

ATA: Adaptive Task Allocation for Efficient Resource Management in Distributed Machine Learning

Artavazd Maranjyan, El Mehdi Saad, Peter Richtárik, Francesco Orabona

Motivation

- Scenario:** Running Minibatch Stochastic Gradient Descent (SGD) on a heterogeneous cluster of 1010 workers.
- Batch Size:** 10.
- Fastest (but wasteful) Method:** Request gradients from all workers asynchronously and use only the first 10 responses that arrive.
- Problem:** This method discards the work of at least 1000 workers, wasting their computational resources.

Problem setup

- Workers: n , Stochastic speeds X_i
- Batch Size: B
- Naive Policy:** Request from all n , use first $B \rightarrow$ Wastes $n-B$ jobs
- Goal:** Find optimal request set size B to balance speed vs. efficiency
 $\mathcal{A} := \{a \in \mathbb{N}^n : \|a\|_1 = B\}$

Objective

$$C_K := \sum_{k=1}^K \mathbb{E}[C(a^k)]$$

$$C(a^k) := \max_{i \in \text{supp}(a^k)} \sum_{u=1}^{a_i^k} X_i^{k,u}, \quad a^k \in \mathcal{A}$$

Sub-exponential random variables

$$\|X_i - \mu_i\|_{\psi_1} \leq \alpha, \quad \text{for all } i \in [n]$$

$$\|X\|_{\psi_1} := \inf\{C > 0 : \mathbb{E}[\exp(|X|/C)] \leq 2\}$$

Confidence Interval

$$s_i^k = \max\{\hat{\mu}_i^k - \text{conf}(i, k), 0\}$$

$$\text{conf}(i, k) = \begin{cases} 2\alpha \left(\sqrt{\frac{\ln(2k^2)}{K_i^k}} + \frac{\ln(2k^2)}{K_i^k} \right), & K_i^k \geq 1, \\ +\infty, & K_i^k = 0 \end{cases}$$

Theoretical Results

Proxy loss

$$\ell(a, \mu) := \max_{i \in [n]} a_i \mu_i$$

$$\ell(a, \mu) \leq \mathbb{E}[C(a)] \leq (1 + 4\eta \ln(B)) \ell(a, \mu)$$

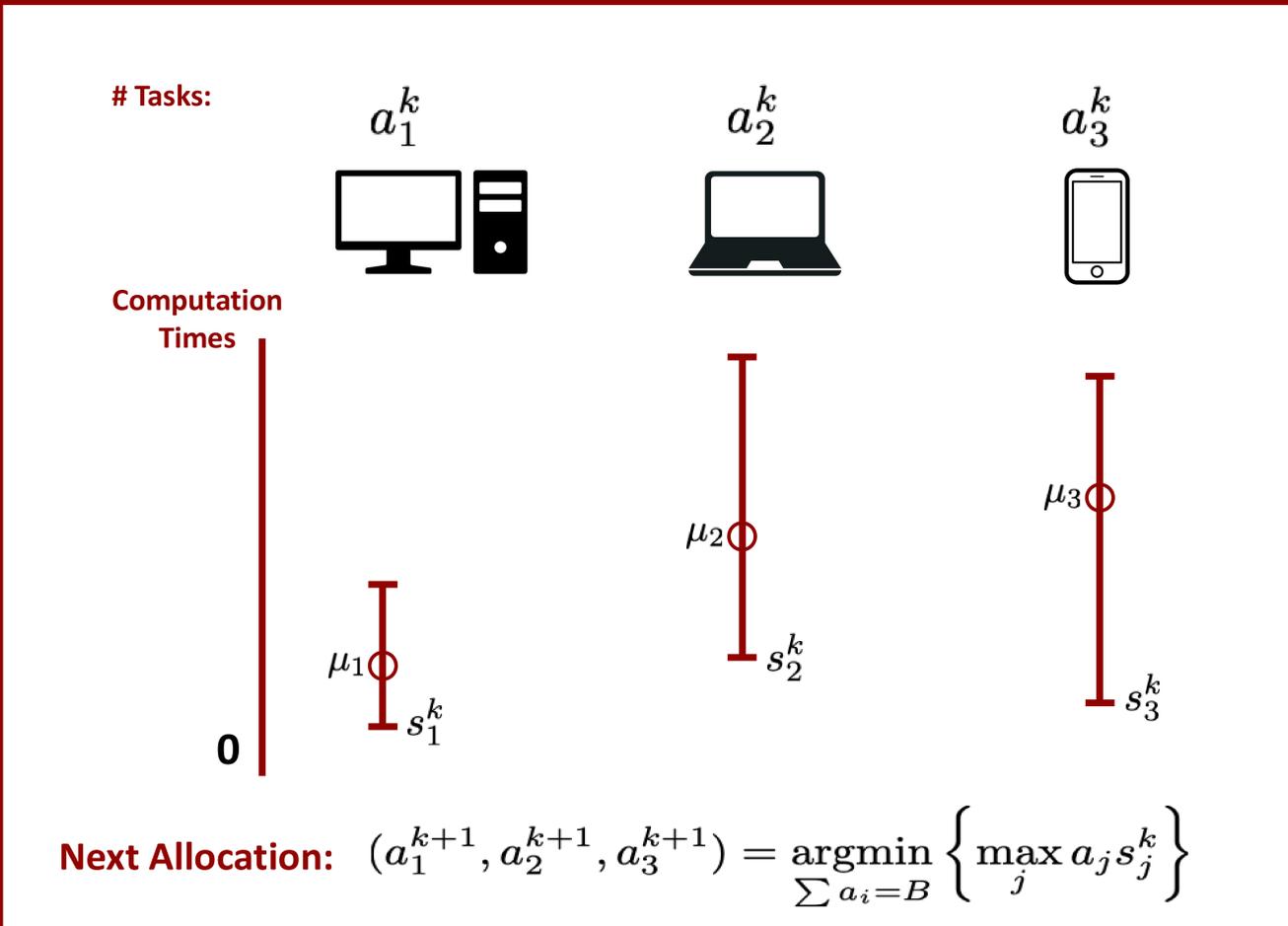
$$\eta := \max_{i \in [n]} \frac{\alpha_i}{\mu_i}$$

Guarantees

$$C_K \leq (1 + 4\eta \ln(B)) C_K^* + \mathcal{O}(\ln K)$$

$$C_K^* := K \mathbb{E}[C(a^*)], \quad a^* \in \underset{a \in \mathcal{A}}{\text{argmin}} \mathbb{E}[C(a)]$$

Multiple devices, unknown speeds? No problem! ATA learns and adapts to minimize computation while maintaining fast distributed ML training.



Experiments

Dataset: CIFAR-100
 Network: CNN with 3 convolutional layers and 2 fully connected layers
 Optimizer: Adam with constant learning rate of 8×10^{-5}

$$B = 23, \quad n = 51$$

$$X_i \sim 29i + \text{Exp}(29i), \quad \text{for all } i \in [n]$$

ATA: Empirical

$$\hat{s}_i^k = \hat{\mu}_i^k \max \left\{ 1 - 2\eta \left(\sqrt{\frac{\ln(2k^2)}{K_i^k}} + \frac{\ln(2k^2)}{K_i^k} \right), 0 \right\}$$

Baselines

FTA: Fixed Task Allocation

GTA: Greedy Task Allocation (asynchronous batch collection)

UTA: Uniform Task Allocation

