

Non-Idle Machine-Aware Worker Placement for Efficient Distributed Training in GPU Clusters

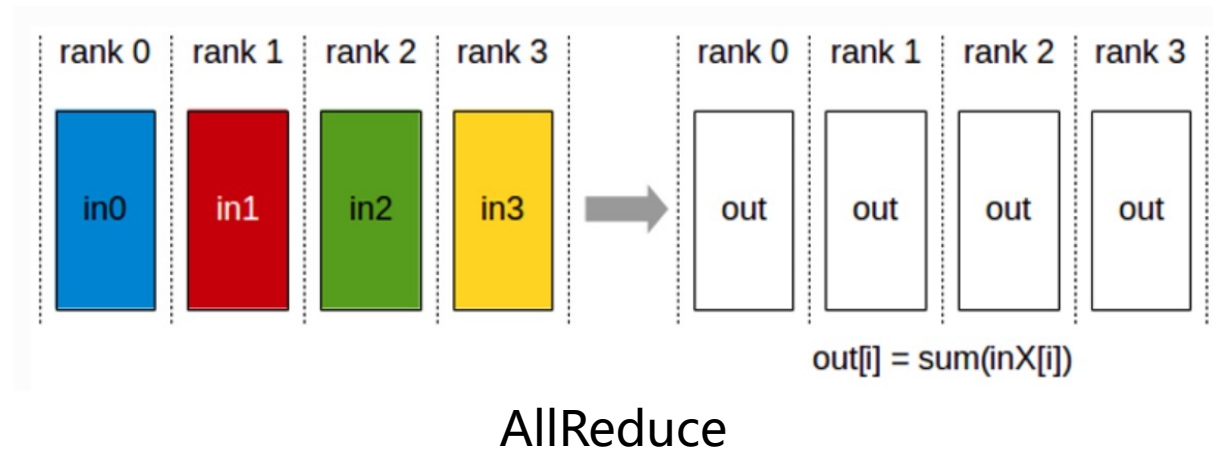
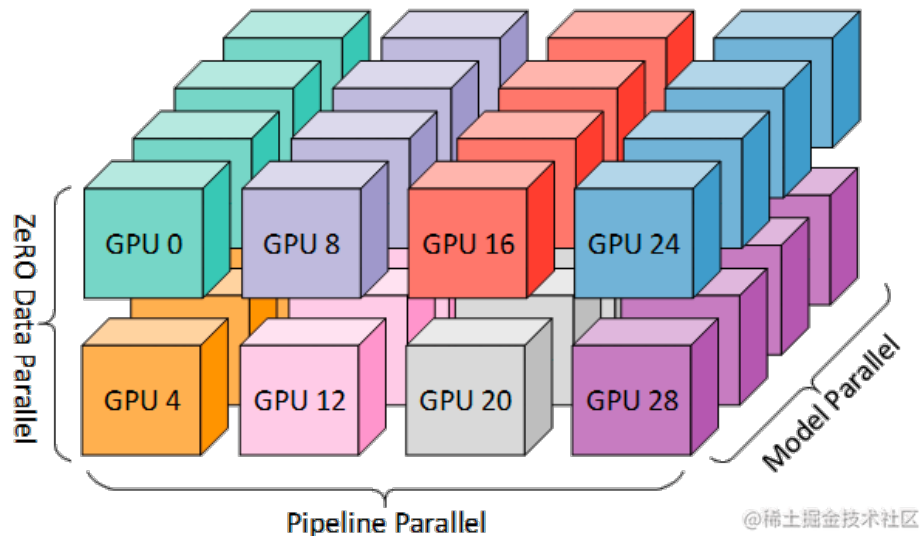
Jin Fang

Gongming Zhao, Hongli Xu, Luyao Luo, Zhen Yao, An Xie



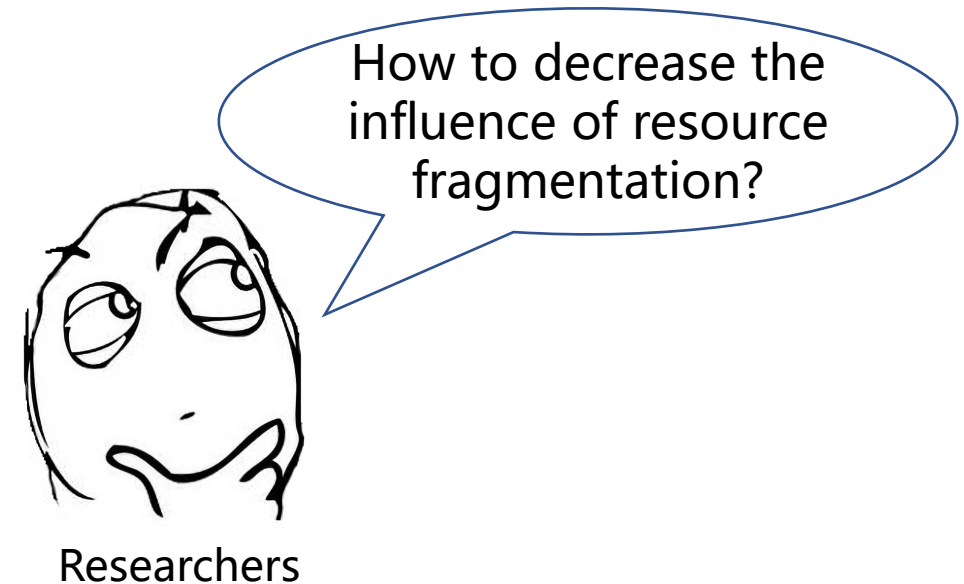
Collective Communication in DT

- With the increasing complexity of machine learning (ML) applications, the scale of ML tasks grows explosively
- **Distributed training** is proposed to speed up the training of large-scale ML tasks
- Placing workers on GPUs of machines to perform one DT task, where workers **communicates collectively** to synchronize gradients/parameters



Problem: Resource Fragmentation

- One machine is equipped with multiple GPUs (e.g., 8 GPUs)
- The number of workers of a DT job varies
 - DP = 2 -> 2 workers
 - TP = 8 -> 8 workers
 - DP = 2 and TP = 8 -> $2 \cdot 8 = 16$ workers
- The arrival time of DT jobs is unpredictable
- There exists a lot of fragmented idle GPUs, leading to low resource utilization



Existing Solution

Consolidation-First Placement

- Elasticflow¹ (ASPLOS 23): allocate the resource at the scale of machine
Resource oversubscription

Fragmentation-First Placement

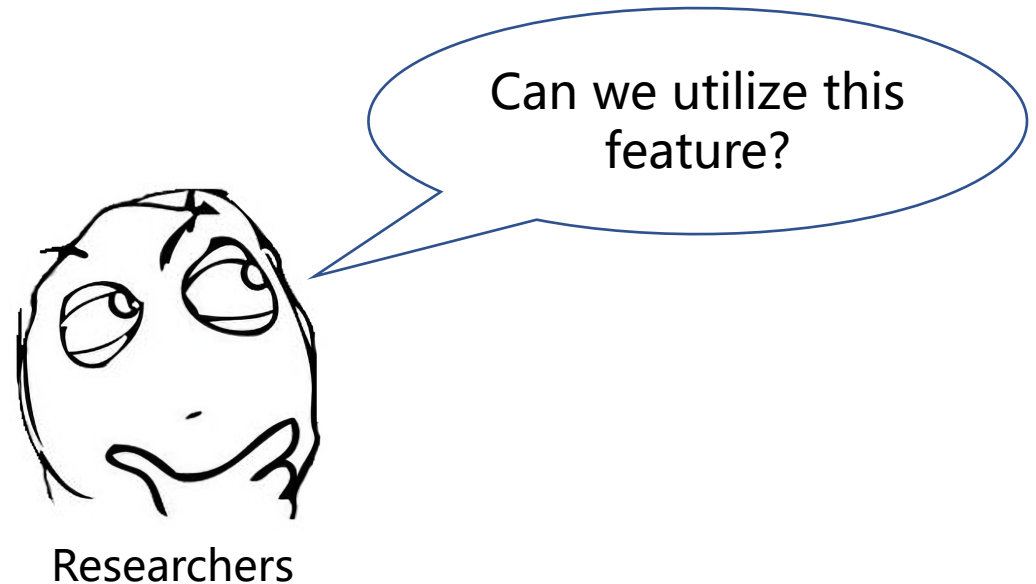
- HiveD² (OSDI 20): priority place placing workers in fragmented machines
Large comm. overhead

1. Gu D, Zhao Y, Zhong Y, et al. ElasticFlow: An elastic serverless training platform for distributed deep learning[C]//Proceedings of the 28th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2. 2023: 266-280.
2. Zhao H, Han Z, Yang Z, et al. {HiveD}: Sharing a {GPU} cluster for deep learning with guarantees[C]//14th USENIX symposium on operating systems design and implementation (OSDI 20). 2020: 515-532.

A Motivating Example

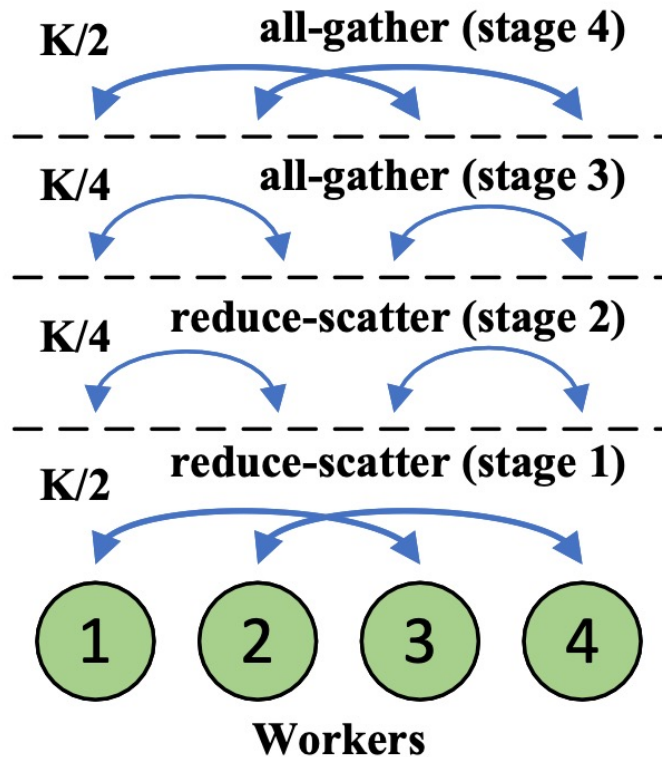
- Collective communication operations take multiple steps, where each step contains different communication pairs and amounts

Workers in one DT job have different communication pattern.



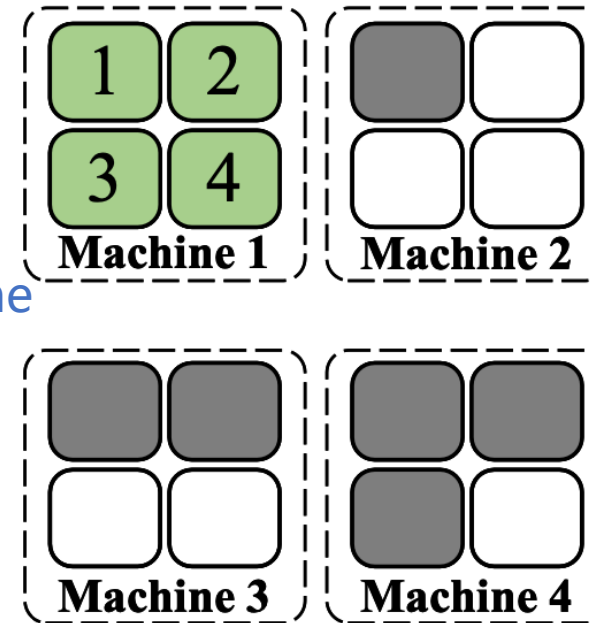
A Motivating Example-Elasticflow

- Consolidation placement to minimize the amount of cross-machine traffic



(a) Communication Pattern Example

Launch an idle machine

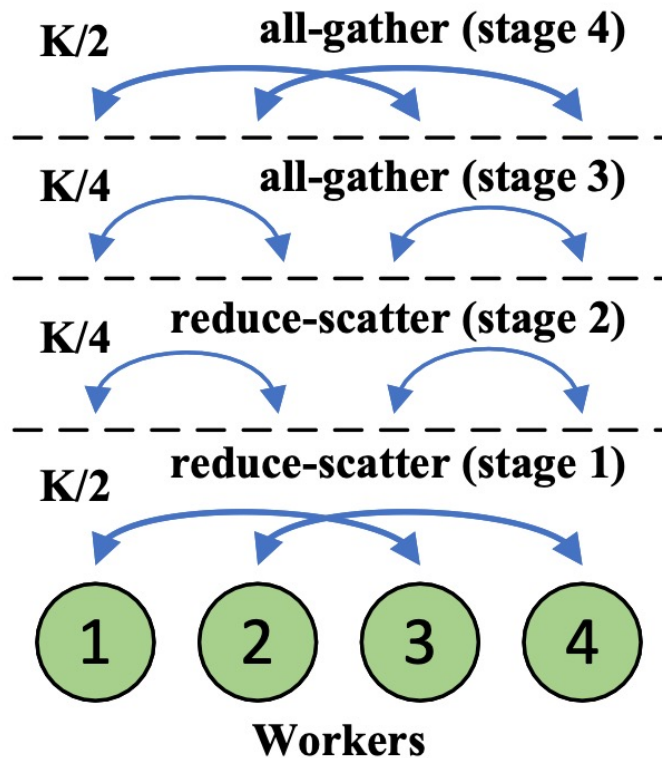


No. of Machines = 4
Cross-machine Traffic = 0

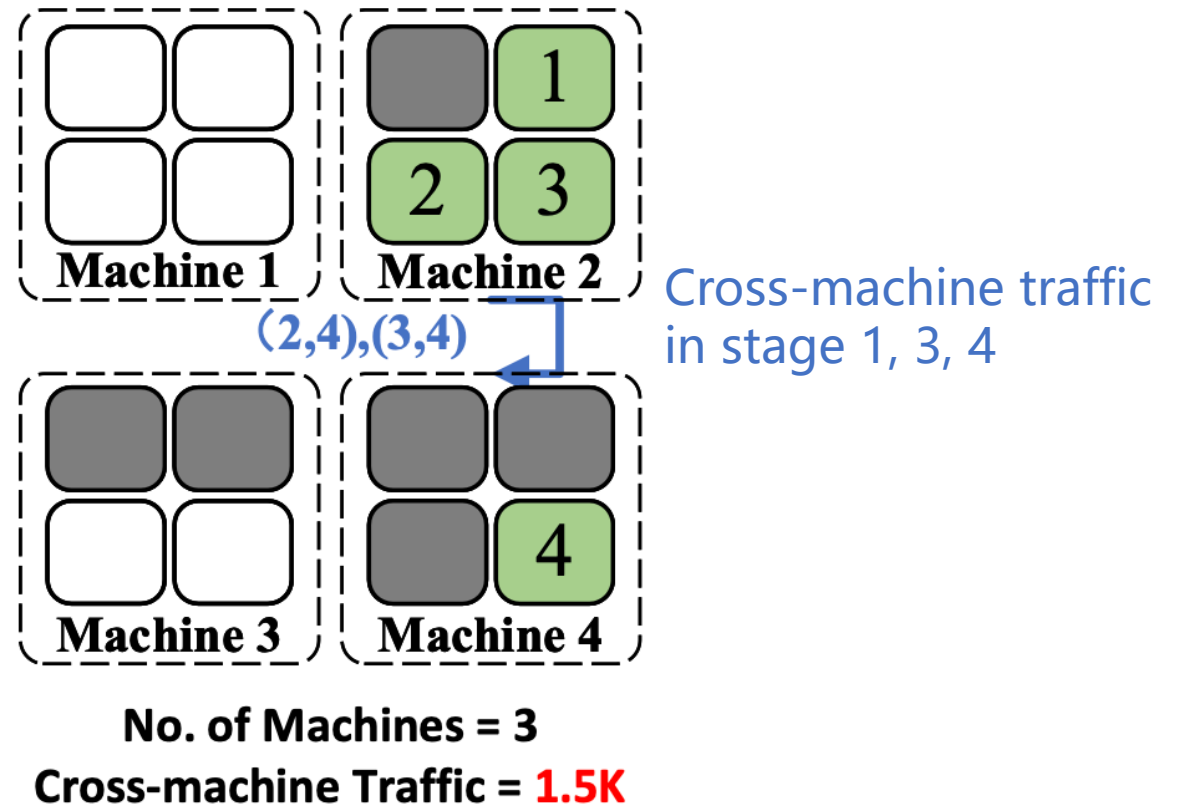
(b) Elasticflow

A Motivating Example-HiveD

- Using fragmented machines to minimize the number of used machines



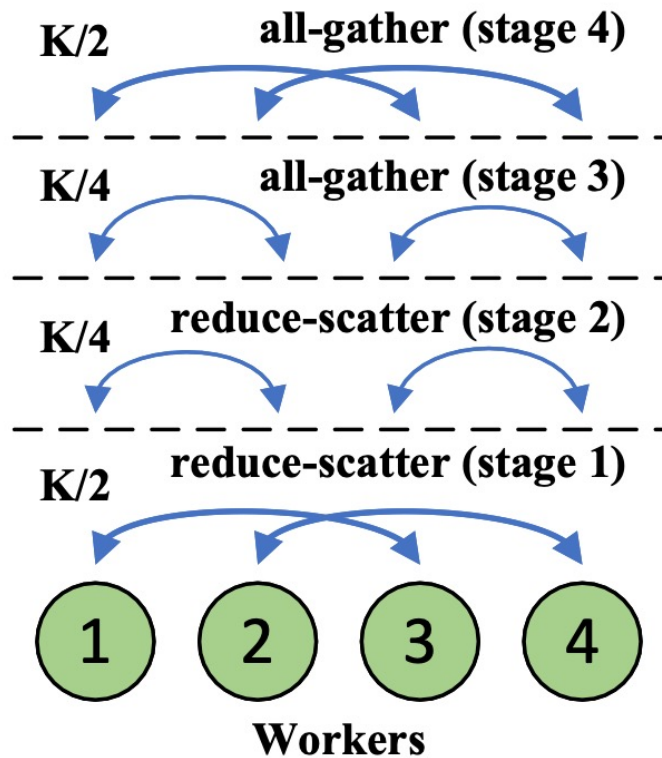
(a) Communication Pattern Example



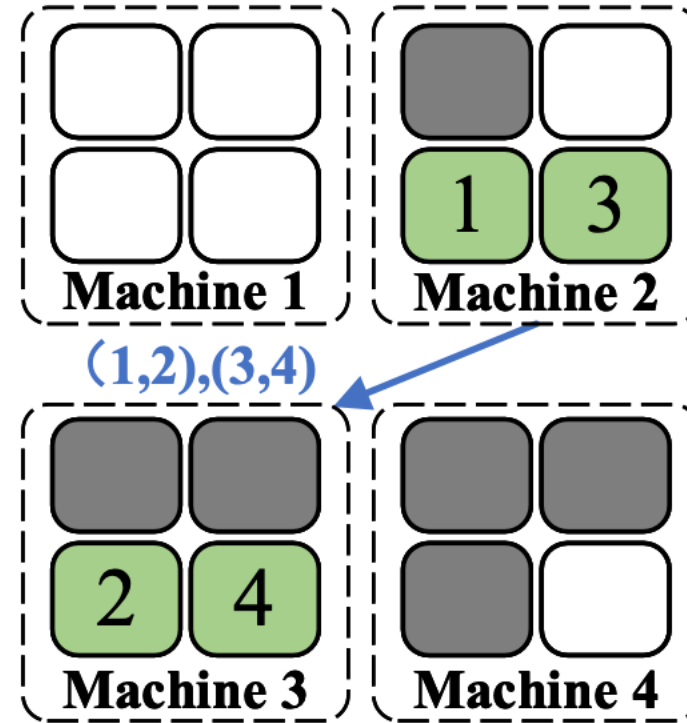
(c) HiveD

A Motivating Example-Titan (Ours)

- Consider the map of workers to GPUs according to collective communication algorithms
- Avoid the use of idle machines



(a) Communication Pattern Example



No. of Machines = 3
Cross-machine Traffic = K

(d) Titan (Ours)

Titan: Problem Formulation

Network Model

- Machine set: $\mathcal{S} = \{s_1, s_2, \dots, s_{|\mathcal{S}|}\}$, each with K GPUs
- Non-Idle machine set: $\mathcal{S}_f \subset \mathcal{S}$
- Idle machine set: $\mathcal{S}_n \subset \mathcal{S}$
- Available bandwidth between machines s and s' : $P_{s,s'} \in \mathbb{Z}$

Communication Pattern

- Worker set: $\mathcal{N} = \{n_1, n_2, \dots, n_{|\mathcal{N}|}\}$
- The traffic amount between worker pair (n, n') in phase t : $C_{n,n'}^t \in \mathbb{Z}$

Titan: Problem Formulation

➤ Objective

1. **Minimize the number of Idle-machines:** reduce resource fragmentation
2. **Minimize the number of total machines:** reduce the amount of cross-machine traffic

$$\min O_1 = \sum_{s \in S_n} y_s$$

$$\min O_2 = \sum_{s \in S} y_s$$

- Placement constraint
- Resource constraint
- Bandwidth constraint

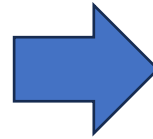
$$S.t. \begin{cases} \sum_{s \in S} x_n^s = 1, & \forall n \in N \\ \sum_{n \in N} x_n^s \leq R_s, & \forall s \in S \\ \sum_{n \in N} \sum_{n' \in N} x_n^s \cdot x_{n'}^{s'} \cdot C_{n,n'}^t \leq P_{s,s'}, & \forall s, s' \in S, t \in T \\ y_s \geq x_n^s, & \forall s \in S, n \in N \\ x_n^s \in \{0, 1\}, & \forall s \in S, n \in N \\ y_s \in \{0, 1\}, & \forall s \in S \end{cases} \quad (1)$$

Titan: Algorithm Design

- Convert the problem into an equivalent maximization problem
 - So we can construct the submodular function for the greedy algorithm

$$\min O_1 = \sum_{s \in S_n} y_s$$

$$\min O_2 = \sum_{s \in S} y_s$$



$$\max O_1 = C_{\max} - \sum_{s \in S_n} C(\mathcal{N}_s)$$

$$\max O_2 = C_{\max} - \sum_{s \in S} C(\mathcal{N}_s)$$

$$S.t. \begin{cases} \sum_{s \in S} x_n^s = 1, & \forall n \in N \\ \sum_{n \in N} x_n^s \leq R_s, & \forall s \in S \\ \sum_{n \in N} \sum_{n' \in N} x_n^s \cdot x_{n'}^{s'} \cdot C_{n,n'}^t \leq P_{s,s'}, & \forall s, s' \in S, t \in T \\ y_s \geq x_n^s, & \forall s \in S, n \in N \\ x_n^s \in \{0, 1\}, & \forall s \in S, n \in N \\ y_s \in \{0, 1\}, & \forall s \in S \end{cases} \quad (1)$$

$$S.t. \begin{cases} \sum_{s \in S} x_n^s = 1, & \forall n \in N \\ \sum_{n \in N} x_n^s \leq R_s, & \forall s \in S \\ \sum_{n \in N} \sum_{n' \in N} x_n^s \cdot x_{n'}^{s'} \cdot C_{n,n'}^t \leq P_{s,s'}, & \forall s, s' \in S, t \in T \\ x_n^s \in \{0, 1\}, & \forall s \in S, n \in N \end{cases} \quad (3)$$

Titan: Algorithm Design

- Solve the converted problem with a **submodular-based greedy algorithm**
- **Search the feasible worker set** for each machine to guarantee bandwidth constraint
- **Merge the worker set with submodular function** to minimize the number of used idle machines
- **Update and merge the worker set** to minimize the number of used machines

Algorithm 1 Search for Feasible Worker Set

- 1: **Step 1: Initialization**
 - 2: Let feasible worker set $A(s) = \emptyset$ for machine s .
 - 3: Let available worker set $N_c = N - \mathcal{N}$.
 - 4: **Step 2: Iterative update feasible worker sets according to collective communication**
 - 5: Use $\mathcal{C}_{f,f'}$ to denote the existing communication overhead between machines s and s' .
 - 6: **for** $n \in N_c$ **do**
 - 7: **for** $n' \in \mathcal{N}$ **do**
 - 8: **for** $t \in T$ **do**
 - 9: **if** $C_{n,n'}^t + \mathcal{C}_{f,f'} \leq P_{s,s'}$ **then**
 - 10: $A(s) \leftarrow A(s) + n$
 - 11: **end if**
 - 12: **end for**
 - 13: **end for**
 - 14: **end for**
 - 15: Output the feasible worker set $A(s)$ for machine s .
-

Algorithm 3 The Overall Algorithm

- 1: **Step 1: Minimizing the Number of Deployed New Racks**
 - 2: Initiate $\mathcal{N}_s, \forall s \in S_n$ by randomly distributing workers.
 - 3: Initiate $\Phi_n \leftarrow \emptyset$.
 - 4: Calculate Φ_n on idle machine set S_n with Alg. 2.
 - 5: **Step 2: Minimizing the Total Number of Deployed Racks**
 - 6: Initiate $\mathcal{N}_s, \forall s \in S_f \cup \Phi_n$ by randomly distributing workers.
 - 7: Initiate $\Phi \leftarrow \emptyset$.
 - 8: Calculate Φ on machine set $S_f \cup \Phi_n$ with Alg. 2.
 - 9: **Step 3: Determining the Deployment of Workers**
 - 10: **for** $\mathcal{N}_s \in \Phi$ **do**
 - 11: Set $x_n^s = 1, \forall n \in \mathcal{N}_s, s \in S$.
 - 12: **end for**
-

Titan: Algorithm Design

- Solve the converted problem with a **submodular-based greedy algorithm**
- **Search the feasible worker set** for each machine to guarantee bandwidth constraint
- **Merge the worker set with submodular function** to minimize the number of used idle machines
- **Update and merge the worker set** to minimize the number of used machines
- Tight approximate ratio (1-1/e)

Algorithm 2 Submodular-based Algorithm

- 1: **Step 1: Initialization**
 - 2: Initiate $\mathcal{N}_s, \forall s \in S$ by randomly distributing workers.
 - 3: Initiate $\Phi \leftarrow \emptyset$.
 - 4: **Step 2: Iterative Merging Worker Subsets**
 - 5: **while** $|\Phi| \leq K - 1$ **do**
 - 6: Set $tmp \leftarrow 0, opt \leftarrow 0$
 - 7: **for** $s \in S_n$ **do**
 - 8: **for** $n \in \mathcal{N}_s - \Phi$ **do**
 - 9: $tmp \leftarrow H(\Phi \cup \{n\})$
 - 10: **if** $tmp > opt$ **then**
 - 11: $opt \leftarrow tmp, \mathcal{N}^* \leftarrow \mathcal{N}^* + \{n\}$
 - 12: **end if**
 - 13: **end for**
 - 14: **end for**
 - 15: $\Phi \leftarrow \Phi + \mathcal{N}^*$
 - 16: Update the feasible worker sets based on the bandwidth constraint of Eq. (5) with Alg. 1.
 - 17: **end while**
 - 18: Deploy the remaining workers on one machine (*i.e.*, $\Phi \leftarrow \Phi + \{N - \bigcup_{\mathcal{N} \in \Phi} \mathcal{N}\}$).
-

$$\sum_{n \in \mathcal{N}_s} \sum_{n' \in \mathcal{N}_{s'}} C_{n,n'}^t \leq P_{s,s'}, \forall s' \in S, t \in T \quad (5)$$

Evaluation: Setup

Simulation Topology

- Fat-tree topology with 64 racks, each of which contains 8 machines
- Each machine is equipped with 8 GPUs
- All connected with 100Gbps links

Real-world Traces

- Microsoft cluster: 2-month trace with 69742 jobs
- Shanghai AI lab cluster: 6-month trace with 880740 jobs

Evaluation: Setup

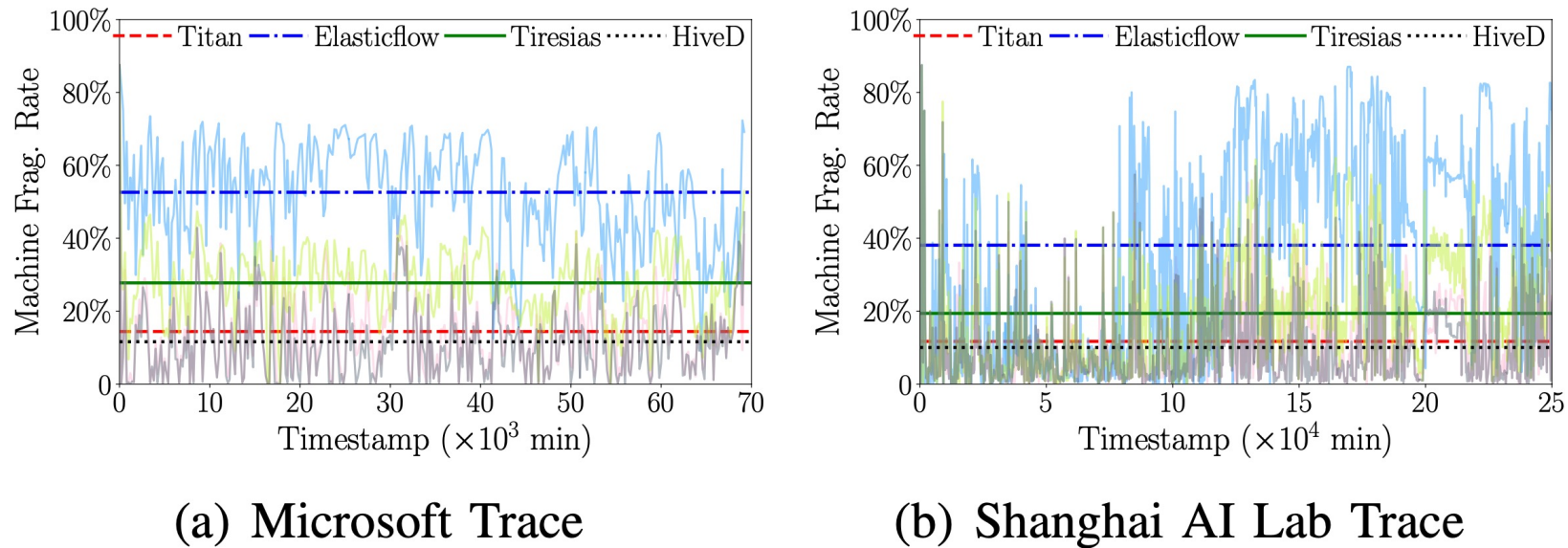
Benchmark

- **Elasticflow**¹ : consolidate the placed GPUs so that the job is allocated with the highest possible bandwidth between its workers
- **HiveD**² : prioritize placing workers in fragmented machines to reduce the machine fragmentation of the cluster
- **Tiresias**³ : minimize the total network traffic and balance the network load across machines in the cluster, by profiling the characteristics of different models

3. Gu J, Chowdhury M, Shin K G, et al. Tiresias: A {GPU} cluster manager for distributed deep learning[C]//16th USENIX Symposium on Networked Systems Design and Implementation (NSDI 19). 2019: 485-500.

Evaluation: Fragmentation Rate

- Record the ratio of non-idle machines in the cluster, after a job arrives



(a) Microsoft Trace

(b) Shanghai AI Lab Trace

Fig. 4: Machine Fragmentation Rate vs. Timestamp

- Titan reduces machine fragmentation rate by **38.1%**

Evaluation: Comm. Overhead

- Use the HD algorithm and calculate the total communication overhead

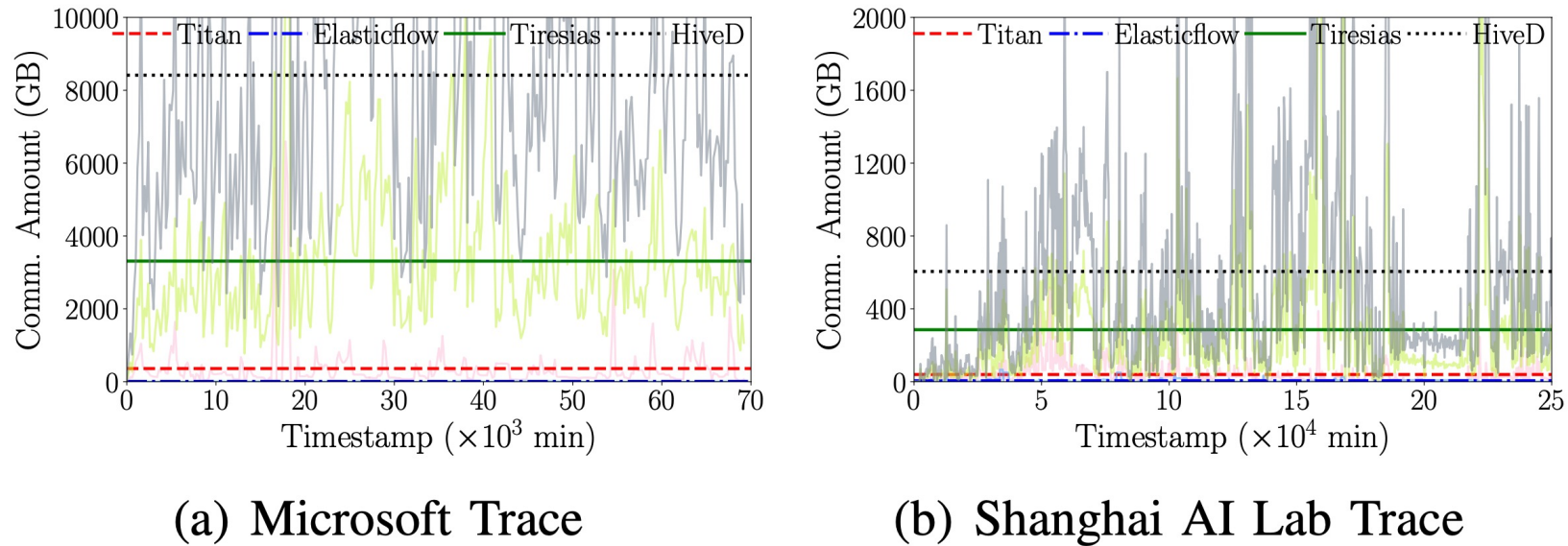


Fig. 5: Communication Overhead vs. Timestamp

- Titan reduces communication overhead of the cluster by **76.4%**

Evaluation: Profit

- Profit relates to No. of jobs, price of the job, No. of used machines and the duration of each job (See the manuscript for details)

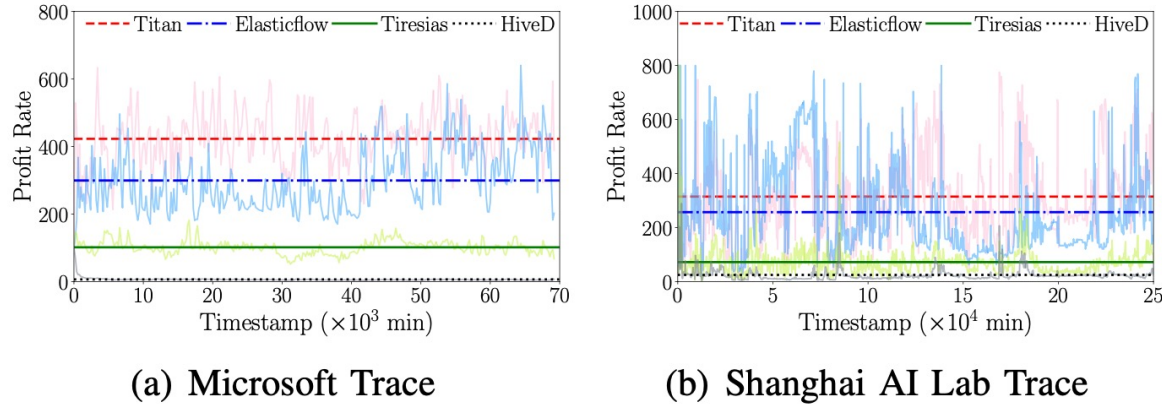


Fig. 6: Profit Rate vs. Timestamp

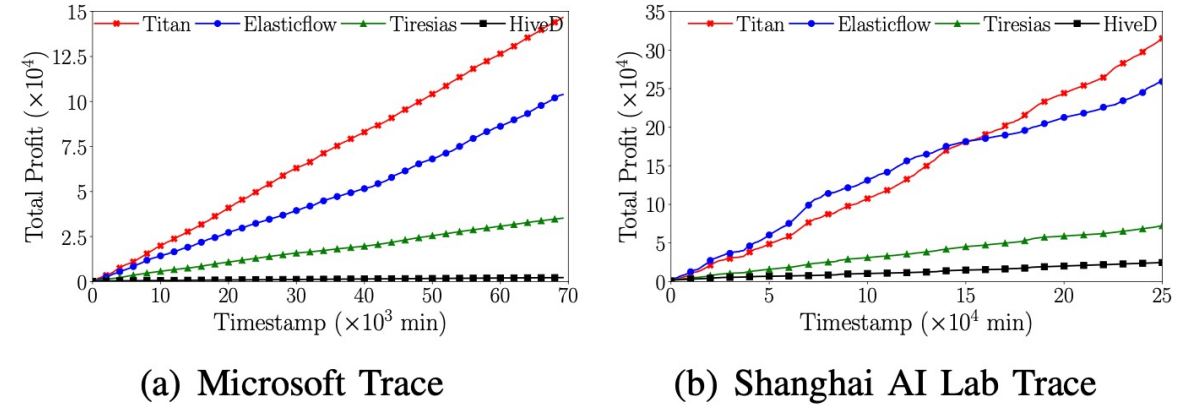


Fig. 7: Total Profit vs. Timestamp

- Titan improves the total profit by **41.2%~65X**

Summary

Goal

- Minimize the resource fragmentation rate of GPU cluster with non-idle machine-aware worker placement

Challenges

- Collective Communication constraint
- Resource fragmentation limitation
- Model a multi-subjective non-linear problem

Solution

- Submodular-based greedy algorithm with a tight approximation ratio

Thank you!

IEEE International Conference on Network Protocols (ICNP 2024)

Committee, Reviewers, Volunteers

My Advisors and Collaborators!

Jin Fang

fangjin98@mail.ustc.edu.cn

www.fangjin.site