in [i, P-1]\})(t^f + t^b + c_n^r t^f)}, \quad (3)$$

$$297 \mathcal{M}_F = \sum_{n=1}^{N-1} \sum_{i=0}^{P-1} (1 + \lceil \frac{P-i-1}{c_{n,i}^a} \rceil - c_{n,i}^o) (|w_i| +$$

$$298 |a_i| - c_n^r \sum_{l=L_i+1}^{L_{i+1}-1} |\hat{a}_l|), \quad (4)$$

299 where  $LCM(\cdot)$  denotes the Least Common Multiple.

### 300 5.1.2 Iterative Gradient Compensation

301 Since fine-grained pipeline parallelism is asynchronous, the  
 302 model will be inevitably updated by stale gradients, leading  
 303 to performance degradation. Moreover, different pipeline

304 stages of the model are updated by gradients with varying  
 305 staleness. To surmount the above challenges, Ferret firstly  
 306 proposes to efficiently approximate  $\nabla \mathcal{L}(D^{t-1}; \theta^t)$  using  
 307  $\nabla \mathcal{L}(D^{t-1}; \theta^{t-1})$  by a cost-effective approximator  $A_{\mathcal{I}}(\cdot)$   
 308 based on Taylor series expansion and the Fisher informa-  
 309 tion matrix. Then, we extend this approximator to approx-  
 310 imate  $\nabla \mathcal{L}(D^{t-1}; \theta^{t+\tau-1})$  using  $\nabla \mathcal{L}(D^{t-1}; \theta^{t-1})$  by iterat-  
 311 ively applying  $A_{\mathcal{I}}(\cdot)$ .

312 **Gradients Compensation via Taylor Series Expan-**  
 313 **sion:** In prior work,  $\nabla \mathcal{L}(D^{t-1}; \theta^t)$  was naively set to  
 314  $\nabla \mathcal{L}(D^{t-1}; \theta^{t-1})$  [57, 58] and can be regarded as a zero-  
 315 order Taylor series expansion, leading to a high approxima-  
 316 tion error  $\|\nabla \mathcal{L}(D^{t-1}; \theta^t) - \nabla \mathcal{L}(D^{t-1}; \theta^{t-1})\|^2$ . To reduce  
 317 the approximation error, we expand  $\nabla \mathcal{L}(D^{t-1}; \theta^t)$  at  $\theta^{t-1}$   
 318 by a first-order Taylor series expansion as follows:

$$319 \nabla \mathcal{L}(D^{t-1}; \theta^t) \approx \nabla \mathcal{L}(D^{t-1}; \theta^{t-1}) + \mathbb{H}(\mathcal{L}(D^{t-1}; \theta^{t-1})) \odot (\theta^t - \theta^{t-1}), \quad (5)$$

320 where  $\mathbb{H}(\cdot)$  denotes the Hessian matrix of  $\cdot$ . Previous works  
 321 have revealed that the Fisher information matrix serves as  
 322 an approximation of the Hessian matrix if the loss function  
 323  $\mathcal{L}(\cdot)$  is a negative log-likelihood loss [27, 62]. Assuming  $\theta^t$   
 324 gradually converges to its optimal value  $\theta^*$  during training,  
 325 we can achieve an unbiased estimation of  $\mathbb{H}(\cdot)$  by

$$326 \epsilon_t \triangleq \mathbb{E}_{D, \theta^*} \|\mathcal{I}(\theta^t) - \mathbb{H}(\mathcal{L}(\cdot; \theta^t))\| \rightarrow 0, t \rightarrow +\infty. \quad (6)$$

327 To further mitigate space complexity,  $\mathcal{I}(\theta)$  is approximated  
 328 by its diagonal elements with a hyper-parameter  $\lambda$  to control  
 329 variance, i.e.,

$$330 \mathbb{H}(\mathcal{L}(\cdot; \theta^t)) \approx \lambda \nabla \mathcal{L}(\cdot; \theta^{t-1}) \odot \mathcal{L}(\cdot; \theta^{t-1})^\top. \quad (7)$$

331 Incorporating Eq.7 into Eq.5, we obtain:

$$332 \nabla \mathcal{L}(D^{t-1}; \theta^t) \approx A_{\mathcal{I}}(\nabla \mathcal{L}(D^{t-1}; \theta^{t-1}), \theta^t, \theta^{t-1}) = \quad (8)$$

$$333 \nabla \mathcal{L}(D^{t-1}; \theta^{t-1}) + \lambda \nabla \mathcal{L}(D^{t-1}; \theta^{t-1}) \odot \nabla \mathcal{L}(D^{t-1}; \theta^{t-1})^\top \odot \Delta \theta,$$

334 where  $A_{\mathcal{I}}(\cdot)$  serves as the approximator to compensate  $\cdot$ .

335 **Iterative Compensation:** More generally, to approx-  
 336 imate  $\nabla \mathcal{L}(D^{t-1}; \theta^{t+\tau-1})$  using  $\nabla \mathcal{L}(D^{t-1}; \theta^{t-1})$ , Ferret

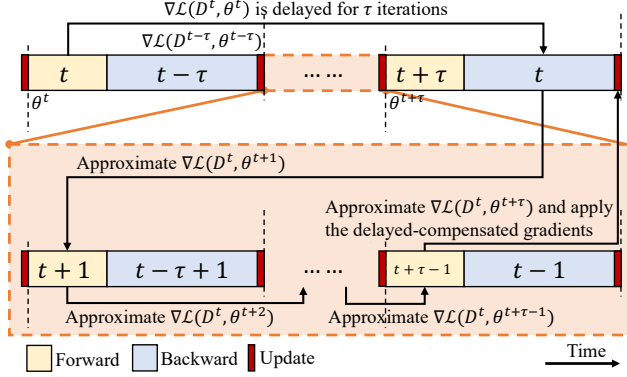


Figure 2. To adapt to different levels of staleness in fine-grained pipeline parallelism,  $\nabla\mathcal{L}(D^t, \theta^{t+\tau})$  is iteratively approximated by  $\nabla\mathcal{L}(D^t, \theta^t)$ .

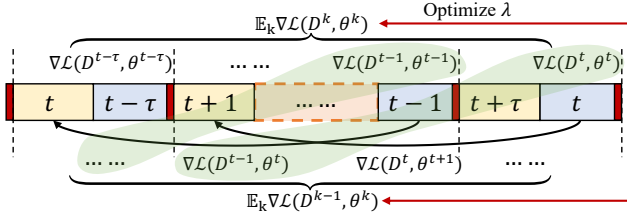


Figure 3. To further reduce approximation errors, we optimize  $\lambda$  automatically by comparing historical approximations ( $\nabla\mathcal{L}(D^t, \theta^t)$ , etc.) and observations ( $\nabla\mathcal{L}(D^{t-1}, \theta^t)$ , etc.)

337 proposes an iterative application of  $A_{\mathcal{I}}(\cdot)$ , as depicted in  
338 Fig. 2. This iterative process is defined as follows:

$$339 \nabla\mathcal{L}(D^{t-1}; \theta^{t+\tau-1}) \approx A_{\mathcal{I}}(\nabla\mathcal{L}(D^{t-1}; \theta^{t+\tau-2}), \theta^{t+\tau-1}, \theta^{t+\tau-2})$$

$$340 \approx A_{\mathcal{I}}(\dots A_{\mathcal{I}}(\nabla\mathcal{L}(D^{t-1}; \theta^{t-1}), \theta^t, \theta^{t-1}) \dots, \theta^{t+\tau-1}, \theta^{t+\tau-2}). \quad (9)$$

341 However, this iterative process introduces a cas-  
342 cade of errors, wherein the approximation error  $\|g^t -$   
343  $A_{\mathcal{I}}(g^{t-1}, \theta^t, \theta^{t-1})\|^2$  is propagated and amplified with each  
344 successive approximation. This arises because each approx-  
345 imation depends on the output of the preceding one.

346 To mitigate this problem, we propose to optimize  $\lambda$  under  
347 the mild assumption that the distributions of  $\mathbb{E}_k D^k$  and  
348  $\mathbb{E}_k D^{k+1}$  are similar, as illustrated in Fig. 3. Thus, the ob-  
349 jective of minimizing the approximation error of iterative  
350 gradient compensation can be formulated as follows:

$$351 \min_{\lambda} \mathbb{E}_k \|\nabla\mathcal{L}(D^k; \theta^k) - A_{\mathcal{I}}(\nabla\mathcal{L}(D^{k-1}, \theta^k, \theta^{k-1}))\|^2 + \nu \|\lambda\|^2$$

$$352 = \min_{\lambda} \|D - E - \lambda F\|^2 + \nu \|\lambda\|^2, \quad (10)$$

$$353 \text{ where } D = \mathbb{E}_k \nabla\mathcal{L}(D^k; \theta^k), \quad E = \mathbb{E}_k \nabla\mathcal{L}(D^{k-1}; \theta^{k-1}),$$

$$354 F = \mathbb{E}_k \nabla\mathcal{L}(D^{k-1}; \theta^{k-1}) \odot \nabla\mathcal{L}(D^{k-1}; \theta^{k-1})^\top \odot (\theta^k - \theta^{k-1}),$$

355 where  $\nu \|\lambda\|^2$  is an  $\ell_2$  regularization term to constrain the  
356 solution of  $\lambda$  for better stability. To reduce memory over-

head,  $D$  and  $E$  can be approximated by Exponential Mov- 357  
ing Average (EMA), *i.e.*, 358

$$\mathbb{E}_k \nabla\mathcal{L}(D^k; \theta^k) = \alpha \mathbb{E}_k \nabla\mathcal{L}(D^{k-1}; \theta^{k-1}) + (1-\alpha) \nabla\mathcal{L}(D^{k-1}; \theta^{k-1}), \quad (11)$$

where  $\alpha$  is the EMA coefficient. Hence, we have 360

$$D - E = (1 - \alpha)(\nabla\mathcal{L}(D^{k-1}; \theta^{k-1}) - \mathbb{E}_k \nabla\mathcal{L}(D^{k-1}; \theta^{k-1})). \quad (12)$$

**Convergence:** Similar to the analyses in [84], our iter- 362  
ative gradient compensation algorithm yields convergence 363  
rates of  $\mathcal{O}(V_1^2 \tau / T)$  and  $\mathcal{O}(V_2 / \sqrt{T})$  for convex and non- 364  
convex case, respectively. Here,  $V_1$  and  $V_2$  represent the 365  
upper-bound of the  $\|\cdot\|^2$  norm and the variance of the delay- 366  
compensated gradient  $A_{\mathcal{I}}(\cdot)$ , accordingly. Compared to the 367  
work in [84], Ferret fixes  $\tau$  to 1, and minimizes  $V_1$  and  $V_2$  368  
by Eq. 10, boosting algorithm’s robustness and accelerating 369  
the convergence of the model. 370

**Algorithm Design:** The algorithm of Ferret’s iterative 371  
gradient compensation is illustrated in Alg. 1 in the ap- 372  
pendix. Since the maximum possible  $\tau$  equals  $(P - 1)$ , 373  
the time complexity of the algorithm is  $\mathcal{O}(P - 1)$ , which 374  
is considered negligible during model training. Moreover, 375  
since two additional variables,  $v_r$  and  $v_a$ , are stored in mem- 376  
ory for optimizing  $\lambda$ , the space complexity of this algorithm 377  
is  $\mathcal{O}(2 \sum_{j=0}^{P-1} |w_j|)$ . However, by setting  $\eta_\lambda = 0$ , the opti- 378  
mization of  $\lambda$  is effectively terminated, and  $\lambda$  remains fixed 379  
at  $\lambda^0$ . This adjustment allows for manual tuning of  $\lambda$  and 380  
eliminates the need for  $v_r$  and  $v_a$ , thereby increasing flexi- 381  
bility and avoiding additional memory overhead. 382

## 5.2. Model Partitioning and Pipeline Planning 383

The objective of model partitioning and pipeline planning is 384  
to find an optimal model partition scheme  $L^*$  and its corre- 385  
sponding pipeline configuration  $C^*$  that maximize  $\mathcal{R}$  within 386  
a given memory constraint  $M$ , namely, 387

$$L^*, C^* = \arg \max_{L, C} \mathcal{R}_F^T \text{ s.t. } \mathcal{M}_F \leq M. \quad (13)$$

This problem can be reformulated as a bi-level optimiza- 389  
tion problem, decomposing it into two interrelated sub- 390  
problems: (1) determining the optimal  $C$  given a  $L$ , and (2) 391  
identifying the optimal  $L$  based on the solution from (1): 392

$$L^* = \arg \max_L \{\mathcal{R}_F^T | C_L^*\}$$

$$\text{s.t. } C_L^* = \arg \max_C \{\mathcal{R}_F^T | L\}, \mathcal{M}_F \leq M. \quad (14)$$

### 5.2.1 Iterative Configuration Search (Sub-problem 1) 395

Given a model partition scheme, the objective of sub- 396  
problem (1) is to solve 397

$$C^* = \arg \max_C \{\mathcal{R}_F^T | L\} \text{ s.t. } \mathcal{M}_F \leq M. \quad (15)$$

With more than  $2^{N(P+1)}$  potential combinations for  $C$ , a 399  
brute-force enumeration of  $C$  is impractical. Observing that 400  
 $d\mathcal{M}_F/d \max_C \{\mathcal{R}_F^T | L\} \geq 0$ , we employ an iterative algo- 401

rithm to determine the optimal  $C$  that maximize  $\mathcal{R}_F^T$  while ensuring  $\mathcal{M}_F$  remains within the memory budget. Specifically, to prevent memory over-consumption, we progressively deploy **T1-T4** as follows to balance  $\mathcal{R}_F^T$  and  $\mathcal{M}_F$ .

**S1. Deploy T1 for all workers:** By setting  $c_n^r = 1$  for all workers, the data processing time increases. Specifically, for the  $n$ -th worker, setting  $c_n^r = 1$  will respectively reduce  $\mathcal{R}_F^T$  and  $\mathcal{M}_F$  by Eq. 19 in the appendix.

**S2. Deploy T2 for the  $j$ -th stage in the  $n$ -th worker:** If  $c_{n,j}^o = 0$ , increasing  $c_{n,j}^a$  by  $\Delta c_{n,j}^a = \lceil \frac{P-j-1}{\lceil (P-j-1)/c_{n,j}^a \rceil - 1} \rceil - c_{n,j}^a$  will lead to a reduced frequency of model parameter updates. Here, the value of  $\Delta c_{n,j}^a$  is determined to prevent  $\Delta_{c_{n,j}^a \rightarrow c_{n,j}^a + 1} \mathcal{M}_F = 0$  due to the ceiling function. Consequently,  $\mathcal{R}_F^T$  and  $\mathcal{M}_F$  will be respectively decreased by Eq. 20 in the appendix.

**S3. Deploy T3 For the  $j$ -th stage in the  $n$ -th worker:** If  $\Delta c_{n,j}^a = +\infty$ , setting  $c_{n,j}^a = 1$  and  $c_{n,j}^o = P - 1 - j$  will completely eliminate the need for the  $j$ -th stage in the  $n$ -th worker to store additional model parameters by bypassing any backward pass that requires previous model parameters. Consequently,  $\mathcal{R}_F^T$  and  $\mathcal{M}_F$  will be respectively reduced by Eq. 21 in the appendix.

**S4. Deploy T4 for the  $n$ -th worker:** If  $c_{n,j}^o \neq 0$  for all  $j \in [0, p - 1)$ , removing the  $n$ -th worker will lead to a decrease in  $\mathcal{R}_F^T$  and  $\mathcal{M}_F$  by Eq. 22 in the appendix.

**Algorithm Design:** The algorithm of the proposed searching is illustrated in Alg. 2 in the appendix. Overall, the time complexity of this algorithm is  $\mathcal{O}(NP^2)$ , and it will be executed only once before fine-grained pipeline parallelism begins.

### 5.2.2 Brute-force Planning (Sub-problem 2)

In Ferret,  $L$  is determined by first establishing an upper bound on the time consumed for each stage ( $t^c$ ), and then solving the following optimization problem:

$$L = \arg \min_L \{P\} \text{ s.t. } t^f + t^b \leq t^c. \quad (16)$$

Namely, minimizing the number of pipeline stages while ensuring the time consumed for each stage is bounded. Since the layers in a stage must be consecutive, this problem can be solved in linear time by iteratively grouping consecutive layers into a stage until no additional adjacent layer can be grouped. Therefore, the solution space for  $L$  is not extensive, being limited to  $(\hat{L}^2 - \hat{L})/2$  at worst. Thus, to solve sub-problem (2), we can simply enumerate all possible model partition schemes, feeding them into Alg. 2 in the appendix to obtain the global optimum  $L^*$ .

**Algorithm Design:** The algorithm of the proposed planning is illustrated in Alg. 3 in the appendix. The time complexity of this algorithm is  $\mathcal{O}(\hat{L}^3)$ . Nevertheless, the algorithm will be executed only once before fine-grained pipeline parallelism begins.

## 6. Experiments

In this section, we seek answers to the following questions. (1) How does Ferret boost online accuracy? (Sec. 6.2) (2) How does Ferret mitigate catastrophic forgetting. (Sec. 6.2) (3) How does our fine-grained pipeline parallelism perform? (Sec. 6.3) (4) What are influences of different pipeline configurations? (Sec. 6.3) (5) How does our iterative gradient compensation algorithm perform?(Sec. 6.4)

### 6.1. Evaluation Setup

**Datasets and models:** Following the conventions of the community [45], 18 image classification datasets, including MNIST [22], FMNIST [78], CIFAR10 [42], CIFAR100 [42], SVHN [59], Tiny-ImageNet [44], CORE50 [49], CORE50-iid, Split-MNIST, Split-FMNIST, Split-CIFAR10, Split-CIFAR100, Split-SVHN, Split-Tiny-ImageNet, Coverttype [8], CLEAR10 [47], CLEAR100 [47], are used in our experiments. More details about the datasets can be found in the appendix. To cover both simple and complicated learning problems, five models including Multi-Layer Perceptron (MLP), MNISTNet, ConvNet, ResNet-18 [33] and MobileNet [34] are used in the experiments. Note that ResNet-18 and MobileNet are pretrained on the ImageNet-1K dataset [21].

**Compared methods:** For question (1): **Oracle**, **1-Skip** [28], **Random- $N$  B-Skip**, **Last- $N$  B-Skip** and **Camel** [45]. Here, Oracle is an ideal method that sequentially processes every streaming data without any delay. **B-Skip** and **Camel** selects a subset from the latest  $B$  unprocessed data using **Random- $N$** , **Last- $N$**  and **Coreset sampler**, respectively.

For question (2): **Vanilla**, **ER** [12], **MIR** [3], **LwF** [46], **MAS** [2]. For question (3) and (4): **DAPPLE** [24], **Pipedream** [57], **Pipedream<sub>2BW</sub>** [58], **Zero-Bubble** [65] and **Hanayo** [48]. For question (5): **None**, **Step-Aware** [32, 40], **Gap-Aware** [7], **Fisher** [14], and **Iter-Fisher** (*i.e.*, iterative gradient compensation).

**Evaluation metrics:** To measure catastrophic forgetting while accounting for memory footprint, Test Accuracy [28, 45] Gain per unit of Memory<sup>1</sup> (the higher the better) is defined as

$$\text{tagm}_{\mathcal{B}}(\mathcal{A}, t) = \log\left(\frac{\exp(\text{tacc}_{\mathcal{A}}(t) - \text{tacc}_{\mathcal{B}}(t))}{\mathcal{M}_{\mathcal{A}}/\mathcal{M}_{\mathcal{B}}}\right), \quad (17)$$

where  $\text{tacc}$  computes the test accuracy and  $\mathcal{B}$  is the baseline

<sup>1</sup>Results in the appendix show standard test accuracy and online accuracy but ignore memory consumption during training.

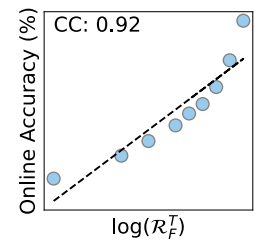


Figure 7. Relation between  $\text{oacc}$  and  $\log(\mathcal{R}_F^T)$

Table 1. Online Accuracy Gain per unit of Memory ( $agm_{\mathcal{B}}(\mathcal{A}, T)$ ) of different algorithms, where  $\mathcal{B}$  is the 1-Skip. "M-", "M", "M+" refer to the ferret method with minimal, medium and maximal memory footprint, respectively.

Setting	Oracle	1-Skip	Random-N	Last-N	Camel	Ferret <sub>M-</sub>	Ferret <sub>M</sub>	Ferret <sub>M+</sub>
MNIST/MNISTNet	27.32±0.71	0±0	-0.43±0.6	-0.26±0.14	-0.71±0.32	5.31±0.7	16.26±0.37	<b>26.34</b> ±0.7
FMNIST/MNISTNet	19.35±0.99	0±0	-0.31±0.47	-0.25±0.5	-0.6±0.4	5.93±0.81	12.69±0.81	<b>18.37</b> ±1.01
EMNIST/MNISTNet	13±0.48	0±0	1.94±0.04	2.02±0.03	1.55±0.1	4.19±0.17	8.8±0.4	<b>12.09</b> ±0.47
CIFAR10/ConvNet	10.57±0.09	0±0	4.71±0.05	4.78±0.03	4.7±0.05	3.21±0.16	6.21±0.15	<b>9.44</b> ±0.12
CIFAR100/ConvNet	5.24±0.01	0±0	0.78±0.07	0.83±0.06	0.75±0.08	1.58±0.04	2.6±0.03	<b>4.39</b> ±0.05
SVHN/ConvNet	15.41±0.23	0±0	7.04±0.08	7.24±0.11	7.39±0.09	5±0.1	11.52±0.23	<b>14.34</b> ±0.31
TinyImagenet/ConvNet	2.13±0.07	0±0	-0.22±0.04	-0.2±0.03	-0.21±0.05	0.48±0.03	0.54±0.03	<b>1.19</b> ±0.07
CORE50/ConvNet	26.01±0.42	0±0	12.13±0.42	12.27±0.44	11.07±0.48	9.08±0.45	17.95±0.45	<b>24.49</b> ±0.43
CORE50-iiid/ConvNet	19.24±2.9	0±0	2.87±5.71	5.74±2.8	5.21±2.49	3.55±2.77	10.74±2.76	<b>17.96</b> ±2.88
SplitMNIST/MNISTNet	18.21±0.76	0±0	2.34±0.43	2.37±0.63	3.3±0.48	6.11±0.84	14.55±0.55	<b>17.05</b> ±0.72
SplitFMNIST/MNISTNet	11.32±1.47	0±0	1.53±0.47	1.49±0.42	1.96±0.39	5.43±0.56	9.37±1.35	<b>10.29</b> ±1.47
SplitCIFAR10/ConvNet	7.49±0.12	0±0	3.05±0.11	3.11±0.11	3.12±0.09	2.91±0.19	4.84±0.2	<b>6.19</b> ±0.07
SplitCIFAR100/ConvNet	10.51±0.15	0±0	2.81±0.07	2.86±0.05	2.74±0.13	3.54±0.03	6.13±0.13	<b>9.61</b> ±0.04
SplitSVHN/ConvNet	6.49±0.33	0±0	2.9±0.19	2.91±0.21	2.89±0.21	2.76±0.16	5±0.28	<b>5.38</b> ±0.35
SplitTinyImagenet/ConvNet	2.14±0.1	0±0	-0.24±0.03	-0.21±0.02	-0.26±0.01	0.47±0.01	0.62±0.01	<b>1.19</b> ±0.06
CLEAR10/ResNet	10.37±0.06	0±0	7.84±0.07	7.93±0.06	-2.9±10.55	2.44±0.06	7.71±0.06	<b>9.26</b> ±0.08
CLEAR10/MobileNet	20.36±0.2	0±0	11.8±0.22	12±0.14	11.85±0.07	-1.77±0.15	14.68±0.5	<b>18.51</b> ±0.35
CLEAR100/ResNet	21.71±0.43	0±0	15.19±0.49	15.36±0.46	14.39±0.46	7.51±0.44	15.53±0.35	<b>20.84</b> ±0.57
CLEAR100/MobileNet	23.51±0.03	0±0	9.16±0.28	9.39±0.15	8.72±0.06	1.05±0.13	15.8±0.39	<b>22.11</b> ±0.59
Coverttype/MLP	7.66±0.27	0±0	-1.33±0.3	-1.3±0.31	-1.34±0.29	0.74±0.21	1.61±0.29	<b>3.38</b> ±0.42

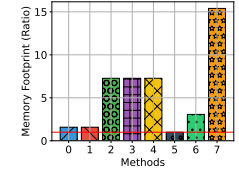
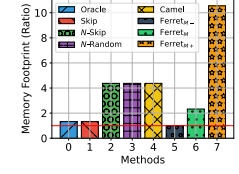


Figure 4. Consumed memory of different stream learning algorithms. Ferret achieves rapid adaptation across varying memory constraints.

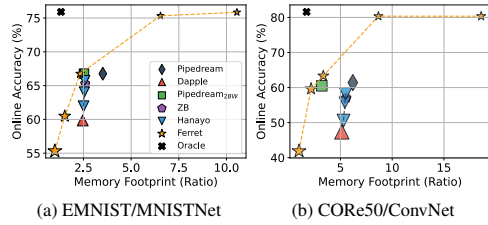


Figure 6. Relationships between online accuracy and memory consumption of different pipeline parallelism strategies, the marker size represents the standard errors of means.

Table 2. Online Accuracy Gain per unit of Memory ( $agm_{\mathcal{B}}(\mathcal{A}, T)$ ) and Test Accuracy Gain per unit of Memory ( $tagm_{\mathcal{B}}(\mathcal{A})$ ) of different integrated OCL algorithms on CORE50/ConvNet, where  $\mathcal{B}$  is the 1-Skip. Camel has its dedicated component to mitigate catastrophic forgetting and cannot be integrated with various OCL algorithm.

	Metric	Oracle	1-Skip	Random-N	Last-N	Camel	Ferret <sub>M-</sub>	Ferret <sub>M</sub>	Ferret <sub>M+</sub>
Vanilla	$agm$	26.01±0.42	0±0	12.13±0.42	12.27±0.44	11.07±0.48	9.08±0.45	17.21±0.45	<b>24.82</b> ±0.43
	$tagm$	2.36±0.64	0±0	1.08±0.62	0.97±0.39	1.48±0.42	1.01±0.45	1.07±0.49	<b>1.73</b> ±0.53
ER [12]	$agm$	24.03±0.26	0±0	7.84±0.17	8.11±0.29	-	7.12±0.09	16.09±0.16	<b>23.5</b> ±0.25
	$tagm$	4.18±0.43	0±0	1.94±0.26	2.34±0.29	-	0.82±0.29	3.1±0.25	<b>4.06</b> ±0.36
MIR [3]	$agm$	24.03±0.26	0±0	7.82±0.21	8.06±0.22	-	7.12±0.09	16.09±0.15	<b>23.5</b> ±0.25
	$tagm$	4.18±0.43	0±0	2.1±0.15	2.2±0.14	-	0.82±0.29	3.1±0.25	<b>4.06</b> ±0.36
LwF [46]	$agm$	26.02±0.42	0±0	12.25±0.42	12.4±0.45	-	9.02±0.4	17.96±0.46	<b>24.67</b> ±0.4
	$tagm$	2.36±0.64	0±0	1.2±0.62	1.09±0.39	-	0.91±0.44	1.88±0.49	<b>1.54</b> ±0.53
MAS [2]	$agm$	25.86±0.25	0±0	12.04±0.23	12.23±0.22	-	8.7±0.17	17.79±0.22	<b>24.46</b> ±0.23
	$tagm$	2.7±0.23	0±0	0.81±0.41	0.91±0.24	-	0.59±0.23	1.66±0.18	<b>1.69</b> ±0.21

500 method for comparison. Similarly, as Fig. 7 shows online  
 501 accuracy can be used to estimate  $\mathcal{R}_F^T$  [11], Online Accuracy  
 502 Gain per unit of Memory<sup>1</sup> (the higher the better) is defined as  
 503

$$504 \quad agm_{\mathcal{B}}(\mathcal{A}, t) = \log\left(\frac{\exp(oacc_{\mathcal{A}}(t) - oacc_{\mathcal{B}}(t))}{\mathcal{M}_{\mathcal{A}}/\mathcal{M}_{\mathcal{B}}}\right). \quad (18)$$

505 We evaluate three versions of Ferret under different  
 506 memory constraints: Ferret<sub>M-</sub> (minimal), Ferret<sub>M</sub> (the  
 507 same memory constraint as Pipedream<sub>2BW</sub>), and Ferret<sub>M+</sub>  
 508 (no constraint). Without clarification, each experiment is  
 509 independently repeated three times to obtain the final re-  
 510 sults. In all tables, the best and second-best performance  
 511 are highlighted by **bold** and underline, respectively. More  
 512 details about the evaluation setup can be found in Sec. 12.

## 513 6.2. Overall Comparisons

514 Table 1 shows  $agm_{\mathcal{B}}(\mathcal{A}, T)$  across 20 different settings to  
 515 evaluate both performance and consumed memory of differ-  
 516 ent frameworks. Here,  $\mathcal{B}$  is chosen to be the 1-Skip due to  
 517 its low memory footprint. From the table, it is evident that

518 Ferret<sub>M</sub> and Ferret<sub>M+</sub> constantly outperform other compet-  
 519 ing algorithms. Notably, Ferret<sub>M+</sub> even achieves compa-  
 520 rable performance compared to Oracle, indicating that Fer-  
 521 ret effectively enables rapid adaptation. On the other hand,  
 522 while Ferret<sub>M-</sub> shows slightly inferior performance com-  
 523 pared to its counterparts, it demands less memory for OCL,  
 524 as depicted in Fig. 4. This implies that in scenarios where  
 525 memory is severely constrained, Ferret is the only method  
 526 capable of learning.

527 Furthermore, various OCL algorithms are integrated on  
 528 CORE50/ConvNet in Table 2. It can be observed that Fer-  
 529 ret not only mitigates catastrophic forgetting (*i.e.*, increased  
 530  $tagm$ ) but also markedly enhances online performance (*i.e.*,  
 531 increased  $agm$ ), validating its orthogonality and superiority  
 532 compared to other OCL frameworks for rapid adaptation.

## 533 6.3. Comparisons on Pipeline Parallelism

534 Table 3 compares  $agm_{\mathcal{B}}(\mathcal{A}, T)$  of different pipeline paral-  
 535 lelism strategies across 20 different settings to evaluate the  
 536 performance of Ferret’s fine-grained pipeline parallelism

Table 3. Online Accuracy Gain per unit of Memory ( $agm_{\mathcal{B}}(\mathcal{A}, T)$ ) of different pipeline parallelism strategies, where  $\mathcal{B}$  is the DAPPLE. Note that "1W", "2W" and "3W" refer to 1, 2 and 3 wave(s) for the Hanayo algorithm, and no gradients compensation is applied to all asynchronous pipeline parallelism strategies for fair comparisons

Setting	Synchronous PP				Asynchronous PP			
	DAPPLE	ZB	Hanayo <sub>1W</sub>	Hanayo <sub>2W</sub>	Hanayo <sub>3W</sub>	Pipedream	Pipedream <sub>2W</sub>	Ferret <sub>M</sub>
MNIST/MnNet	0±0	6.79±0.4	2.44±0.3	5±0.16	7.12±0.38	8.16±0.35	8.23±0.39	8.35±0.35
FMNIST/MnNet	0±0	4.06±0.36	1.52±0.45	2.8±0.65	4.26±0.64	5.29±0.53	5.36±0.54	5.48±0.53
EMNIST/MnNet	0±0	2.33±0.08	0.9±0.02	1.81±0.05	2.55±0.04	2.84±0.09	2.99±0.07	3.02±0.09
C10/CNet	0±0	1.76±0.08	0.96±0.14	1.51±0.04	1.93±0.12	2.53±0.04	2.78±0.06	3.05±0.15
C100/CNet	0±0	0.71±0.04	0.05±0.05	0.56±0.06	0.74±0.06	0.87±0.09	1.11±0.06	1.72±0.01
SVHN/CNet	0±0	2.13±0.32	0.36±0.15	1.52±0.23	2.21±0.26	3.3±0.24	3.32±0.19	3.61±0.16
TinyI/CNet	0±0	0.18±0.02	0.03±0.01	0.19±0.02	0.19±0.04	0.26±0.01	0.52±0.03	0.5±0.06
CORE50/CNet	0±0	4.18±0.23	1.38±0.12	3.6±0.16	4.69±0.11	5.91±0.22	7.13±0.18	7.13±0.18
CORE50-iid/CNet	0±0	3.58±0.02	1.19±0.18	2.94±0.02	3.74±0.07	5.07±0.13	5.24±0.07	6.18±0.05
S-MNIST/MnNet	0±0	2.94±0.29	1.33±0.26	2.97±0.21	3.69±0.23	4.3±0.29	4.11±0.33	4.47±0.29
S-FMNIST/MnNet	0±0	1.56±0.29	0.91±0.16	1.48±0.2	1.89±0.31	2.06±0.28	2.09±0.27	2.24±0.28
S-C10/CNet	0±0	0.96±0.14	0.44±0.13	1.19±0.03	1.42±0.14	2.21±0.08	2.16±0.05	2.58±0.1
S-C100/CNet	0±0	1.57±0.05	0.54±0.14	1.25±0.12	1.67±0.12	2.48±0.12	2.49±0.1	3.49±0.06
S-SVHN/CNet	0±0	0.86±0.05	0.48±0.08	0.88±0.06	1.13±0.03	1.39±0.04	1.58±0.06	1.75±0.03
S-TinyI/CNet	0±0	0.27±0.05	0.08±0.03	0.14±0.02	0.22±0.03	0.29±0.03	0.47±0.01	0.66±0.04
CLEAR10/CNet	0±0	0.38±0.13	0.46±0.08	1.04±0.06	1.4±0.04	1.8±0.06	1.92±0.05	2.12±0.05
CLEAR10/MoNet	0±0	1.03±0.64	0.65±0.23	2.31±0.15	2.65±0.49	4.25±0.09	3.82±0.26	5.34±0.11
CLEAR100/RNet	0±0	2.76±0.1	1.36±0.22	2.52±0.21	3.3±0.23	3.85±0.2	3.98±0.19	4.24±0.22
CLEAR100/MoNet	0±0	3.11±0.53	1.26±0.12	3.03±0.52	4.24±0.12	5.66±0.19	5.88±0.58	7.42±0.69
Covertype/MLP	0±0	0.62±0.14	0.24±0.12	0.6±0.16	0.83±0.16	0.92±0.16	0.82±0.13	0.89±0.08

Table 4.

Online Accuracy differences between Ferret with and without gradients compensation algorithms .

Ferret <sub>M+</sub>				Ferret <sub>M</sub>			
Step-Aware	Gap-Aware	Fisher	Iter-Fisher	Step-Aware	Gap-Aware	Fisher	Iter-Fisher
-56.04±2.78	-14.03±1.24	-0.02±0.01	0.01±0	-43.53±2.36	-12.11±1.03	-0.01±0.01	0.02±0
-37.75±2.17	-9.07±0.63	-0.02±0.03	0.05±0.01	-37.74±2.53	-7.07±0.55	-0.01±0.01	0.02±0
-20.5±0.05	-4.81±0.15	0.01±0.02	0.04±0.02	-33.36±0.13	-3.46±0.33	0.01±0.01	0.04±0.02
-10.12±0.19	-1.71±0.43	-0.32±0.07	0.25±0.06	-9.6±0.19	-1.22±0.39	-0.14±0.3	0.42±0.21
-8.08±0.17	-2.17±0.05	-0.04±0.08	0.13±0.05	-5.4±0.07	-1.04±0.04	-0.01±0.07	0.1±0.02
-14.92±0.42	-2.63±0.04	0.02±0.04	0.31±0.2	-24.7±1.18	-2.91±0.11	-0.22±0.12	0.3±0.07
-3.72±0.16	-1.06±0.09	-0.01±0.02	0.06±0.03	-1.32±0.17	-0.11±0.19	0.17±0.19	0.35±0.15
-23.93±0.16	-3.27±0.05	-0.22±0.11	0.1±0.08	-33.41±0.24	-4.18±0.34	0.02±0.12	0.34±0.07
-24.56±0.22	-3.91±0.2	0.22±0.08	0.32±0.06	-23.77±0.3	-3.15±0.53	0.23±0.12	0.39±0.1
-26.8±3.2	-4.21±0.21	-0.05±0.02	0.03±0.01	-46.24±1.41	-4.82±0.48	-0.03±0.01	0.02±0
-14.43±2.83	-2.37±0.14	0.01±0.02	0.03±0.02	-46.07±4.03	-2.1±0.33	-0.01±0.01	0±0
-6.93±0.12	-1±0.14	-0.19±0.09	0.1±0.03	-8.11±0.75	-0.99±0.24	-0.12±0.12	0.23±0.12
-14.14±0.37	-3.05±0.1	-0.23±0.18	0.38±0.23	-12.56±0.23	-1.92±0.09	-0.1±0.19	0.24±0.09
-5.69±0.46	-1.18±0.12	-0.03±0.03	0.05±0.03	-12.13±0.87	-1.27±0.09	-0.05±0.02	0.03±0.01
-3.61±0.14	-1.01±0.05	0.05±0.1	0.12±0.1	-1.7±0.12	-0.37±0.03	0.08±0.07	0.18±0.08
-6.26±0.18	-0.72±0.08	0.02±0.06	0.14±0.04	-11.25±0.28	-0.52±0.05	-0.02±0.04	0.08±0.04
-18.17±0.26	-1.7±0.05	-0.14±0.32	0.5±0.15	-20.69±0.58	-1.88±0.46	-0.28±0.37	0.36±0.18
-17.45±0.02	-0.38±0.54	-0.07±0.04	0.65±0.32	-24.62±0.24	0.18±0.32	-0.2±0.18	1.27±0.21
-30.65±0.93	-0.3±0.17	0.63±0.46	1.2±0.48	-26.85±1.22	-0.3±0.33	0.58±0.78	1.03±0.7
-9.2±1.13	-2.8±0.1	-0.02±0.05	0.05±0.02	-10.27±0.39	-1.52±0.14	-0.01±0	0.09±0.01

537 under memory constraints. Specifically,  $\mathcal{B}$  is selected as  
538 DAPPLE, and no gradients compensation is applied to any  
539 asynchronous pipeline parallelism strategies. Additionally,  
540 Hanayo<sub>1W</sub>, Hanayo<sub>2W</sub> and Hanayo<sub>3W</sub> are three variants  
541 with 1, 2, and 3 waves, respectively.

542 In general, all asynchronous pipeline parallelism strategies  
543 significantly outperform synchronous pipeline parallelism  
544 strategies, even ZB, which claims to eliminate pipeline  
545 bubbles. This is because synchronous pipeline parallelism  
546 strategies, in an effort to achieve higher hardware utilization  
547 rates and avoid conflicting model versions, must design complex  
548 workflows that stage gradients and update model parameters  
549 synchronously, resulting in delays in data processing and wasted  
550 data value. Conversely, asynchronous pipeline parallelism strategies  
551 process data and update model parameters immediately, thereby  
552 minimizing processing latency. Among all asynchronous pipeline  
553 parallelism strategies, Ferret<sub>M</sub>'s fine-grained pipeline parallelism  
554 strategy consistently surpasses the others due to its more efficient  
555 memory utilization.

556 To investigate the impact of different pipeline configurations  
557 for Ferret, we select five different memory constraints ranging  
558 from minimum to maximum to simulate learning under varying  
559 memory budgets Fig. 6 shows that Ferret successfully solves Eq. 2  
560 for obtaining optimal pipeline configurations under dynamic environments,  
561 scaling effectively as we increase the memory constraint. Specifically,  
562 lack of precise control over each pipeline stage to balance between  
563 performance and memory footprint prevents competing strategies  
564 from scaling well.

## 567 6.4. Comparisons on Gradients Compensation

568 To evaluate the effectiveness of Iter-Fisher, we apply various  
569 gradients compensation algorithms to Ferret<sub>M+</sub> and

570 Ferret<sub>M</sub>, and compare the final online accuracy gain. The  
571 results are shown in Table 4. From the table, we can observe  
572 that applying Step-Aware and Gap-Aware algorithms for  
573 compensating stale gradients significantly reduces the online  
574 accuracy. This is because these algorithms mitigate the gradient  
575 staleness problem by simply penalizing the step size of stale  
576 gradients, leading to a slow convergence rate when the system  
577 is highly parallelized. Although Fisher leverages first-order  
578 information for better compensation, it does not consider varying  
579 levels of staleness at different stages of pipeline parallelism,  
580 resulting in a marginal decrease in accuracy compared to no  
581 compensation. On the other hand, Iter-Fisher consistently  
582 improves online accuracy across all settings, without requiring  
583 manual hyperparameters tuning. This indicates that Iter-Fisher  
584 effectively adapts to different levels of staleness in parallelism,  
585 and automatically optimizes  $\lambda$  for better compensation, demonstrating  
586 its robustness and effectiveness.

## 587 7. Conclusion

588 This paper introduces Ferret, a novel framework designed to  
589 boost online accuracy of OCL algorithms under varying memory  
590 constraints. Ferret employs a fine-grained pipeline parallelism  
591 strategy to adapt to varying distributions of incoming streaming  
592 data rapidly. To mitigate the gradient staleness problem in parallel  
593 processing, Ferret integrates an iterative gradient compensation  
594 algorithm to prevent performance degradation. Additionally,  
595 pipelines are automatically scheduled to improve performance  
596 under any memory scenario by optimizing a bi-level optimization  
597 problem. Extensive experiments conducted on 18 datasets and 5  
598 models confirm Ferret's superior efficiency and robustness  
599 compared to existing methods, demonstrating its potential as a  
600 scalable solution for adaptive, memory-efficient OCL.

604

605

**References**

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

*Workshop on Multi-Task and Lifelong Reinforcement Learning*, 2019. 1, 2, 6, 7, 5

659

660

[13] Tianqi Chen, Bing Xu, Chiyuan Zhang, and Carlos Guestrin. Training deep nets with sublinear memory cost. *arXiv preprint arXiv:1604.06174*, 2016. 3

661

662

663

[14] Yangrui Chen, Cong Xie, Meng Ma, Juncheng Gu, Yanghua Peng, Haibin Lin, Chuan Wu, and Yibo Zhu. Sapipe: Staleness-aware pipeline for data parallel dnn training. *Advances in Neural Information Processing Systems*, 35: 17981–17993, 2022. 6

664

665

666

667

668

[15] François Chollet. keras. <https://github.com/fchollet/keras>, 2015. 1

669

670

[16] Gregory Cohen, Saeed Afshar, Jonathan Tapson, and Andre Van Schaik. Emnist: Extending mnist to handwritten letters. In *2017 international joint conference on neural networks (IJCNN)*, pages 2921–2926. IEEE, 2017. 4

671

672

673

674

[17] Bitu Darvish Rouhani, Huili Chen, and Farinaz Koushanfar. Deepsigns: An end-to-end watermarking framework for ownership protection of deep neural networks. In *Proceedings of the twenty-fourth international conference on architectural support for programming languages and operating systems*, pages 485–497, 2019. 1, 2

675

676

677

678

679

680

[18] Marcos Dias de Assuncao, Alexandre da Silva Veith, and Rajkumar Buyya. Distributed data stream processing and edge computing: A survey on resource elasticity and future directions. *Journal of Network and Computer Applications*, 103: 1–17, 2018. 1

681

682

683

684

685

[19] Matthias De Lange and Tinne Tuytelaars. Continual prototype evolution: Learning online from non-stationary data streams. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8250–8259, 2021. 1, 2

686

687

688

689

[20] Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Aleš Leonardis, Gregory Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *IEEE transactions on pattern analysis and machine intelligence*, 44(7):3366–3385, 2021. 3

690

691

692

693

694

695

[21] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 6

696

697

698

699

[22] Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012. 6, 4

700

701

702

[23] Matthew F Dixon, Igor Halperin, and Paul Bilokon. *Machine learning in finance*. Springer, 2020. 1

703

704

[24] Shiqing Fan, Yi Rong, Chen Meng, Zongyan Cao, Siyu Wang, Zhen Zheng, Chuan Wu, Guoping Long, Jun Yang, Lixue Xia, et al. Dapple: A pipelined data parallel approach for training large models. In *Proceedings of the 26th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, pages 431–445, 2021. 6, 1

705

706

707

708

709

710

[25] Chrisantha Fernando, Dylan Banarse, Charles Blundell, Yori Zwols, David Ha, Andrei A Rusu, Alexander Pritzel, and Daan Wierstra. Pathnet: Evolution channels gradient descent in super neural networks. *arXiv preprint arXiv:1701.08734*, 2017. 1, 2

711

712

713

714

715

- 716 [26] Chelsea Finn, Aravind Rajeswaran, Sham Kakade, and  
717 Sergey Levine. Online meta-learning. In *International  
718 conference on machine learning*, pages 1920–1930. PMLR,  
719 2019. 1, 2
- 720 [27] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The  
721 elements of statistical learning*. Springer series in statistics  
722 New York, 2001. 4
- 723 [28] Yasir Ghunaim, Adel Bibi, Kumail Alhamoud, Motasem  
724 Alfarra, Hasan Abed Al Kader Hammoud, Ameya Prabhu,  
725 Philip HS Torr, and Bernard Ghanem. Real-time evaluation  
726 in online continual learning: A new hope. In *Proceedings of  
727 the IEEE/CVF Conference on Computer Vision and Pattern  
728 Recognition*, pages 11888–11897, 2023. 1, 2, 6
- 729 [29] Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat  
730 Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan  
731 Misra. Imagebind: One embedding space to bind them all.  
732 In *Proceedings of the IEEE/CVF Conference on Computer  
733 Vision and Pattern Recognition*, pages 15180–15190, 2023.  
734 1
- 735 [30] Yue Guan, Yuxian Qiu, Jingwen Leng, Fan Yang, Shuo Yu,  
736 Yunxin Liu, Yu Feng, Yuhao Zhu, Lidong Zhou, Yun Liang,  
737 et al. Amanda: Unified instrumentation framework for deep  
738 neural networks. 2023. 1, 2
- 739 [31] Nuwan Gunasekara, Bernhard Pfahringer, Heitor Murilo  
740 Gomes, and Albert Bifet. Survey on online streaming con-  
741 tinual learning. In *Proceedings of the Thirty-Second Inter-  
742 national Joint Conference on Artificial Intelligence, IJCAI  
743 2023, 19th-25th August 2023, Macao, SAR, China*, pages  
744 6628–6637. ijcai. org, 2023. 1
- 745 [32] Corentin Hardy, Erwan Le Merrer, and Bruno Sericola. Dis-  
746 tributed deep learning on edge-devices: feasibility via adap-  
747 tive compression. In *2017 IEEE 16th international symposium  
748 on network computing and applications (NCA)*, pages  
749 1–8. IEEE, 2017. 6, 1
- 750 [33] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun.  
751 Deep residual learning for image recognition. In *Proceed-  
752 ings of the IEEE conference on computer vision and pattern  
753 recognition*, pages 770–778, 2016. 6
- 754 [34] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry  
755 Kalenichenko, Weijun Wang, Tobias Weyand, Marco An-  
756 dreetto, and Hartwig Adam. Mobilenets: Efficient convolu-  
757 tional neural networks for mobile vision applications. *arXiv  
758 preprint arXiv:1704.04861*, 2017. 6
- 759 [35] Qinghao Hu, Zhisheng Ye, Meng Zhang, Qiaoling  
760 Chen, Peng Sun, Yonggang Wen, and Tianwei Zhang.  
761 Hydro: {Surrogate-Based} hyperparameter tuning service in  
762 datacenters. In *17th USENIX Symposium on Operating Sys-  
763 tems Design and Implementation (OSDI 23)*, pages 757–777,  
764 2023. 1, 2
- 765 [36] Xing Hu, Ling Liang, Shuangchen Li, Lei Deng, Pengfei  
766 Zuo, Yu Ji, Xinfeng Xie, Yufei Ding, Chang Liu, Timothy  
767 Sherwood, et al. Deepsniffer: A dnn model extraction frame-  
768 work based on learning architectural hints. In *Proceedings  
769 of the Twenty-Fifth International Conference on Architec-  
770 tural Support for Programming Languages and Operating  
771 Systems*, pages 385–399, 2020. 1, 2
- 772 [37] Yanping Huang, Youlong Cheng, Ankur Bapna, Orhan Firat,  
773 Dehao Chen, Mia Chen, HyoukJoong Lee, Jiquan Ngiam,  
Quoc V Le, Yonghui Wu, et al. Gpipe: Efficient training of  
giant neural networks using pipeline parallelism. *Advances  
in neural information processing systems*, 32, 2019. 1
- [38] David Isele and Akansel Cosgun. Selective experience re-  
play for lifelong learning. In *Proceedings of the AAAI Con-  
ference on Artificial Intelligence*, 2018. 1, 2
- [39] Anand Jayarajan, Kimberly Hau, Andrew Goodwin, and  
Gennady Pekhimenko. Lifestream: a high-performance  
stream processing engine for periodic streams. In *Procee-  
dings of the 26th ACM International Conference on Architec-  
tural Support for Programming Languages and Operating  
Systems*, pages 107–122, 2021. 1
- [40] Jiawei Jiang, Bin Cui, Ce Zhang, and Lele Yu.  
Heterogeneity-aware distributed parameter servers. In *Pro-  
ceedings of the 2017 ACM International Conference on Man-  
agement of Data*, pages 463–478, 2017. 6, 1
- [41] Jin Kyu Kim, Qirong Ho, Seunghak Lee, Xun Zheng, Wei  
Dai, Garth A Gibson, and Eric P Xing. Strads: A dis-  
tributed framework for scheduled model parallel machine  
learning. In *Proceedings of the Eleventh European Confer-  
ence on Computer Systems*, pages 1–16, 2016. 1, 2
- [42] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple  
layers of features from tiny images. 2009. 6, 4
- [43] Young D Kwon, Jagmohan Chauhan, Hong Jia, Stylianos I  
Venieris, and Cecilia Mascolo. Lifelearner: Hardware-aware  
meta continual learning system for embedded computing  
platforms. *arXiv preprint arXiv:2311.11420*, 2023. 2
- [44] Ya Le and Xuan Yang. Tiny imagenet visual recognition  
challenge. *CS 231N*, 7(7):3, 2015. 6, 3, 4
- [45] Yiming Li, Yanyan Shen, and Lei Chen. Camel: Manag-  
ing data for efficient stream learning. In *Proceedings of  
the 2022 International Conference on Management of Data*,  
pages 1271–1285, 2022. 2, 6, 1, 3
- [46] Zhizhong Li and Derek Hoiem. Learning without forgetting.  
*IEEE transactions on pattern analysis and machine intelli-  
gence*, 40(12):2935–2947, 2017. 6, 7, 1, 5
- [47] Zhiqiu Lin, Jia Shi, Deepak Pathak, and Deva Ramanan. The  
clear benchmark: Continual learning on real-world imagery.  
In *Thirty-fifth conference on neural information processing  
systems datasets and benchmarks track (round 2)*, 2021. 1,  
2, 6, 3, 4
- [48] Ziming Liu, Shenggan Cheng, Haotian Zhou, and Yang You.  
Hanayo: Harnessing wave-like pipeline parallelism for en-  
hanced large model training efficiency. In *Proceedings of the  
International Conference for High Performance Computing,  
Networking, Storage and Analysis*, pages 1–13, 2023. 6, 1
- [49] Vincenzo Lomonaco and Davide Maltoni. Core50: a new  
dataset and benchmark for continuous object recognition. In  
*Conference on robot learning*, pages 17–26. PMLR, 2017. 6,  
3, 4
- [50] David Lopez-Paz and Marc’Aurelio Ranzato. Gradient  
episodic memory for continual learning. *Advances in neu-  
ral information processing systems*, 30, 2017. 1
- [51] Zheda Mai, Ruiwen Li, Jihwan Jeong, David Quispe, Hyun-  
woo Kim, and Scott Sanner. Online continual learning in  
image classification: An empirical survey. *Neurocomputing*,  
469:28–51, 2022. 1

- 831 [52] Gaurav Menghani. Efficient deep learning: A survey on  
832 making deep learning models smaller, faster, and better.  
833 *ACM Computing Surveys*, 55(12):1–37, 2023. 1 888
- 834 [53] Baharan Mirzasoleiman, Jeff Bilmes, and Jure Leskovec.  
835 Coresets for data-efficient training of machine learning mod-  
836 els. In *International Conference on Machine Learning*, pages  
837 6950–6960. PMLR, 2020. 1, 2 889
- 838 [54] Philipp Moritz, Robert Nishihara, Stephanie Wang, Alexey  
839 Tumanov, Richard Liaw, Eric Liang, Melih Elibol, Zongheng  
840 Yang, William Paul, Michael I Jordan, et al. Ray: A dis-  
841 tributed framework for emerging {AI} applications. In *13th*  
842 *USENIX symposium on operating systems design and imple-*  
843 *mentation (OSDI 18)*, pages 561–577, 2018. 1, 2 890
- 844 [55] Khan Muhammad, Amin Ullah, Jaime Lloret, Javier Del Ser,  
845 and Victor Hugo C de Albuquerque. Deep learning for safe  
846 autonomous driving: Current challenges and future direc-  
847 tions. *IEEE Transactions on Intelligent Transportation Sys-*  
848 *tems*, 22(7):4316–4336, 2020. 1 891
- 849 [56] MG Sarwar Murshed, Christopher Murphy, Daqing Hou,  
850 Nazar Khan, Ganesh Ananthanarayanan, and Faraz Hussain.  
851 Machine learning at the network edge: A survey. *ACM Com-*  
852 *puting Surveys (CSUR)*, 54(8):1–37, 2021. 1 892
- 853 [57] Deepak Narayanan, Aaron Harlap, Amar Phanishayee,  
854 Vivek Seshadri, Nikhil R Devanur, Gregory R Ganger,  
855 Phillip B Gibbons, and Matei Zaharia. Pipedream: gener-  
856 alized pipeline parallelism for dnn training. In *Proceedings*  
857 *of the 27th ACM symposium on operating systems principles*,  
858 pages 1–15, 2019. 4, 6, 1 893
- 859 [58] Deepak Narayanan, Amar Phanishayee, Kaiyu Shi, Xie  
860 Chen, and Matei Zaharia. Memory-efficient pipeline-parallel  
861 dnn training. In *International Conference on Machine Learn-*  
862 *ing*, pages 7937–7947. PMLR, 2021. 3, 4, 6, 1 894
- 863 [59] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bis-  
864 sacco, Baolin Wu, Andrew Y Ng, et al. Reading digits in  
865 natural images with unsupervised feature learning. In *NIPS*  
866 *workshop on deep learning and unsupervised feature learn-*  
867 *ing*, page 7. Granada, Spain, 2011. 6, 4 895
- 868 [60] Shuaicheng Niu, Jiayang Wu, Yifan Zhang, Yafo Chen,  
869 Shijian Zheng, Peilin Zhao, and Mingkui Tan. Efficient  
870 test-time model adaptation without forgetting. In *Interna-*  
871 *tional conference on machine learning*, pages 16888–16905.  
872 PMLR, 2022. 1 896
- 873 [61] NVIDIA, Péter Vingelmann, and Frank H.P. Fitzek. Cuda,  
874 release: 10.2.89, 2020. 2 897
- 875 [62] Alecos Papadopoulos. The information matrix equality:  
876 proof, misspecification, and the quasi-maximum likelihood  
877 estimator. *Athens University of Economics and Business*,  
878 2014. 4 898
- 879 [63] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer,  
880 James Bradbury, Gregory Chanan, Trevor Killeen, Zeming  
881 Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An im-  
882 perative style, high-performance deep learning library. *Ad-*  
883 *vances in neural information processing systems*, 32, 2019.  
884 1 899
- 885 [64] Ameya Prabhu, Philip HS Torr, and Puneet K Dokania.  
886 Gdumb: A simple approach that questions our progress in  
887 continual learning. In *Computer Vision–ECCV 2020: 16th*  
*European Conference, Glasgow, UK, August 23–28, 2020,*  
*Proceedings, Part II 16*, pages 524–540. Springer, 2020. 1,  
2 900
- [65] Penghui Qi, Xinyi Wan, Guangxing Huang, and Min  
Lin. Zero bubble pipeline parallelism. *arXiv preprint*  
*arXiv:2401.10241*, 2023. 6, 1 901
- [66] Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and  
Yuxiong He. Deepspeed: System optimizations enable train-  
ing deep learning models with over 100 billion paramet-  
ers. In *Proceedings of the 26th ACM SIGKDD International*  
*Conference on Knowledge Discovery & Data Mining*, pages  
3505–3506, 2020. 1 902
- [67] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg  
Sperl, and Christoph H Lampert. icarl: Incremental classifier  
and representation learning. In *Proceedings of the IEEE con-*  
*ference on Computer Vision and Pattern Recognition*, pages  
2001–2010, 2017. 1, 2 903
- [68] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins,  
Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Raz-  
van Pascanu, and Raia Hadsell. Progressive neural networks.  
*arXiv preprint arXiv:1606.04671*, 2016. 1, 2 904
- [69] Doyen Sahoo, Quang Pham, Jing Lu, and Steven CH Hoi.  
Online deep learning: Learning deep neural networks on the  
fly. *arXiv preprint arXiv:1711.03705*, 2017. 1, 2 905
- [70] David Sayce. The number of tweets per day in 2022. 2022.  
1 906
- [71] Alexander Sergeev and Mike Del Balso. Horovod: fast and  
easy distributed deep learning in tensorflow. *arXiv preprint*  
*arXiv:1802.05799*, 2018. 1 907
- [72] Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick  
LeGresley, Jared Casper, and Bryan Catanzaro. Megatron-  
lm: Training multi-billion parameter language models using  
model parallelism. *arXiv preprint arXiv:1909.08053*, 2019.  
1 908
- [73] Kai Sheng Tai, Vatsal Sharan, Peter Bailis, and Gregory  
Valiant. Sketching linear classifiers over data streams. In  
*Proceedings of the 2018 international conference on man-*  
*agement of data*, pages 757–772, 2018. 1 909
- [74] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier  
Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste  
Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al.  
Llama: Open and efficient foundation language models.  
*arXiv preprint arXiv:2302.13971*, 2023. 1 910
- [75] Ehsan Valavi, Joel Hestness, Newsha Ardalani, and Marco  
Iansiti. Time and the value of data. *arXiv preprint*  
*arXiv:2203.09118*, 2022. 1, 3 911
- [76] Siniša Veseli, John Hammonds, Steven Henke, Hannah Par-  
raga, and Nicholas Schwarz. Streaming data from exper-  
imental facilities to supercomputers for real-time data pro-  
cessing. In *Proceedings of the SC’23 Workshops of The Inter-*  
*national Conference on High Performance Computing, Net-*  
*work, Storage, and Analysis*, pages 2110–2117, 2023. 1 912
- [77] Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. A  
comprehensive survey of continual learning: Theory, method  
and application. *IEEE Transactions on Pattern Analysis and*  
*Machine Intelligence*, 2024. 1 913
- [78] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-  
mnist: a novel image dataset for benchmarking machine  
914 915

- 946 learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.  
947 6, 4
- 948 [79] Minhui Xie, Kai Ren, Youyou Lu, Guangxu Yang, Qingx-  
949 ing Xu, Bihai Wu, Jiazhen Lin, Hongbo Ao, Wanhong Xu,  
950 and Jiwu Shu. Kraken: memory-efficient continual learning  
951 for large-scale real-time recommendations. In *SC20: Inter-  
952 national Conference for High Performance Computing, Net-  
953 working, Storage and Analysis*, pages 1–17. IEEE, 2020. 2
- 954 [80] Yu Yang, Hao Kang, and Baharan Mirzasoleiman. Towards  
955 sustainable learning: Coresets for data-efficient deep learn-  
956 ing. In *International Conference on Machine Learning*,  
957 pages 39314–39330. PMLR, 2023. 1
- 958 [81] Jaehong Yoon, Eunho Yang, Jeongtae Lee, and Sung Ju  
959 Hwang. Lifelong learning with dynamically expandable net-  
960 works. *arXiv preprint arXiv:1708.01547*, 2017. 1, 2
- 961 [82] Jaehong Yoon, Divyam Madaan, Eunho Yang, and Sung Ju  
962 Hwang. Online coreset selection for rehearsal-based contin-  
963 ual learning. *arXiv preprint arXiv:2106.01085*, 2021. 1, 2
- 964 [83] Quanlu Zhang, Zhenhua Han, Fan Yang, Yuge Zhang, Zhe  
965 Liu, Mao Yang, and Lidong Zhou. Retiarrii: A deep learning  
966 {Exploratory-Training} framework. In *14th USENIX Sym-  
967 posium on Operating Systems Design and Implementation  
968 (OSDI 20)*, pages 919–936, 2020. 1, 2
- 969 [84] Shuxin Zheng, Qi Meng, Taifeng Wang, Wei Chen, Nenghai  
970 Yu, Zhi-Ming Ma, and Tie-Yan Liu. Asynchronous stochas-  
971 tic gradient descent with delay compensation. In *Internat-  
972 ional conference on machine learning*, pages 4120–4129.  
973 PMLR, 2017. 5, 1
- 974 [85] Zhen Zheng, Chanyoung Oh, Jidong Zhai, Xipeng Shen,  
975 Youngmin Yi, and Wenguang Chen. Hiwaylib: A software  
976 framework for enabling high performance communications  
977 for heterogeneous pipeline computations. In *Proceedings  
978 of the Twenty-Fourth International Conference on Architec-  
979 tural Support for Programming Languages and Operating  
980 Systems*, pages 153–166, 2019. 1, 2
- 981 [86] Yuhao Zhou, Qing Ye, and Jiancheng Lv. Communication-  
982 efficient federated learning with compensated overlap-  
983 fedavg. *IEEE Transactions on Parallel and Distributed Sys-  
984 tems*, 33(1):192–205, 2021. 1
- 985 [87] Yuhao Zhou, Mingjia Shi, Yuanxi Li, Yanan Sun, Qing Ye,  
986 and Jiancheng Lv. Communication-efficient federated learn-  
987 ing with single-step synthetic features compressor for faster  
988 convergence. In *Proceedings of the IEEE/CVF International  
989 Conference on Computer Vision*, pages 5031–5040, 2023. 1