# Powering Multi-Task Federated Learning with Competitive GPU Resource Sharing

**Yongbo Yu, Fuxun Yu, Zirui Xu, Xiang Chen**
**Department of Electrical and Computer Engineering, George Mason University**

GEORGE MASON UNIVERSITY
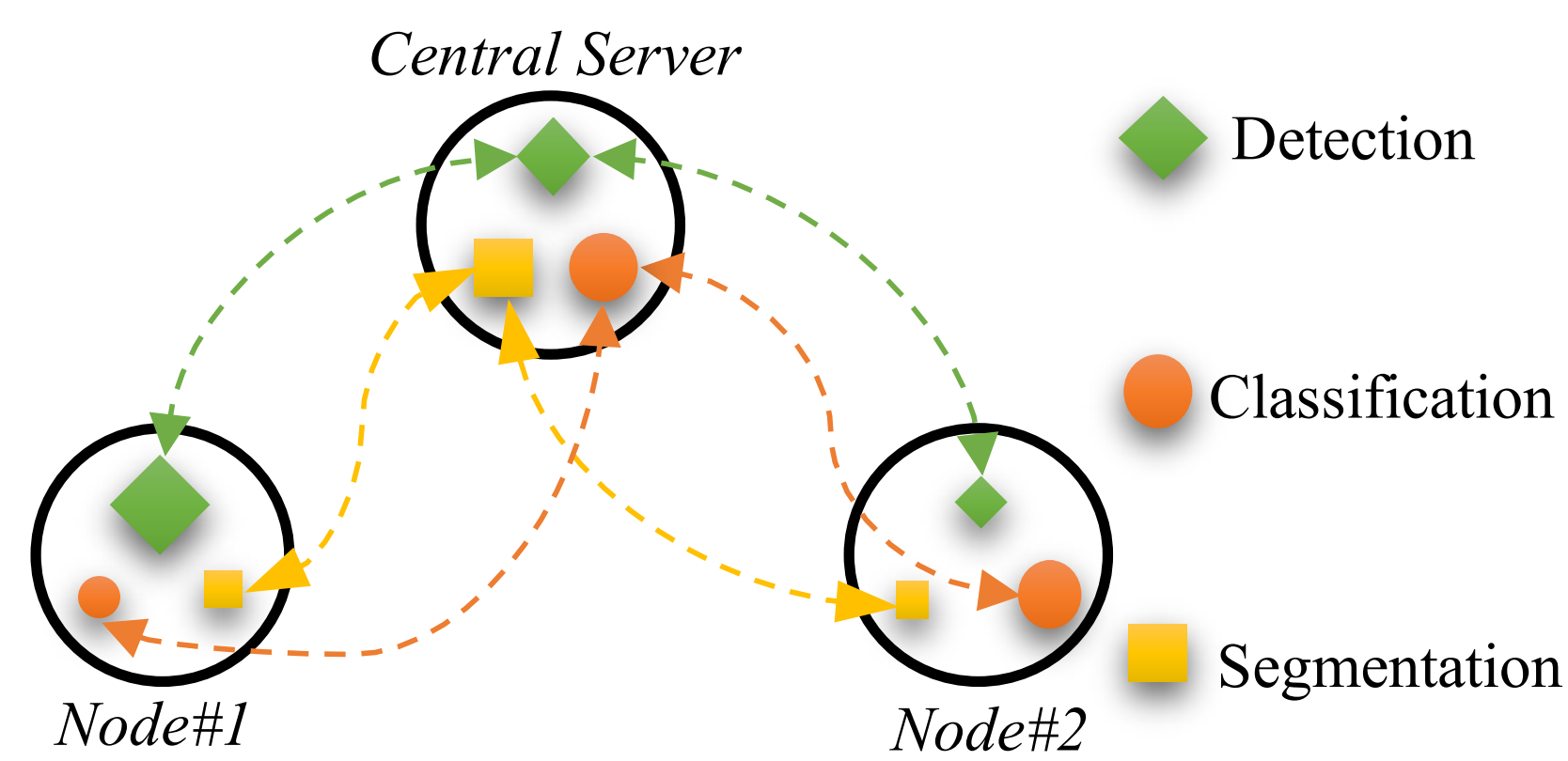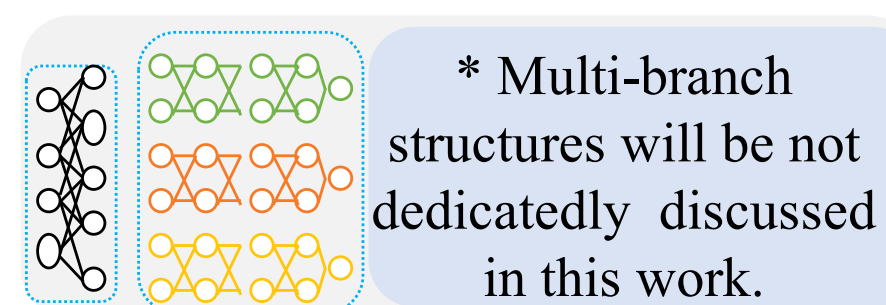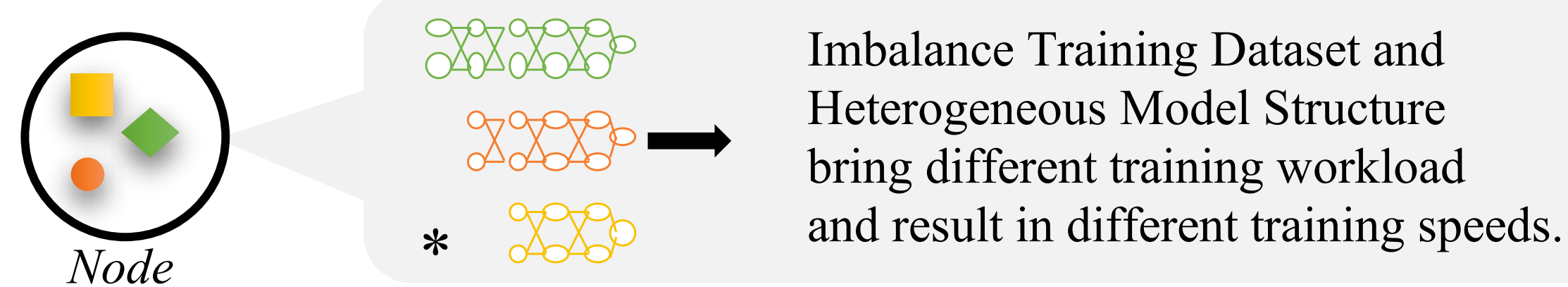
Intelligence-Fusion Laboratory
www.if-lab.org

## 1. Multi-Task Federated Learning

*Multi-Task Federated Learning*
- ❖ Different Learning tasks with
  Imbalance Datasets
  Heterogeneous Models
- ❖ Learning Speed
  Multi-Task Synchronization



*Central Server*

Detection
Classification
Segmentation

Node#1    Node#2

*Multi-Task Deployment with Parallel Model Structures*



*Node*

Imbalance Training Dataset and Heterogeneous Model Structure bring different training workload and result in different training speeds.

\* Multi-branch structures will be not dedicatedly discussed in this work.

*New Challenges for Computing and Coordination*
- ❖ Multi-Task Parallel Computing Deployment
- ❖ Different Learning Speed Multi-Task Coordination

*We propose a full-stack multi-task FL optimization scheme:*
- intra-device GPU scheduling with a competitive resource sharing scheme;
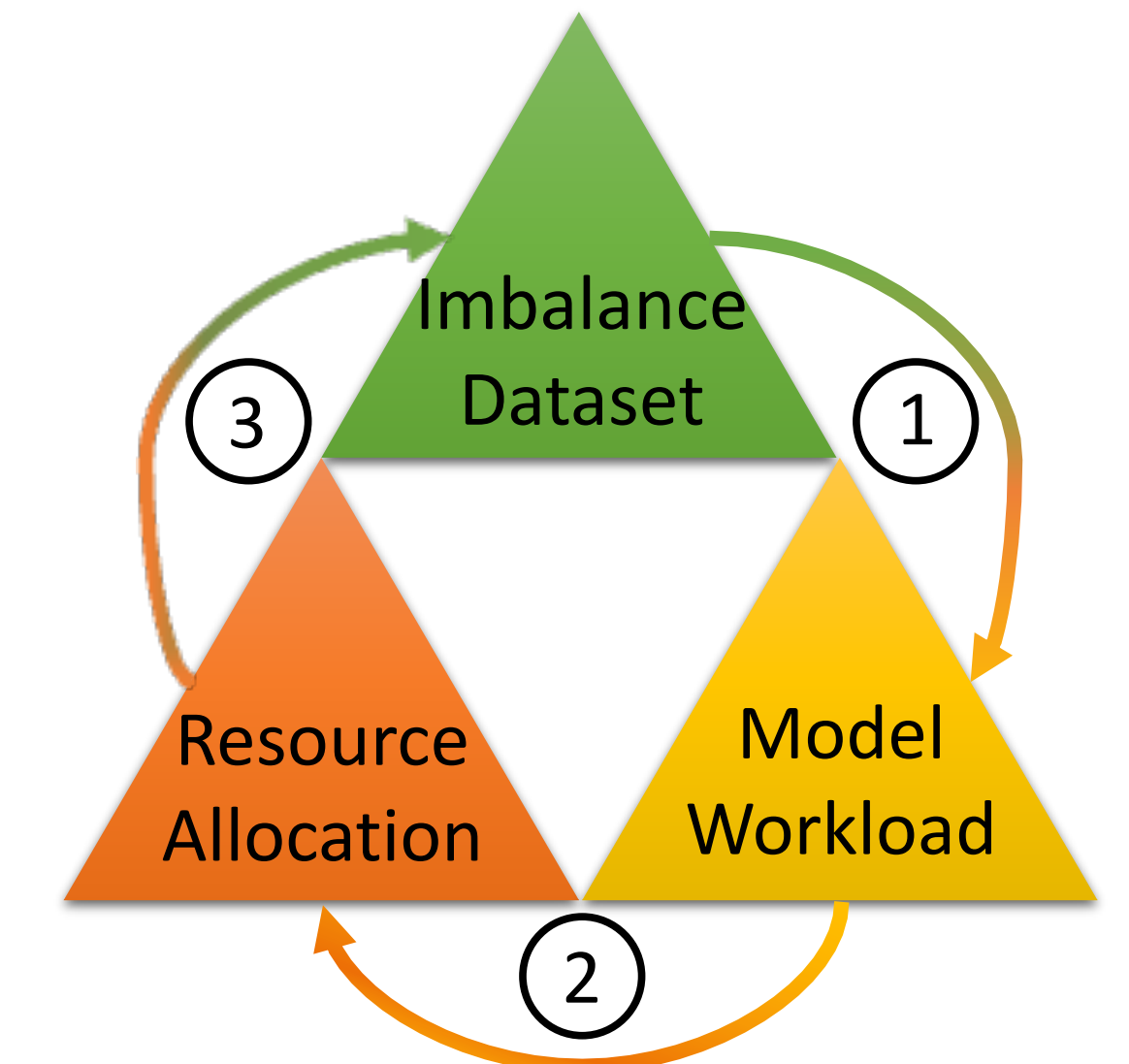- inter- device FL coordination with realistic GPU runtime synchronization.

## 2. Intra-Device Multi-Task Resource Allocation

*Multi-Task Parallel Deployment Mechanism*
How to achieve effective multi-task dedicated GPU resource allocation.

**Challenges:**
- Identify the parallel computing issue;
- Establish the relationship between
- performance and resource allocation.



CUDA MPS Control
Volta MPS Controller
Streaming Multiprocessors (SMs)

**Comprehensive Analysis with Competitive Resource Sharing Cases**

Task-A    Task-B
Competitive Resource Sharing    Contention Delay



① No share — *under utilization*

② 20% share — *higher util. / light competition* — -25ms

③ 60% share — *severe competition* — -10ms

④ 100% share — *sequential execute* — -5ms

**Isolated Spatial Resource Allocation**
Exclusive resource assignment causes under-utilization.

**Spatial Sharing with Light Competition**
Higher resource allocation flexibility and better resource utilization.

**Sharing with Excessive Competition**
Leading to resource competition and considerable contention overhead.

**Extreme Contention Kills Parallelism**
Push back to temporal scheduling when fully sharing the resource without partitioning.

**Optimization Insights!**
Controlling an appropriate degree of resource competition and sharing is the key to achieve the optimal performance.

**"*Virtual Resource*"** management to establish the relationship between performance and resource allocation.

*Definition of Virtual Resource* :
Virtual resource is a number between $(0\% \sim 100\%) \times (I$ of tasks$)$.

As shown in the right figure, the four points correspond to the above four cases.

We propose a machine learning approach to estimate the GPU **throughput P and achieve the maximum throughput.**



*Virtual Resource*

*Physical Resource*    ②  ③
①    *Throughput*    ④

*Competition*

0%    100%    200%

## 3. Inter-device Multi-Task FL coordination

**Coordination Design Motivation:**
We rethink the FL coordination from a GPU scheduling perspective under imbalance dataset, model workload and resource allocation.

**Our Objectives:**
(1) Each device could be fully utilized during each synchronization cycle;
(2) Maximize the overall GPU throughput P.



Imbalance Dataset
Resource Allocation    Model Workload

$$\begin{cases} \text{LocalTraining}: \min_{\{W_{0,j}\}_{j=0}} Loss(D_{0,j}, W_{0,j}), \\ \text{Fusing}: W_{i,j}^{\{k^{th}cycle\}} = \sum_{j=0}^{J} \frac{|D_{0,j}|}{\sum_{k=0}^{J}|D_{0,k}|} W_{i,j}^{\{k-1^{th}cycle\}} \end{cases}$$

*Multi-Task Federated Learning with J devices in a federation cluster, each device has I tasks*

$$\begin{cases} \text{Objective } 1: \min \sum_i \sum_j \frac{|D_i|}{|D_j|} - \frac{O_i}{O_j}, \\ \text{Objective } 2: \max \sum_i P_1, \dots, P_i. \end{cases}$$

**Tri-party Optimization to Achieve Objective 1:**
1. Larger Dataset D need larger Workload O;
2. Larger Workloads O compete for more Resources;
3. More Resources train more data each cycle.

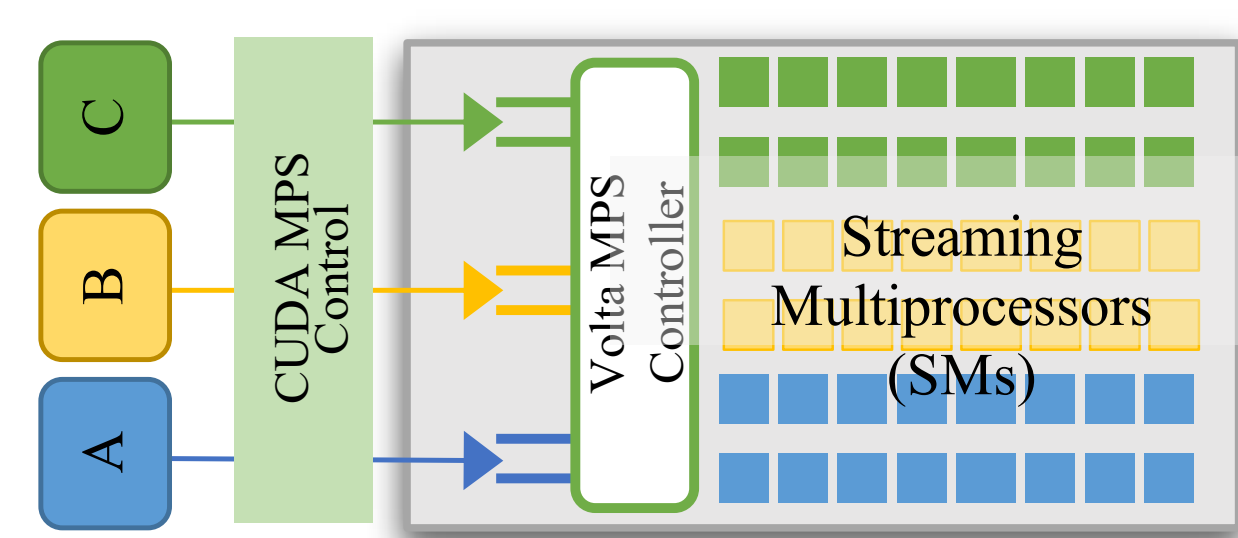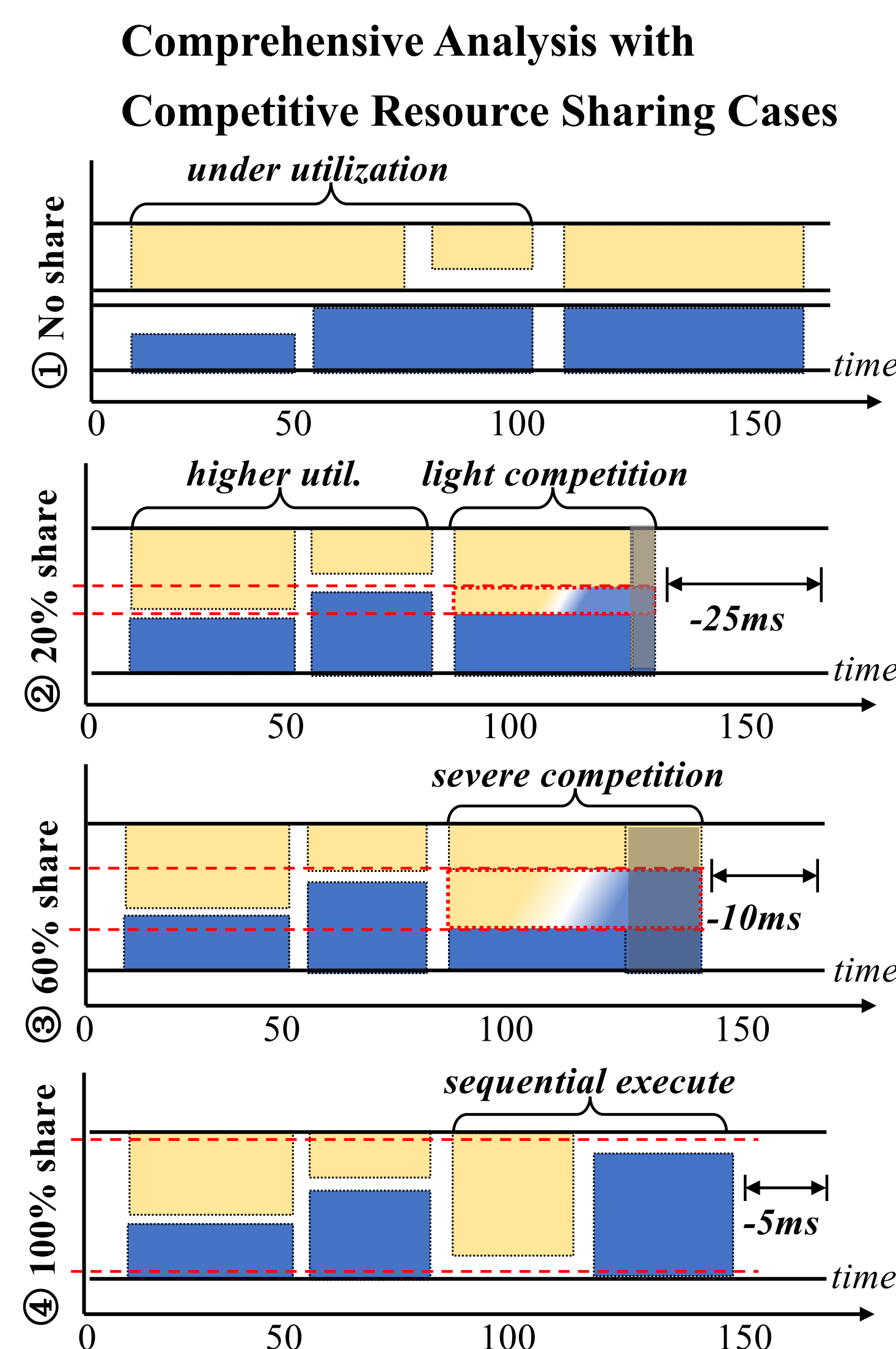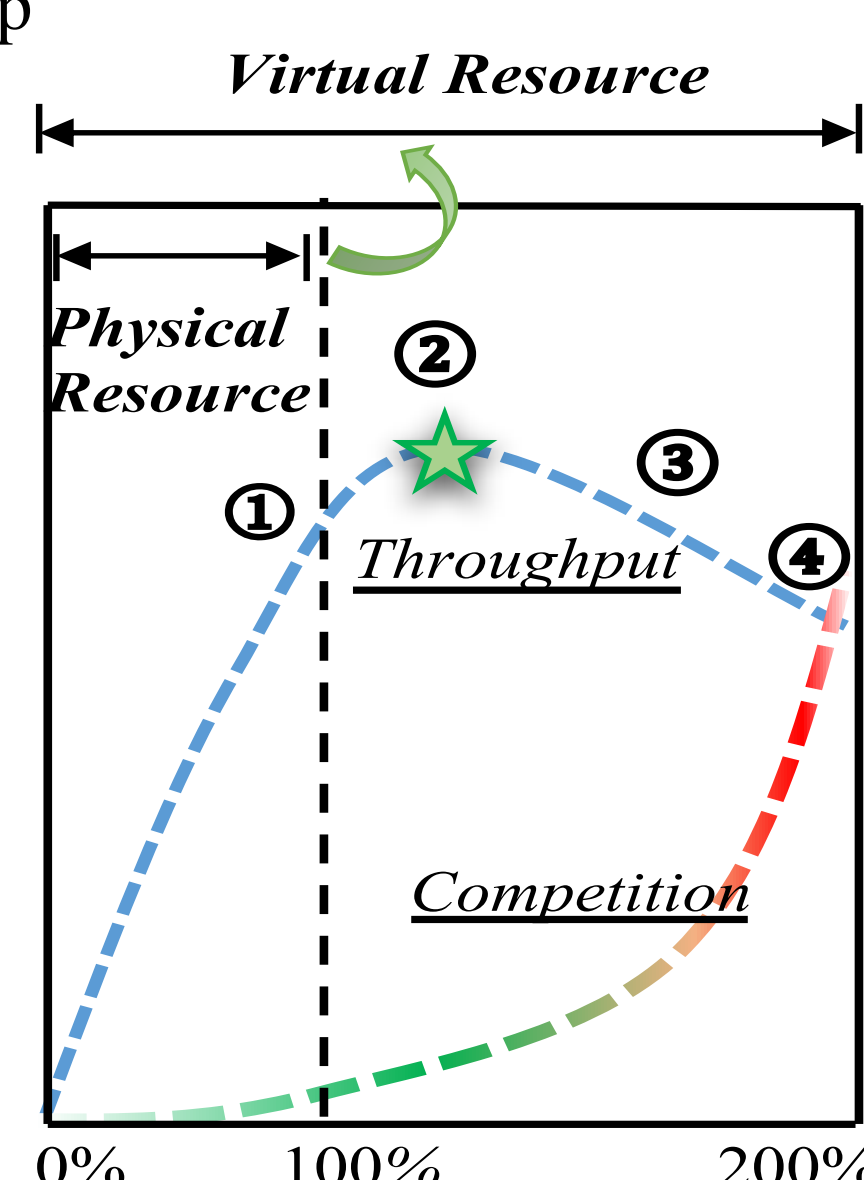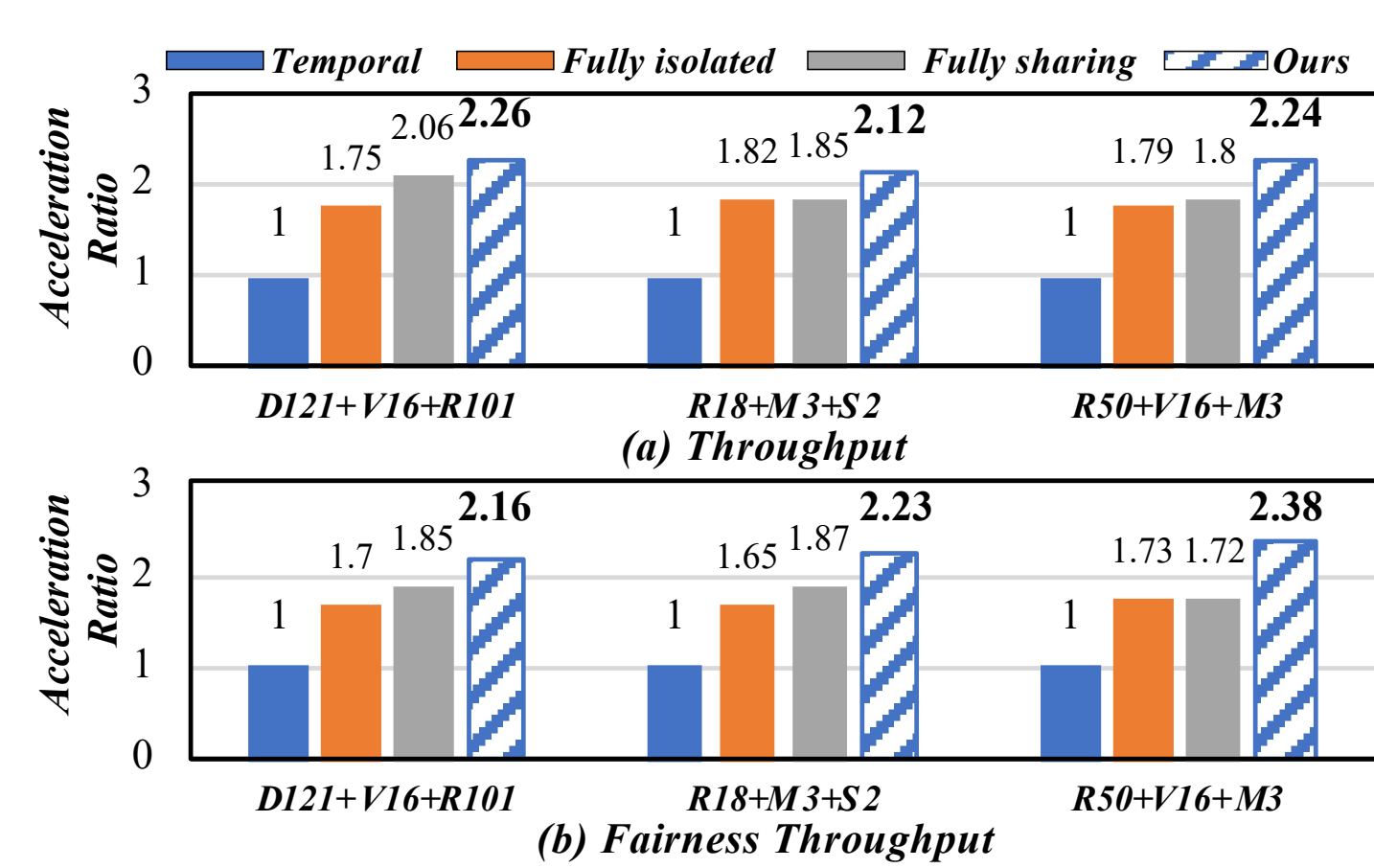We use greed optimization method to find the optimal batch size and resource allocation.

## 4. Experimental Results



Temporal    Fully isolated    Fully sharing    Ours

(a) Throughput
D121+V16+R101: 1, 1.75, 2.06, 2.26
R18+M3+S2: 1, 1.82, 1.85, 2.12
R50+V16+M3: 1, 1.79, 1.8, 2.24

(b) Fairness Throughput
D121+V16+R101: 1, 1.7, 1.85, 2.16
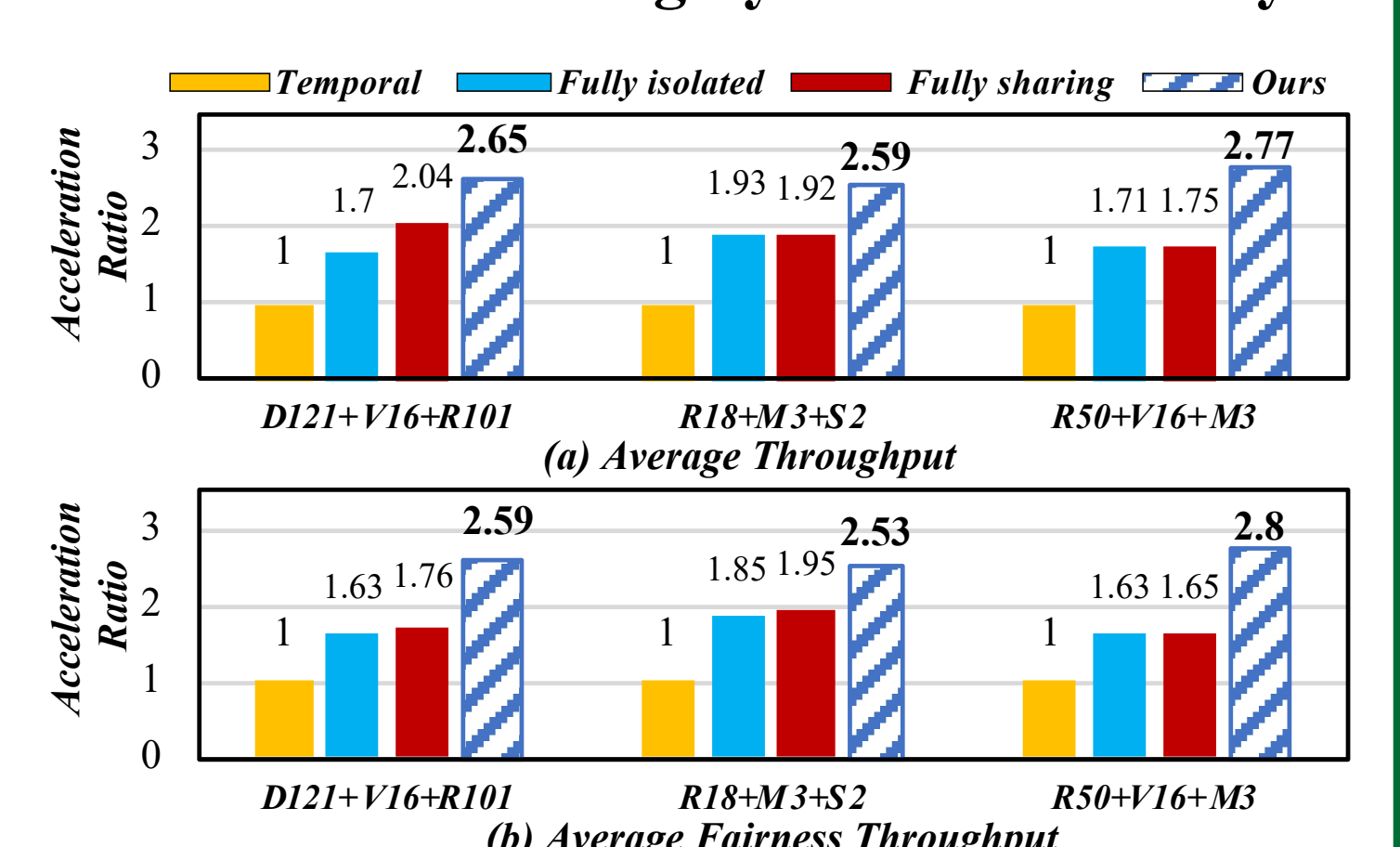R18+M3+S2: 1, 1.65, 1.87, 2.23
R50+V16+M3: 1, 1.73, 1.72, 2.38

**Intra-Device Multi-Task GPU Throughput**

We use a FL system with several devices, each devices have three tasks with different model structure and imbalance data volume.

**Joint optimization of workload and resource allocation allows our method to utilize GPU resources during the whole synchronization cycle.**

We can apply our resources allocation method on various multi-task training scenarios under multiple DNN models' combinations.

**Inter-Device Average Throughput in a Federated Learning Synchronization Cycle**



Temporal    Fully isolated    Fully sharing    Ours

(a) Average Throughput
D121+V16+R101: 1, 1.7, 2.04, 2.65
R18+M3+S2: 1, 1.93, 1.92, 2.59
R50+V16+M3: 1, 1.71, 1.75, 2.77

(b) Average Fairness Throughput
D121+V16+R101: 1, 1.63, 1.76, 2.59
R18+M3+S2: 1, 1.85, 1.95, 2.53
R50+V16+M3: 1, 1.63, 1.65, 2.8

## 5. Contribution and Conclusion

- We analyze the competitive resource sharing mechanism propose an intra-device multi-task dedicated GPU scheduling method.
- We further bring competitive resource sharing mechanism into the inter-device FL cluster and rethink the FL coordination from a GPU scheduling perspective.

Beyond some FL methods which resolve the heterogeneous problems from algorithmic level, FL's heterogeneous research should also dive into computation level due to imbalance dataset, heterogeneous model and different learning speed per task.

## 6. Reference

[1] Shinpei Kato et al.TimeGraph: GPU scheduling for real-time multi-tasking environments. In Proceedings of the USENIX ATC. 17–30.
[2] TianLietal.2020. Federated learning: Challenges, methods, and future directions. IEEE Signal Processing Magazine 3 (2020), 50–60.
[3] Nvidia. MPS. https://docs.nvidia.com/deploy/pdf/CUDA_Multi_Process_Service_Overview.pdf
[4]. Mehdi Salehi Heydar Abad et al. Hierarchical federated learning across heterogeneous cellular networks. In Proceedings of the ICASSP IEEE International Conference on Acoustics, Speech and Signal Processing. IEEE, 8866–8870.