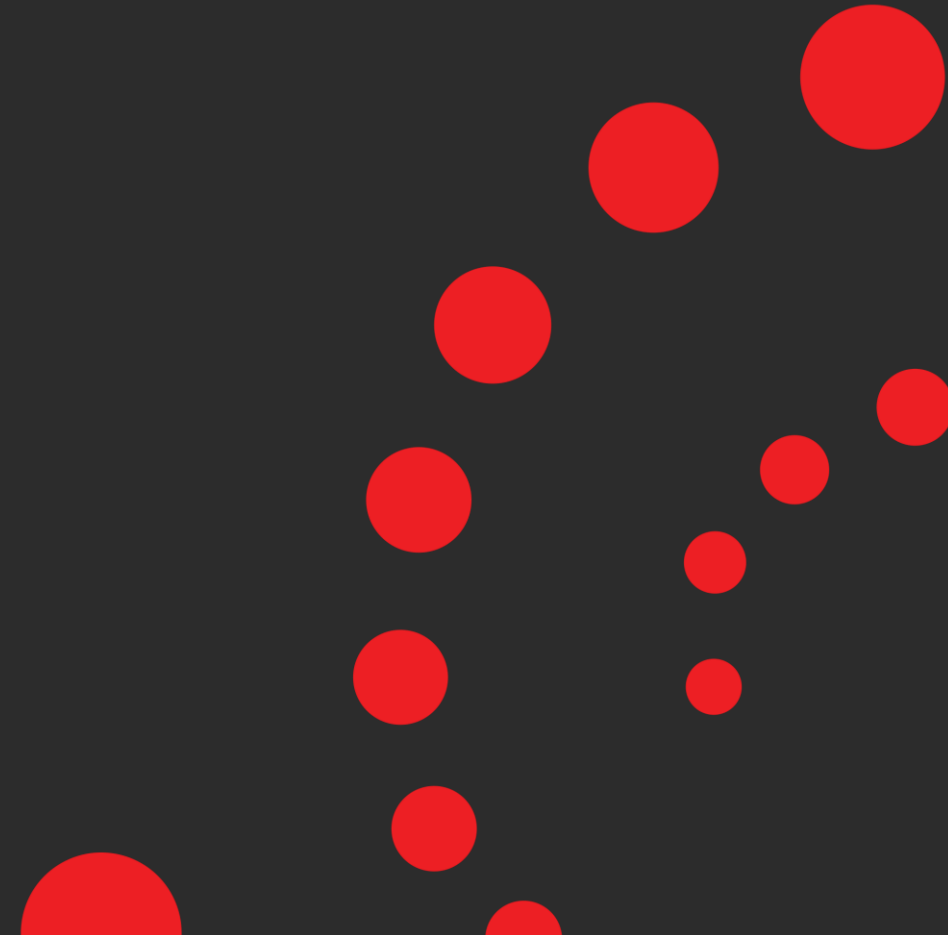# HPDC'22 Summary
June 27 – July 1 2022

Norbert Egi
Infrastructure Software

# Agenda

## *Session 1: Data Centers and HPC Systems*

- DAOS: Data Access-aware Operating System

- FPVM: Towards a Floating Point Virtual Machine

- Lifting and Dropping VMs to Dynamically Transition Between Time- and Space-sharing for Large-Scale HPC Systems

## *Session 2: HPC Memory, I/O, and Storage Systems*

- Access Patterns and Performance Behaviors of Multi-layer Supercomputer I/O Subsystems under Production Load

- NVMe-oAF: Towards Adaptive NVMe-oF for IO-Intensive Workloads on HPC Cloud

- Capri: Compiler and Architecture Support for Whole-System Persistence

## *Session 3: Reliability and Scheduling*

- Understanding Memory Failures on a Petascale Arm System

- **SchedInspector: A Batch Job Scheduling Inspector Using Reinforcement Learning**

- Holmes: SMT Interference Diagnosis and CPU Scheduling for Job Co-location

- TLPGNN: A Lightweight Two-Level Parallelism Paradigm for Graph Neural Network Computation on GPU

**FUTUREWEI**
Technologies

# Agenda

## Session 4: HPC Algorithms

- TAC: Optimizing Error-Bounded Lossy Compression for Three Dimensional Adaptive Mesh Refinement Simulations
- Communication-aware Sparse Patterns for the Factorized Approximate Inverse Preconditioner
- Ultra-fast Error-bounded Lossy Compression for Scientific Dataset
- Optimizing the Bruck Algorithm for Non-uniform All-to-all Communication

## Session 5: HPC Toolchains, Traces, and More

- SciStream: Architecture and Toolkit for Data Streaming between Federated Science Instruments
- Machine Learning Assisted HPC Workload Trace Generation for Leadership Scale Storage Systems
- PROV-IO: An I/O-Centric Provenance Framework for Scientific Data on HPC Systems

## Session 6: Cloud Computing and Machine Learning

- Locality-aware Load-Balancing For Serverless Clusters
- Practical Efficient Microservice Autoscaling with QoS Assurance
- **Hare: Exploiting Inter-job and Intra-job Parallelism of Distributed Machine Learning on Heterogeneous GPUs**
- **Efficient Design Space Exploration for Sparse Mixed Precision Neural Architectures**

**FUTUREWEI** Technologies

# Selected Papers of Key Interest

➤ SchedInspector: A Batch Job Scheduling Inspector Using Reinforcement Learning

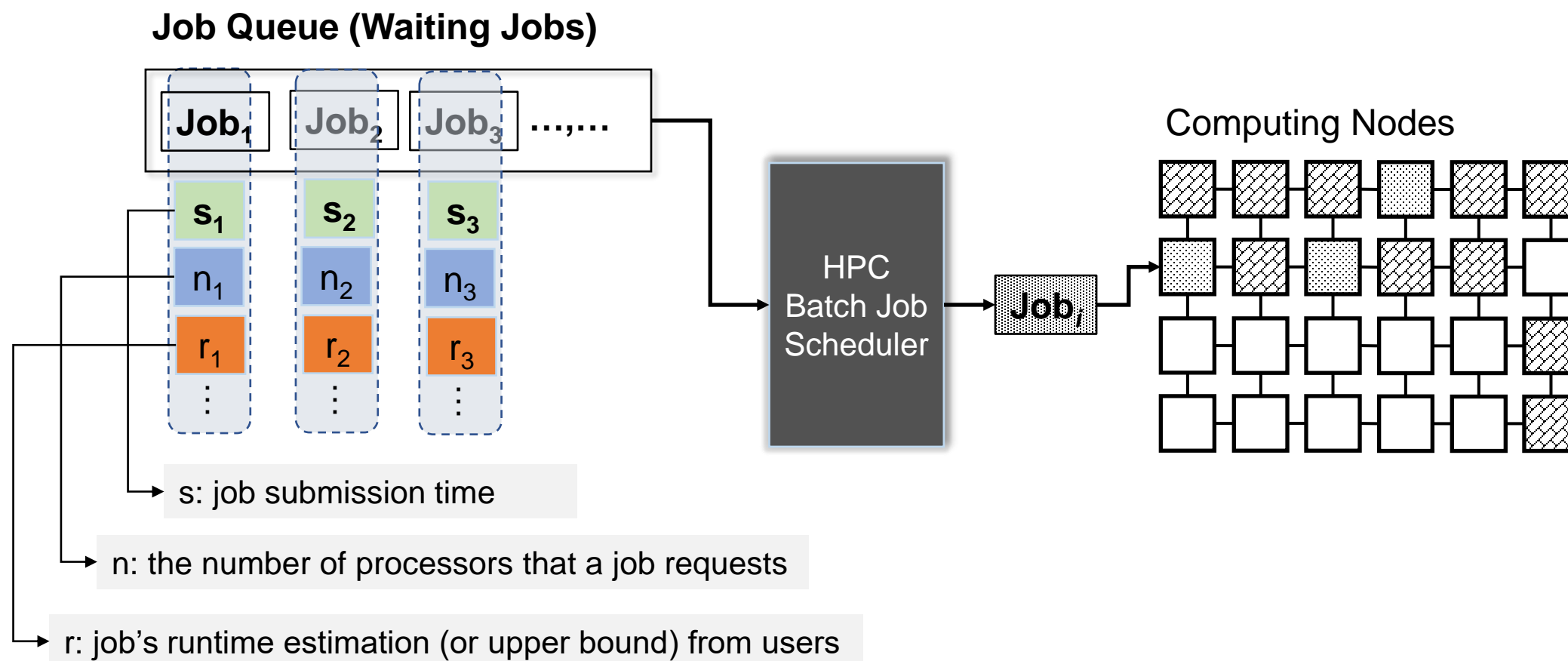➤ Efficient Design Space Exploration for Sparse Mixed Precision Neural Architectures

# SchedInspector: A Batch Job Scheduling Inspector Using Reinforcement Learning
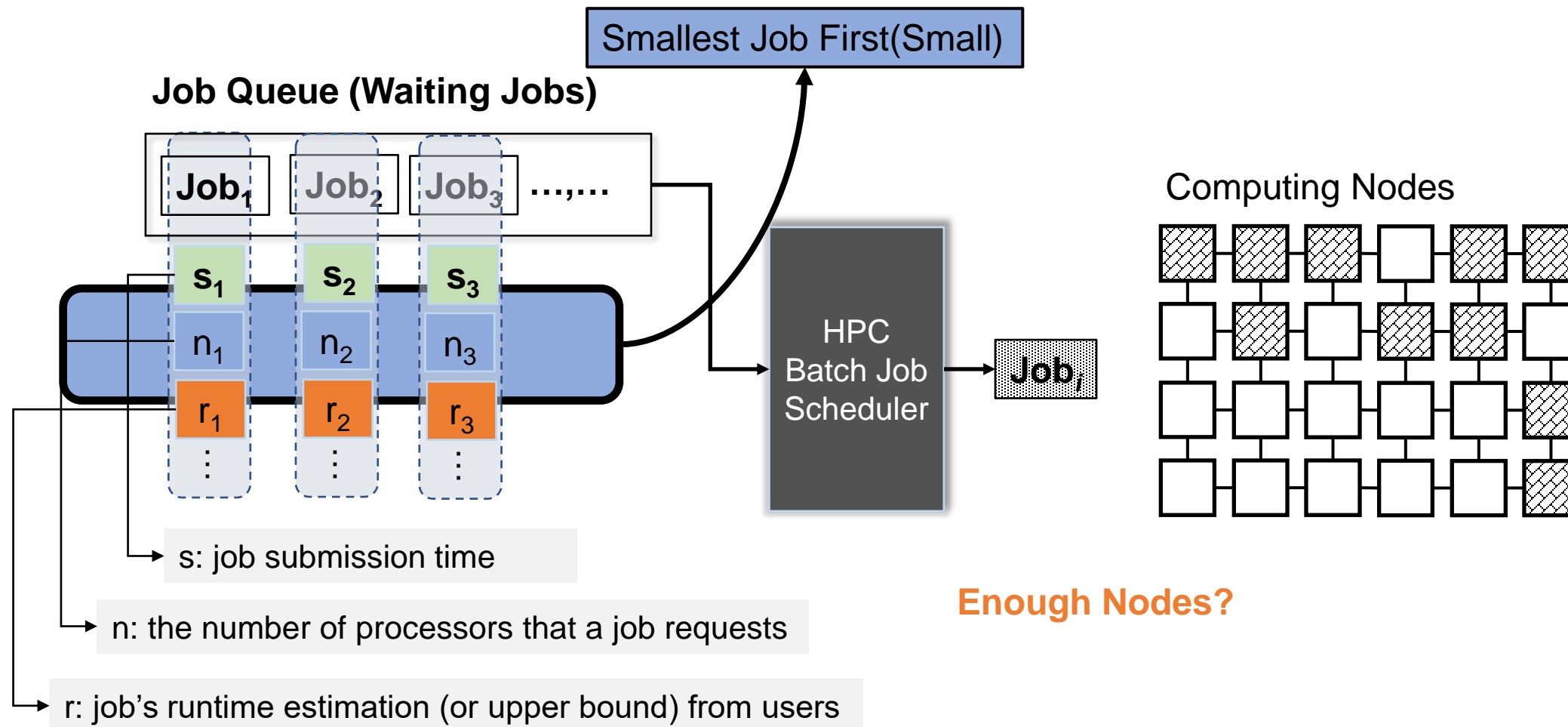
Di Zhang[1], Dong Dai[1], Bing Xie[2]

[1]University of North Carolina at Charlotte
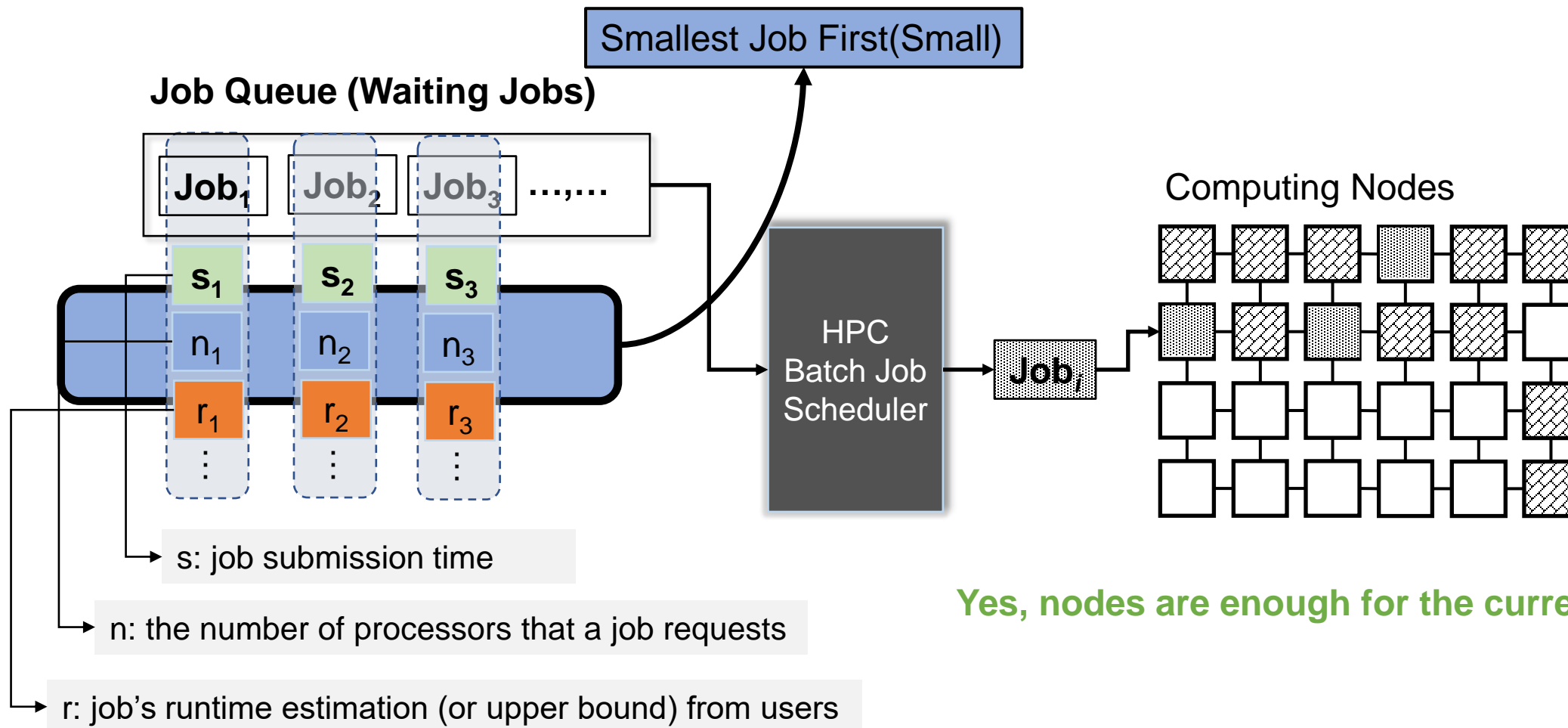
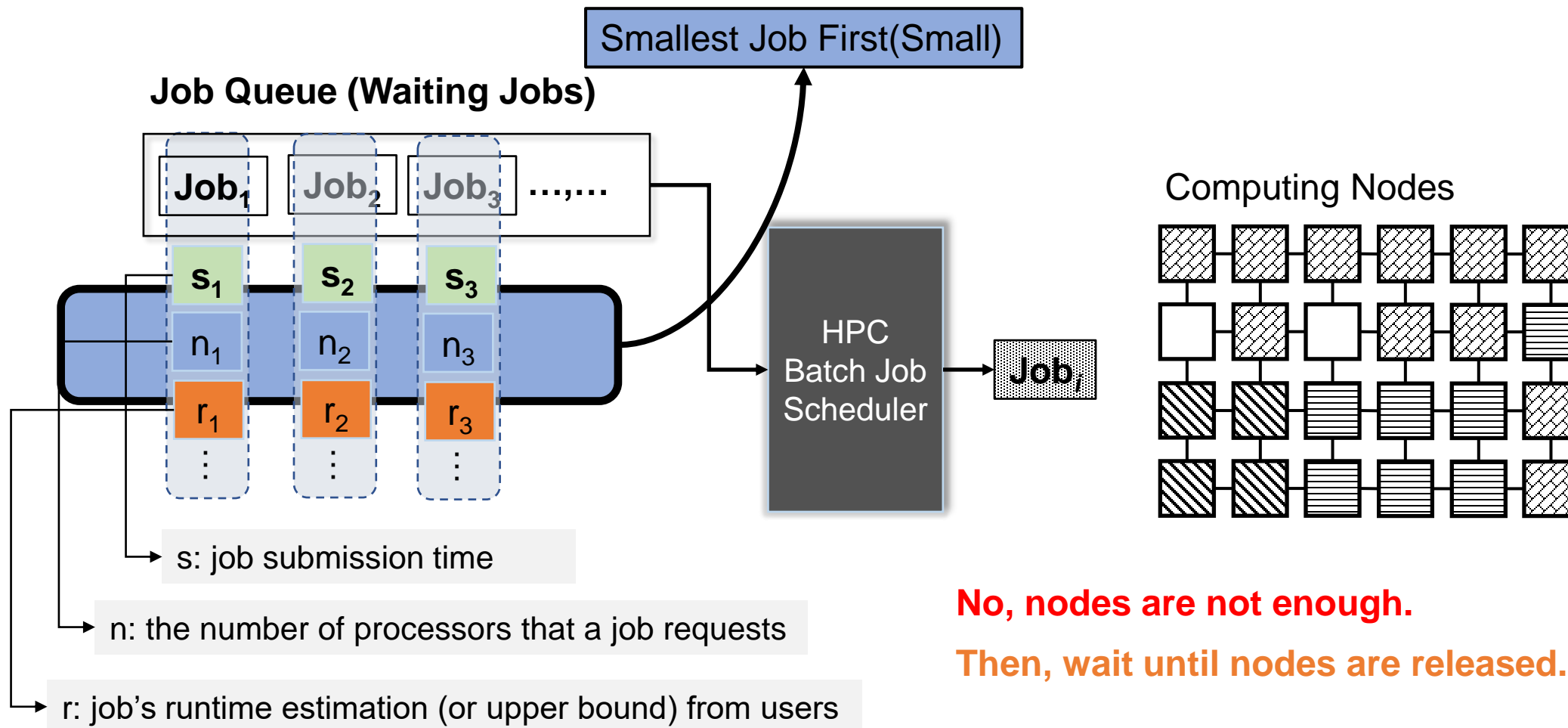[2]Oak Ridge National Laboratory

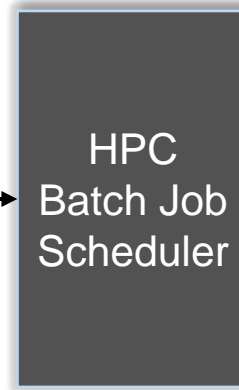# HPC Batch Job Scheduler

**Job Queue (Waiting Jobs)**



s: job submission time

n: the number of processors that a job requests

r: job's runtime estimation (or upper bound) from users

# HPC Batch Job Scheduler

Smallest Job First(Small)

**Job Queue (Waiting Jobs)**

| Job$_1$ | Job$_2$ | Job$_3$ | ...,... |

| $s_1$ | $s_2$ | $s_3$ |
| $n_1$ | $n_2$ | $n_3$ |
| $r_1$ | $r_2$ | $r_3$ |
| ... | ... | ... |

HPC Batch Job Scheduler

Job$_i$

Computing Nodes

s: job submission time

n: the number of processors that a job requests

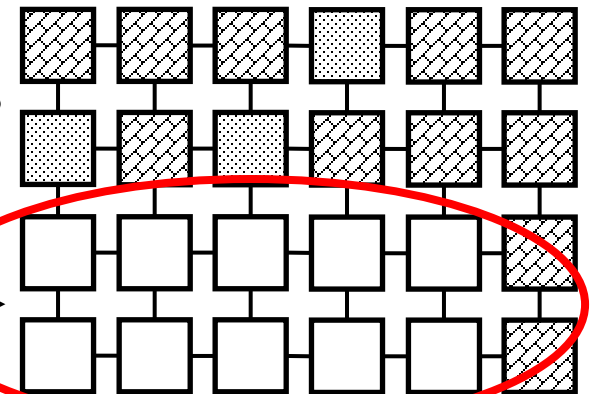r: job's runtime estimation (or upper bound) from users

**Enough Nodes?**

# HPC Batch Job Scheduler
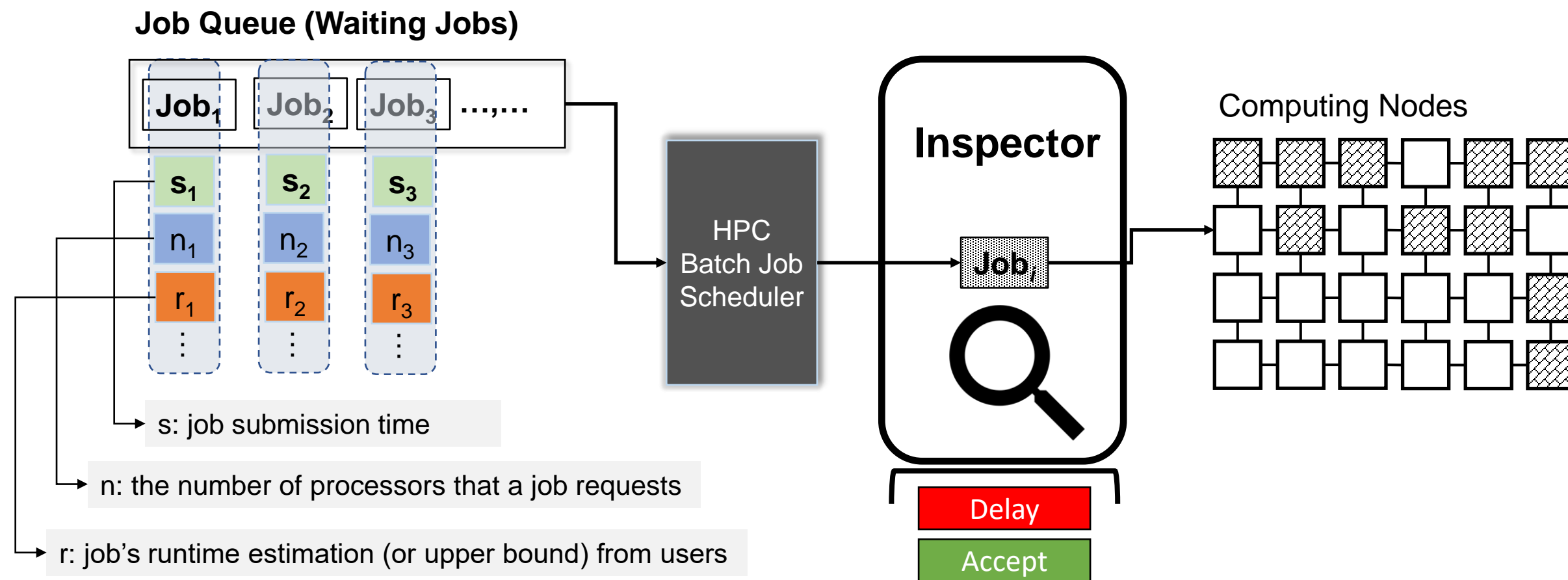


Smallest Job First(Small)

Job Queue (Waiting Jobs)

$Job_1$ $Job_2$ $Job_3$ …,…

$s_1$ $s_2$ $s_3$

$n_1$ $n_2$ $n_3$

$r_1$ $r_2$ $r_3$

HPC Batch Job Scheduler

$Job_i$

Computing Nodes

s: job submission time

n: the number of processors that a job requests

r: job's runtime estimation (or upper bound) from users

**Yes, nodes are enough for the current job.**

DIRLAB

# HPC Batch Job Scheduler

**DIRLAB**

Smallest Job First(Small)

**Job Queue (Waiting Jobs)**

| Job$_1$ | Job$_2$ | Job$_3$ | …,… |
|---------|---------|---------|-----|
| $s_1$ | $s_2$ | $s_3$ | |
| $n_1$ | $n_2$ | $n_3$ | |
| $r_1$ | $r_2$ | $r_3$ | |
| ⋮ | ⋮ | ⋮ | |

HPC Batch Job Scheduler → Job$_i$

Computing Nodes

$s$: job submission time

$n$: the number of processors that a job requests

$r$: job's runtime estimation (or upper bound) from users

**No, nodes are not enough.**

**Then, wait until nodes are released.**

# HPC Batch Job Scheduler



Smallest Job First(Small)

Job Queue (Waiting Jobs)

$Job_1$  $Job_2$  $Job_3$  …,…

$s_1$  $s_2$  $s_3$

$n_1$  $n_2$  $n_3$

$r_1$  $r_2$  $r_3$

HPC Batch Job Scheduler

$Job_j$

Computing Nodes

s: job submission time

n: the number of processors that a job requests

r: job's runtime estimation (or upper bound) from users

**Now, nodes are enough for the current job.**

# Motivation Example

Smallest Job First(Small)

**Job Queue (Waiting Jobs)**

| Job$_1$ | Job$_2$ | Job$_3$ | Job$_4$ |
|---------|---------|---------|---------|
| s$_1$ | s$_2$ | s$_3$ | s$_4$ |
| 30 | 20 | 12 | 10 |
| r$_1$ | r$_2$ | r$_3$ | r$_4$ |
| ⋮ | ⋮ | ⋮ | ⋮ |

HPC Batch Job Scheduler

Request 12 nodes

Job$_3$ ✗

Job$_4$

Computing Nodes

10 nodes are available

s: job submission time

n: the number of processors that a job requests

r: job's runtime estimation (or upper bound) from users

Hold Until Enough Nodes for Job$_3$

Delay Job$_3$ Re-decide Job$_4$

# Motivation

**Job Queue (Waiting Jobs)**

| Job$_1$ | Job$_2$ | Job$_3$ | ...,... |

$s_1$ $\quad$ $s_2$ $\quad$ $s_3$

$n_1$ $\quad$ $n_2$ $\quad$ $n_3$

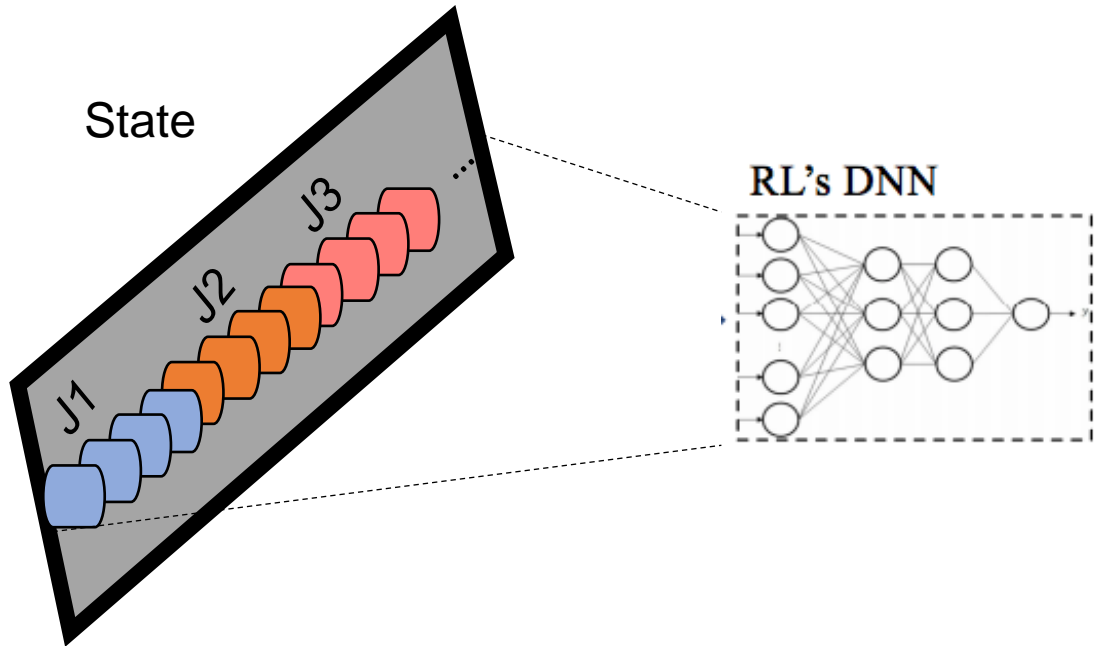$r_1$ $\quad$ $r_2$ $\quad$ $r_3$

$\vdots$ $\quad$ $\vdots$ $\quad$ $\vdots$

s: job submission time

n: the number of processors that a job requests

r: job's runtime estimation (or upper bound) from users

HPC Batch Job Scheduler
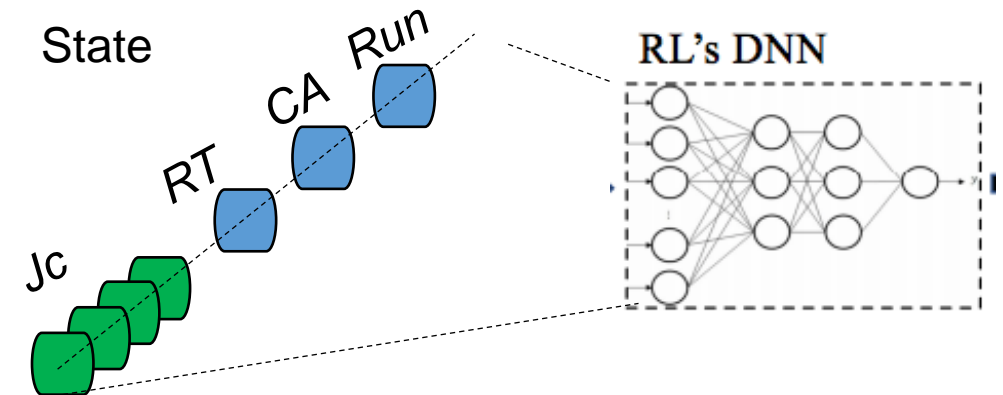
**Inspector**

Job$_i$

Delay

Accept

Computing Nodes

# Challenges

Current Status

Future Job Arrival

| | | |
|---|---|---|
| Understanding of historical data | Impact of the rejection | Whether the job is runnable |
| Attributes of the selected job | Number of rejected times | … |

# Reinforcement Learning



David Silver, et. al. Mastering the game of Go with deep neural networks and tree search, Nature vol. 529 (2016)

Volodymyr Mnih, et. al. Playing Atari with Deep Reinforcement Learning arXiv:1312.5602 (cs)

From https://www.selfdrivingcars360.com/how-autonomous-vehicles-fit-into-our-ai-enabled-future/

# Our Contribution

- The first scheduling inspector for HPC systems.

- New optimizations of the state and reward to enable efficient RL training.

- Extensively evaluations on efficiency, stability and interpretability of SchedInspector.

# SchedInspector



**Job Queue (Waiting Jobs)**

Job$_1$  Job$_2$  Job$_3$ …,…

s$_1$  s$_2$  s$_3$

n$_1$  n$_2$  n$_3$

r$_1$  r$_2$  r$_3$

s: job submission time

n: the number of processors that a job requests

r: job's runtime estimation (or upper bound) from users

HPC Batch Job Scheduler

**RL-based Inspector**

Job$_i$

Reward

Computing Nodes

Agent

state $S_t$   reward $R_t$   action $A_t$

$R_{t+1}$
$S_{t+1}$

Environment

DIRLAB

# Design of State

**Naïve Features**

State

J1 J2 J3 ...

RL's DNN

**Compacted Features**

*Jc: Scheduled Job*
*RT: Rejected Times*
*CA: Cluster Avail.*
*Run: Runnable*

State

Jc RT CA Run

RL's DNN

# Design of State

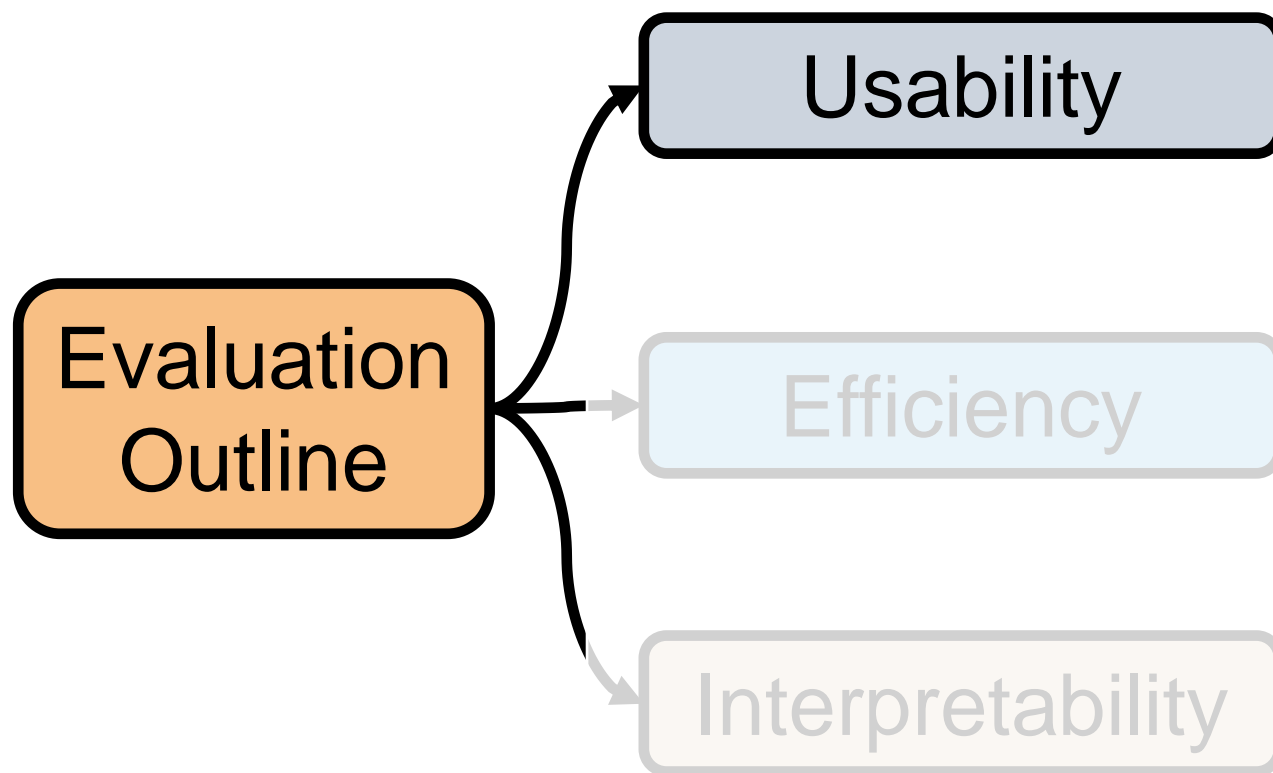**SchedInspector**

*QD: Queue Delay*
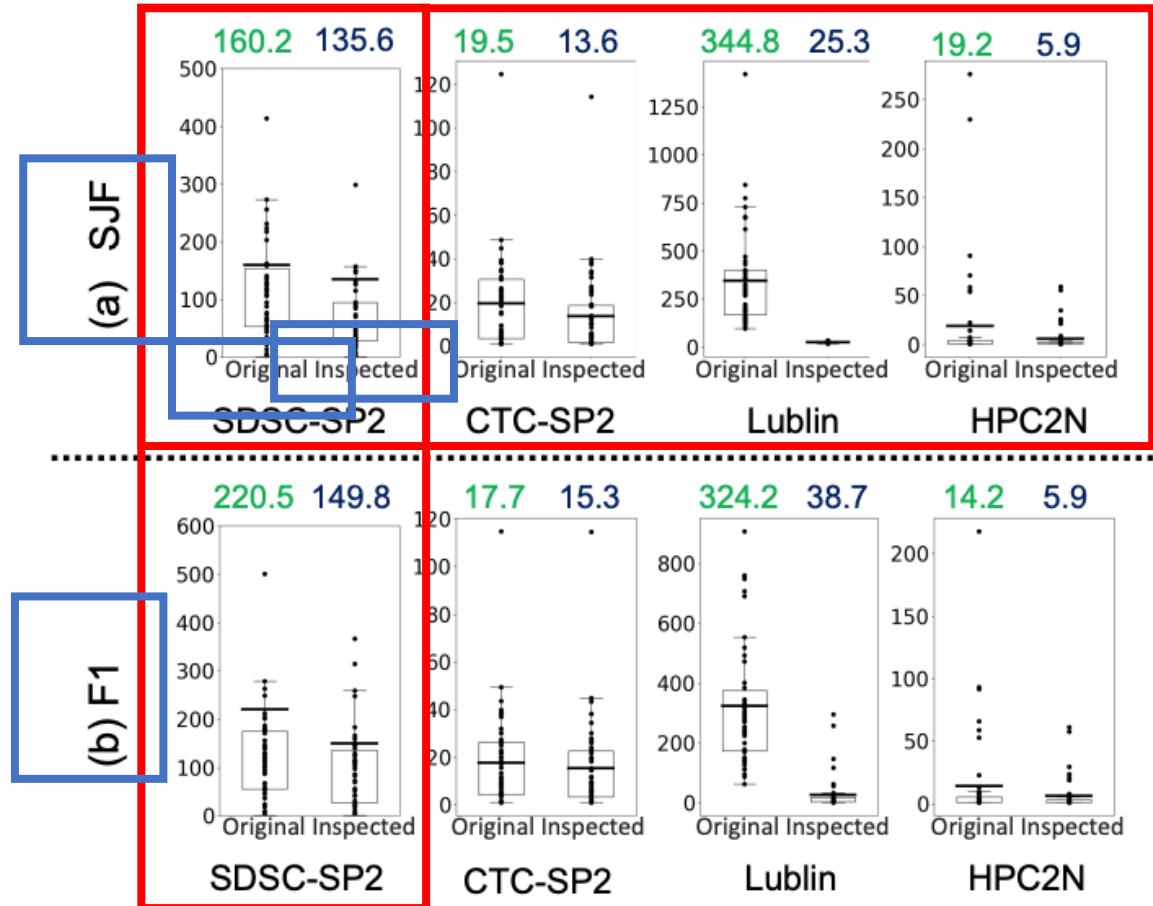*BF: Backfilling*

# Design of Reward

Naïve: $Metric_{inspect} - Metric_{orig}$

Win/Loss: $Integer(Metric_{inspect} > Metric_{orig})$

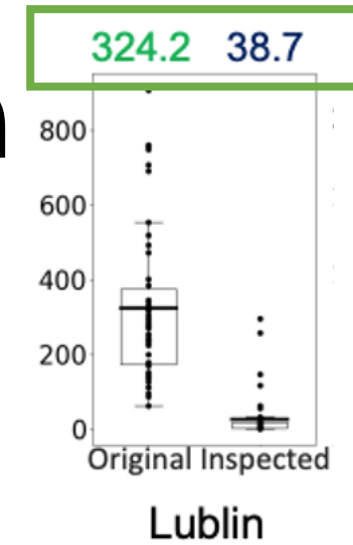✓ Percentage: $(Metric_{inspect} - Metric_{orig})/ Metric_{orig}$

# Testing for **Different Job Traces and Policies**



SchedInspector has significant improvement for the two scheduling policies on all job traces.
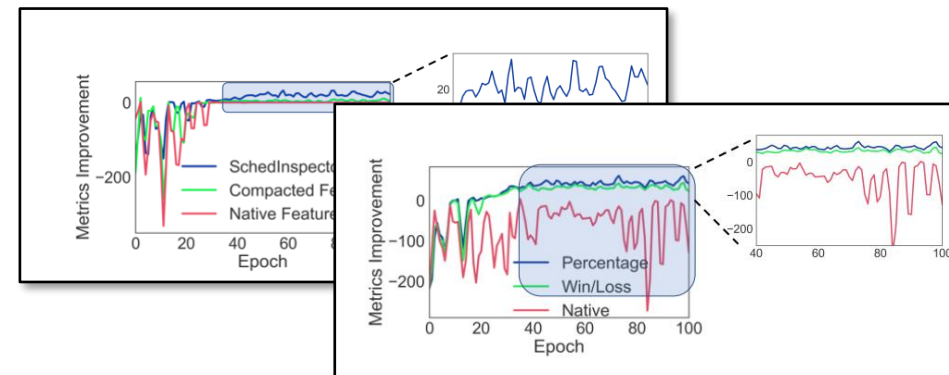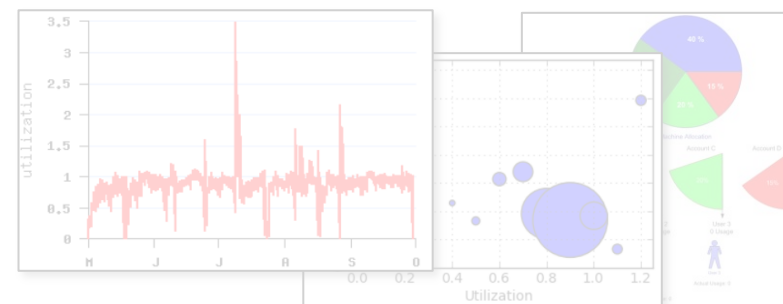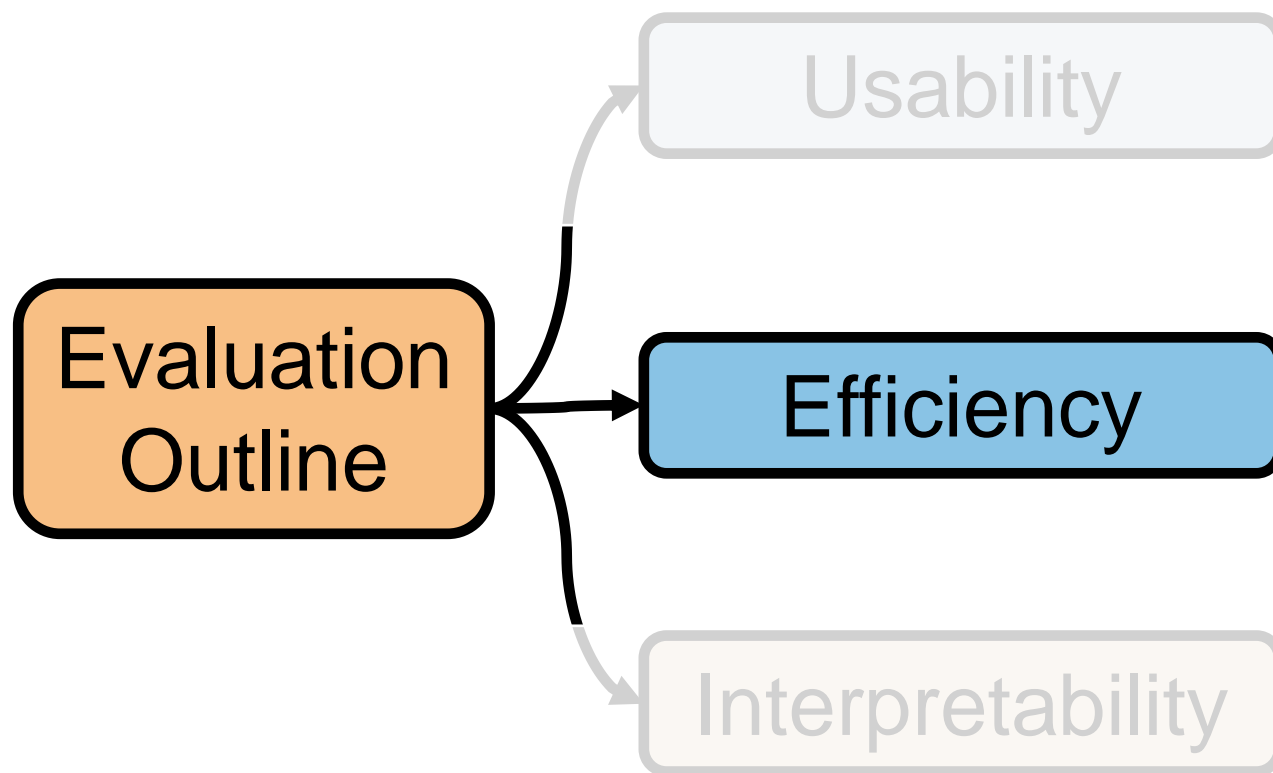
# Impact on **System Utilization**

324.2   38.7

**88% Improvement**

| | SJF | | | F1 | | |
|---|---|---|---|---|---|---|
| | BASE | INSP | Δ | BASE | INSP | Δ |
| *Scheduling without Backfilling* | | | | | | |
| SDSC-SP2 | 59.64% | 59.37% | **-0.27%** | 60.18% | 60.59% | **+0.41%** |
| CTC-SP2 | 51.35% | 49.92% | **-1.43%** | 54.40% | 54.23% | **-0.17%** |
| Lublin | 61.49% | 61.06% | **-0.43%** | 67.37% | 63.04% | **-4.33%** |
| HPC2N | 23.72% | 23.47% | **-0.25%** | 24.00% | 23.79% | **-0.21%** |
| *Scheduling with Backfilling* | | | | | | |
| SDSC-SP2 | 78.45% | 78.37% | **-0.08%** | 76.71% | 76.93% | **+0.22%** |
| CTC-SP2 | 74.98% | 74.89% | **-0.09%** | 75.47% | 76.05% | **+0.58%** |
| Lublin | 79.38% | 77.71% | **-1.67%** | 80.38% | 78.08% | **-2.30%** |
| HPC2N | 56.81% | 57.10% | **+0.29%** | 57.11% | 56.57% | **-0.54%** |

Lublin

SchedInspector has barely noticeable reduction (1% difference) on system utilization

## Evaluation Outline

Usability

**Efficiency**

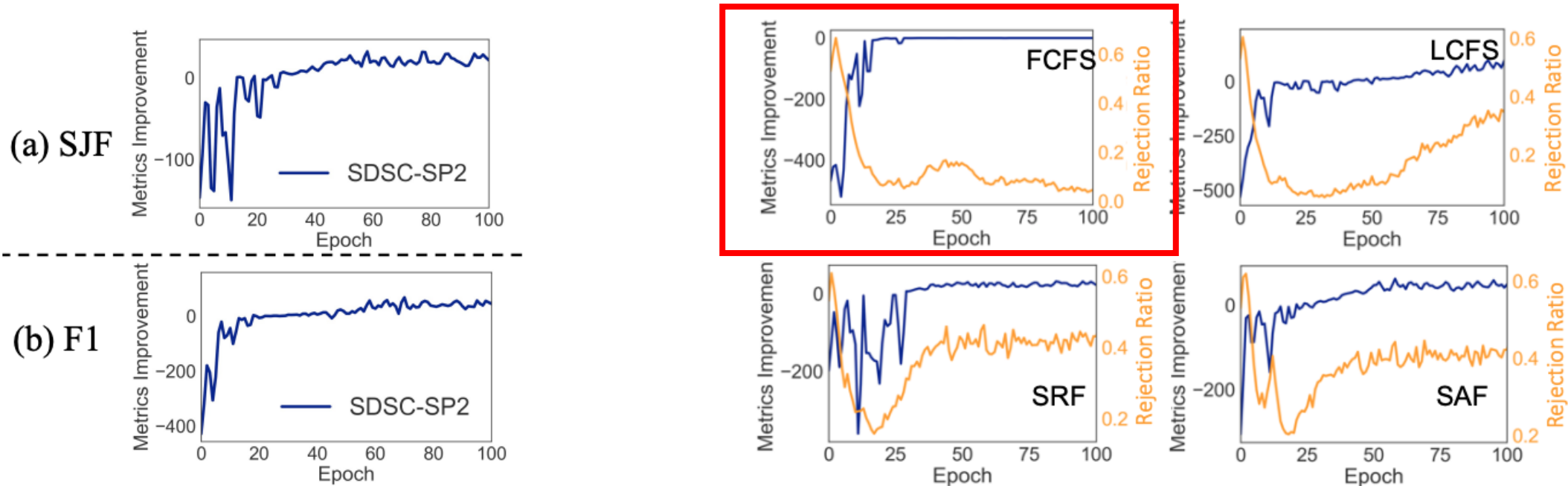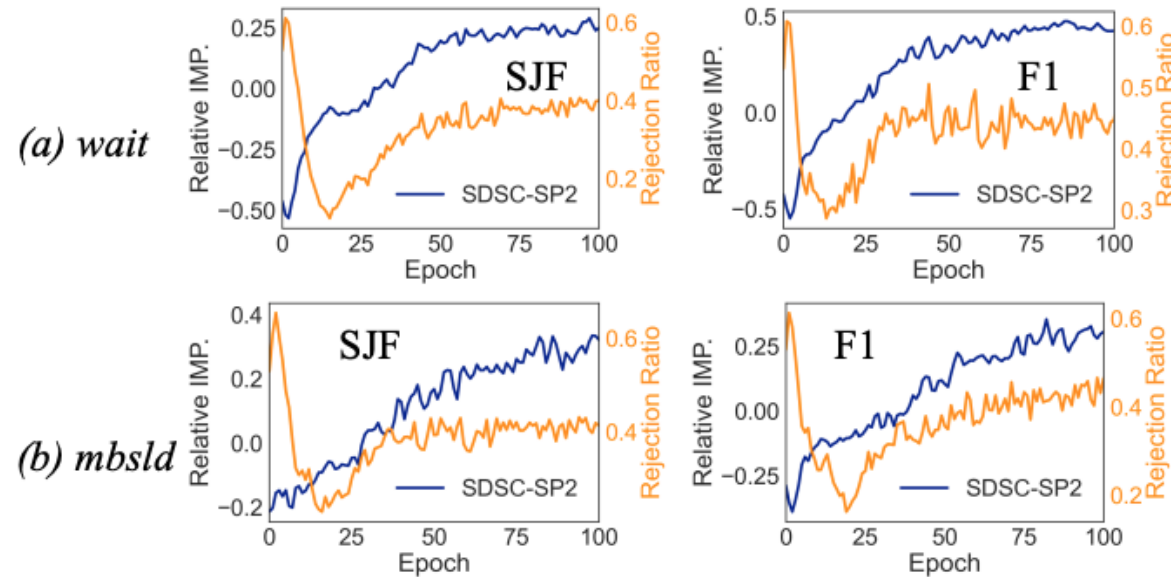Interpretability

# Training on **Different Job Traces**



SchedInspector converges in all of the workloads within 100 training epochs and different job traces have different converge pattern.

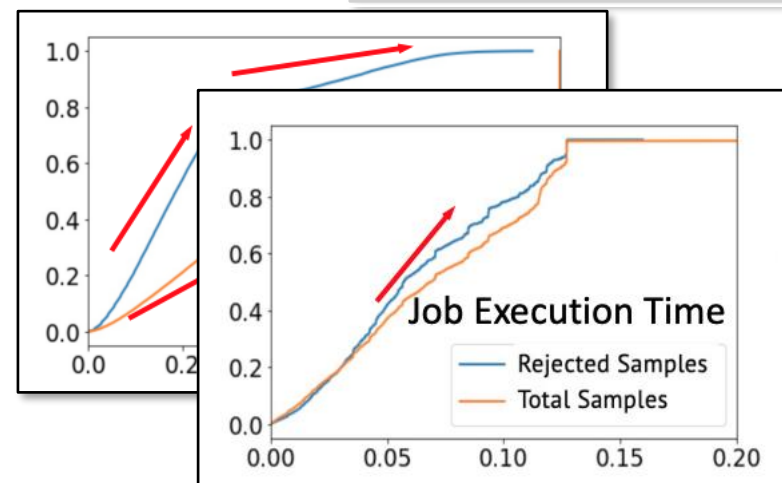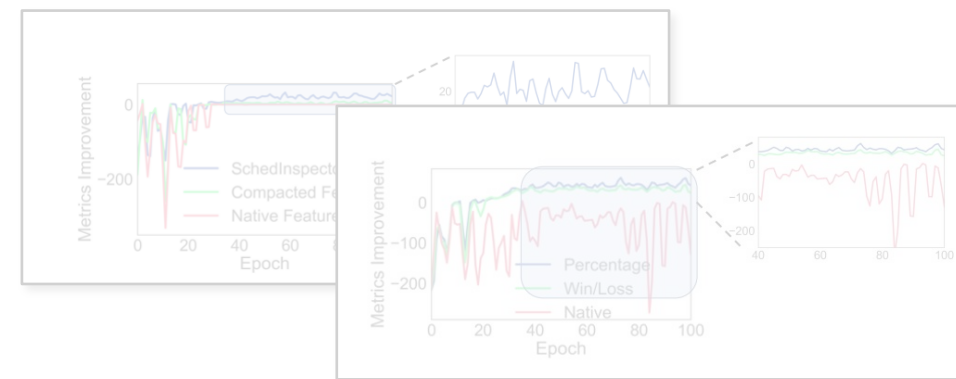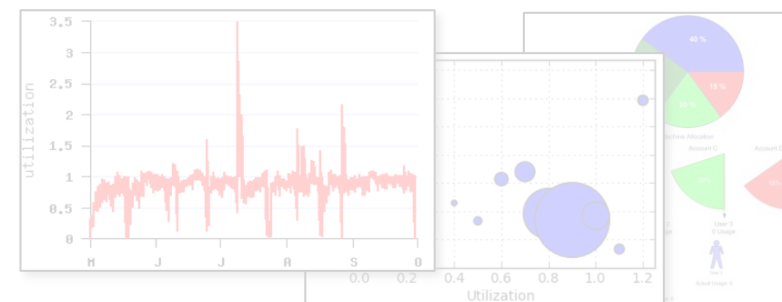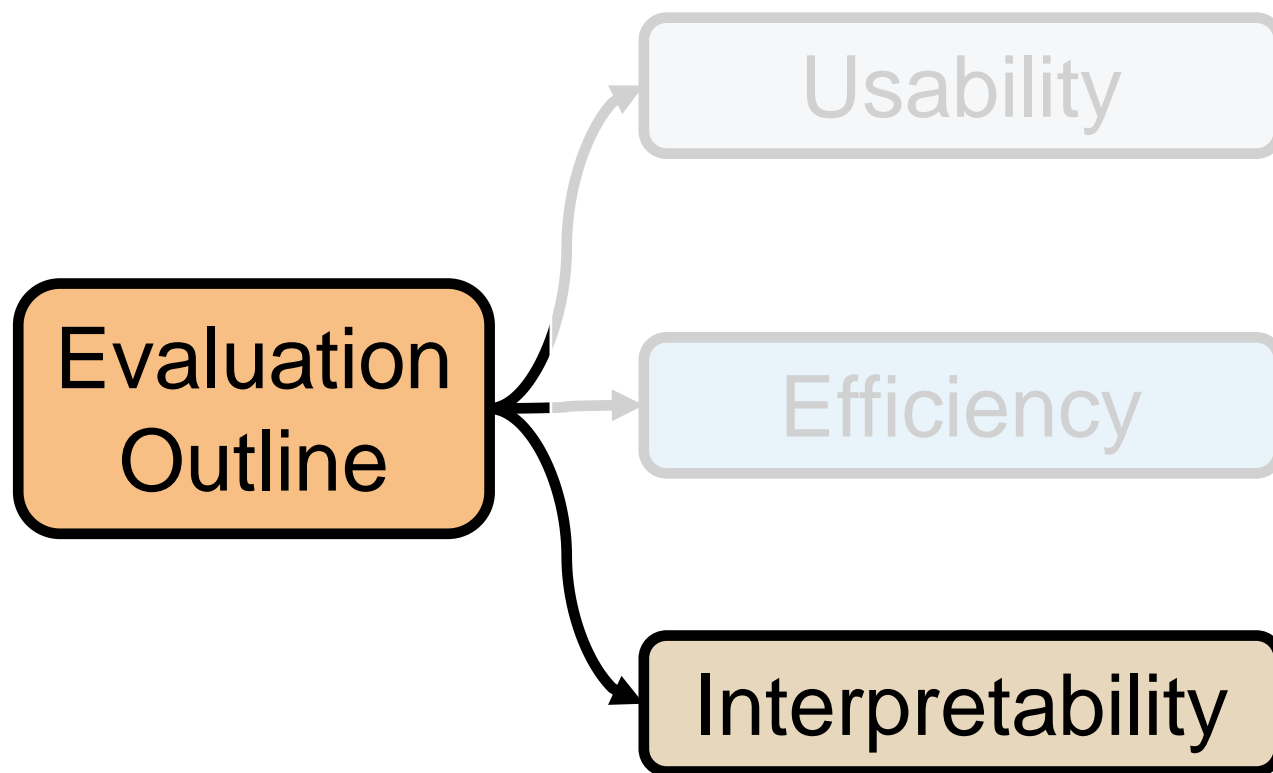# Training on **Different Scheduling Policies**



SchedInspector converges in all scheduling policies.
For some scheduling policies, the converged value is near
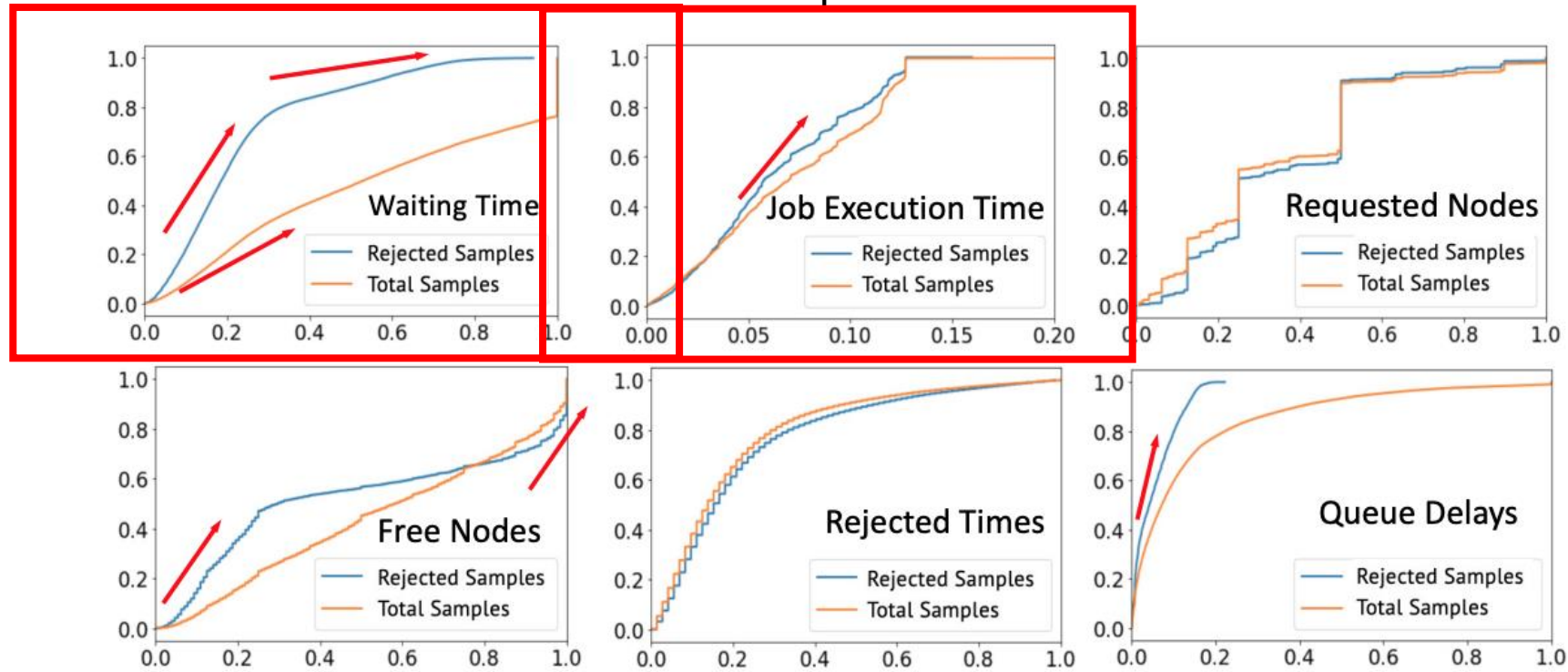0 and the rejection ratio is low.

# Training for **Different Metrics**



SchedInspector converges towards two new metrics but with different patterns.

# Evaluation Outline

- Usability
- Efficiency
- **Interpretability**

# What SchedInspector Learns

CDF of input features.



SchedInspector has obvious patterns for different features which indicates the effectiveness of feature selection

# Hare: Exploiting Inter-job and Intra-job Parallelism of Distributed Machine Learning on Heterogeneous GPUs

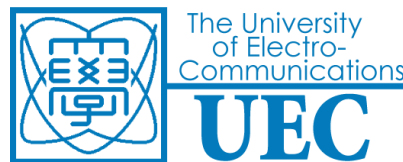Fahao Chen[1], Peng Li[1], Celimuge Wu[2], Song Guo[3]

[1]The University of Aizu

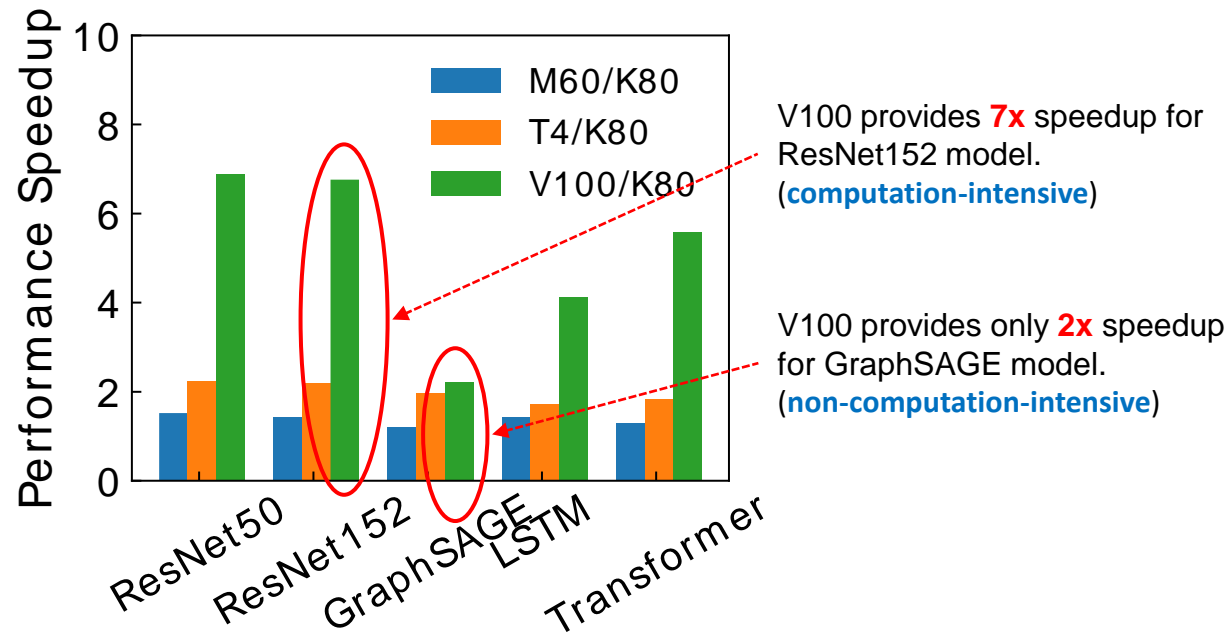[2]University of Electro-Communications

[3]The Hong Kong Polytechnic University & The Hong Kong Polytechnic University Shenzhen Research Institute
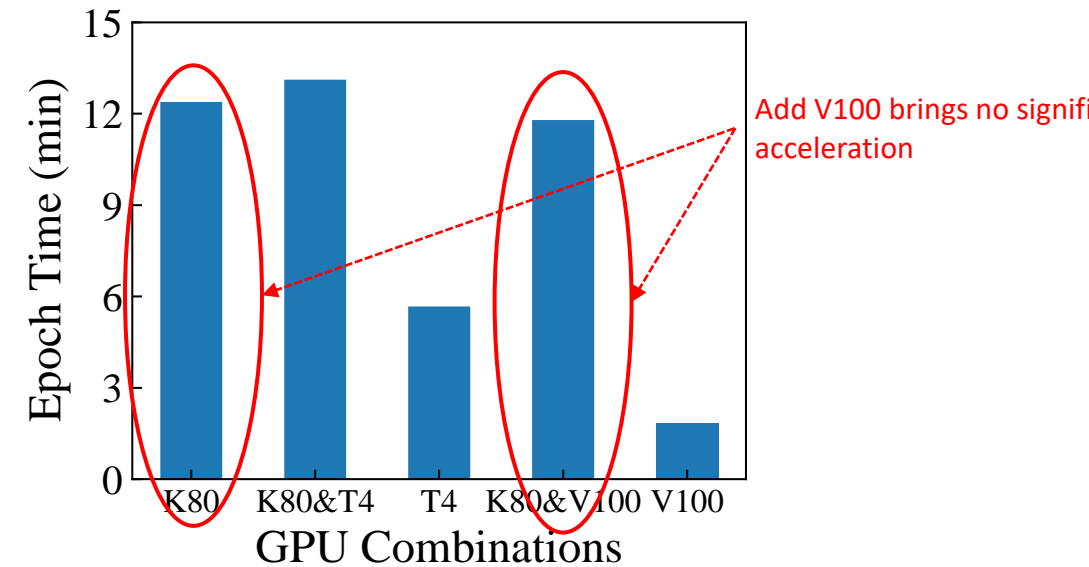
**ACM HPDC 2022**
**Minneapolis, Minnesota, United States**

# Why Do We Need to Consider GPU Heterogeneity?



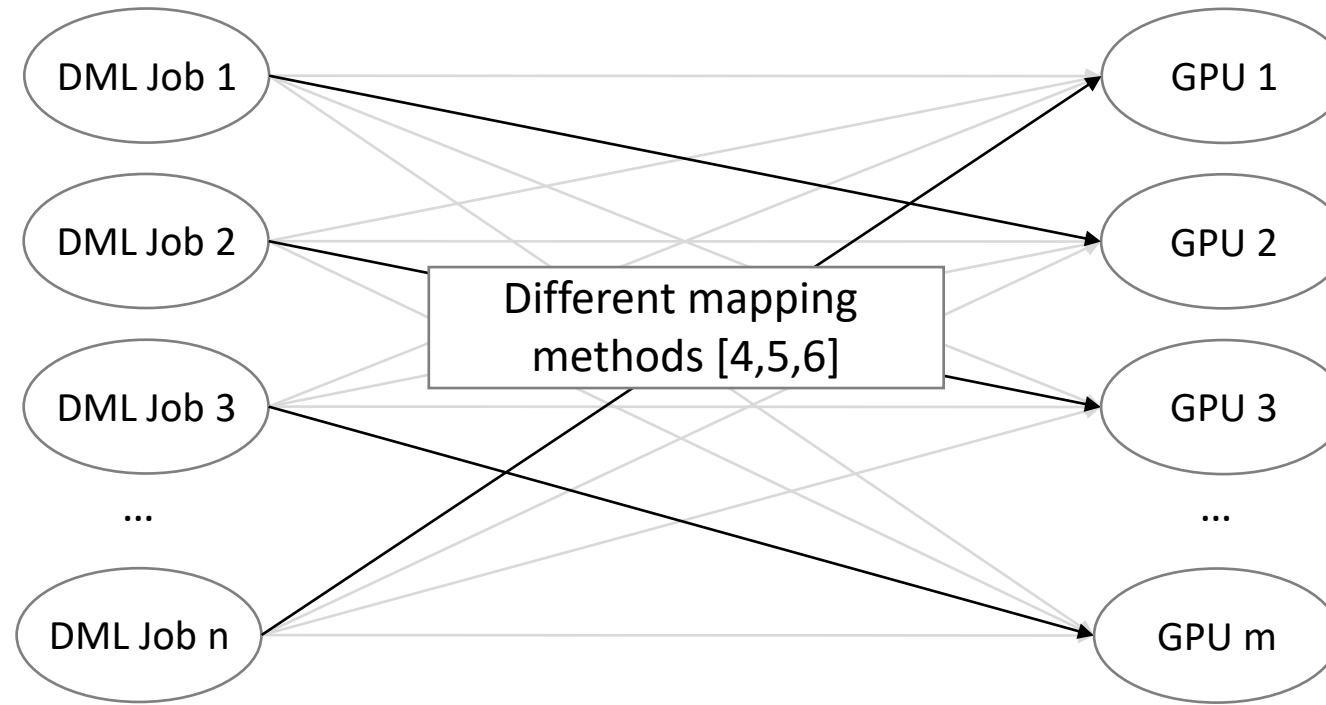(a) Training speedup of different jobs on different GPUs.

V100 provides **7x** speedup for ResNet152 model. (**computation-intensive**)

V100 provides only **2x** speedup for GraphSAGE model. (**non-computation-intensive**)



(b) Round time of ResNet152 under different GPU combinations.

Add V100 brings no signifi acceleration

- Different GPUs provide different performance speedups for different jobs. (**inter-job parallelism**)

- Different GPU combinations provide different performance speedups. (**intra-job parallelism**)

# Existing Works on Heterogeneous GPUs



**They treat DML jobs as unsplittable units and ignore the intra-job parallelism.**
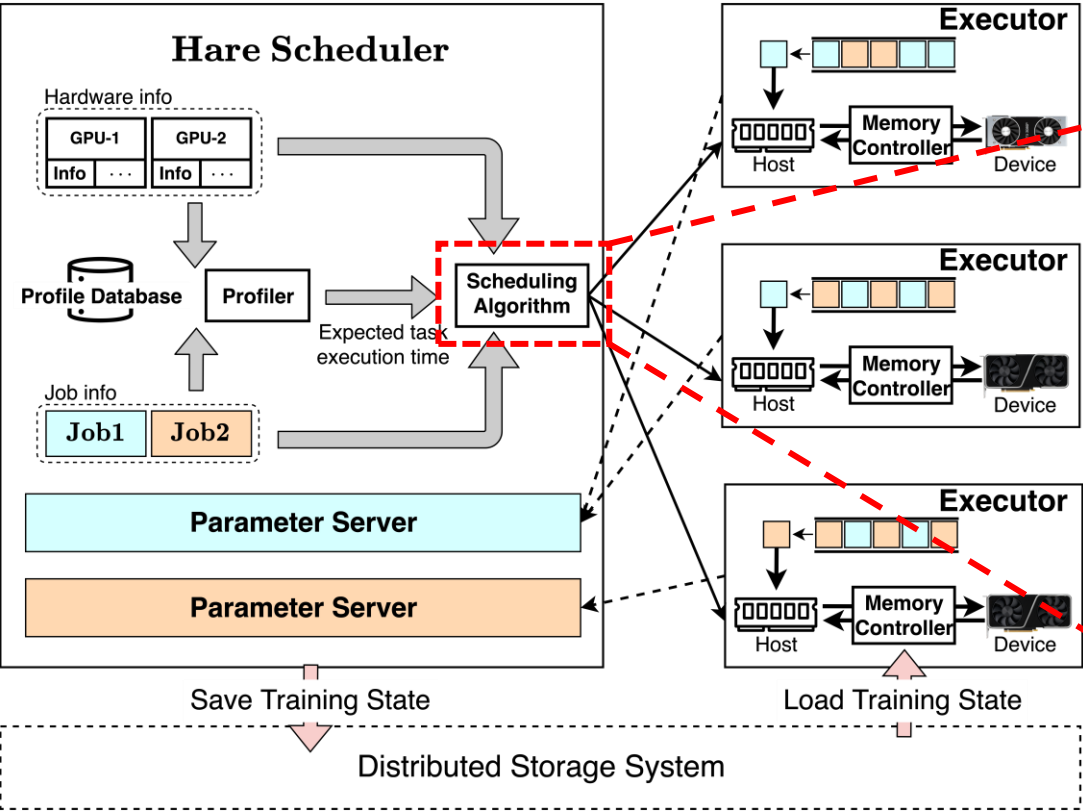
[6] Chaudhary, Shubham, et al. "Balancing efficiency and fairness in heterogeneous GPU clusters for deep learning." Proceedings of the Thirteenth EuroSys Conference. 2020.
[7] Le, Tan N., et al. "Allox: compute allocation in hybrid clusters." Proceedings of the Thirteenth EuroSys Conference. 2020.
[8] Narayanan, Deepak, et al. "Heterogeneity-Aware Cluster Scheduling Policies for Deep Learning Workloads." 14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20).

# Hare: An Efficient DML Job Training System

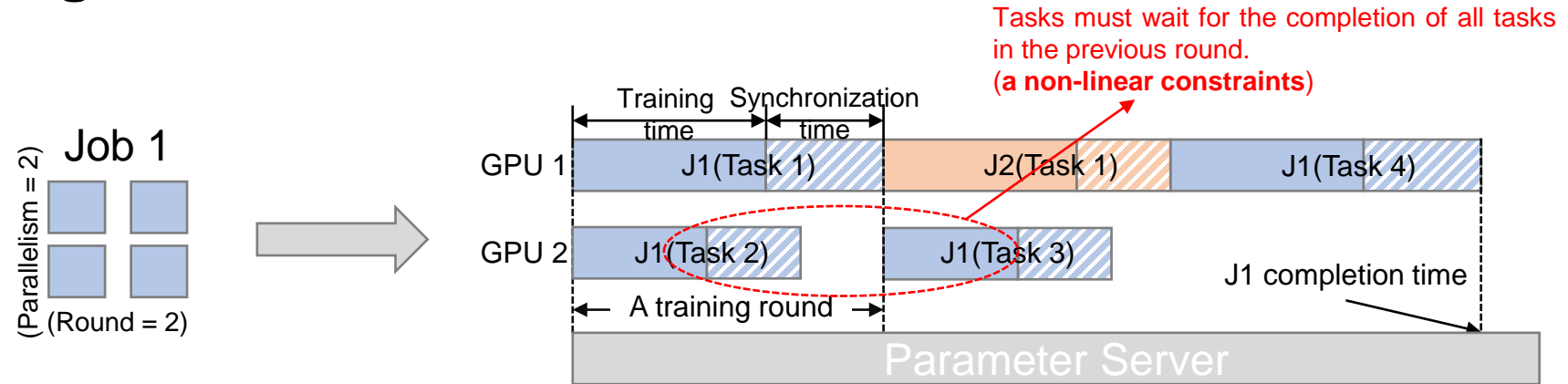**System Goals: *High training efficiency*, *High GPU utilization*, and *Starvation-free***



Overview of Hare

| Scheduling Algorithm | Intra-job Parallelism | Inter-job Parallelism | Theoretical bound |
|---|---|---|---|
| Gandiva Tiresias Optimus Themis Antman | ✓ | ✗ | ✗ |
| Allox | ✗ | ✓ | ✗ |
| Hare | ✓ | ✓ | ✓ |

Compared with other scheduling algorithms

# Scheduling Algorithm: Problem Statement

1. Model:



2. Task scheduling problem
- **Objective:** minimize the total weighted job completion time
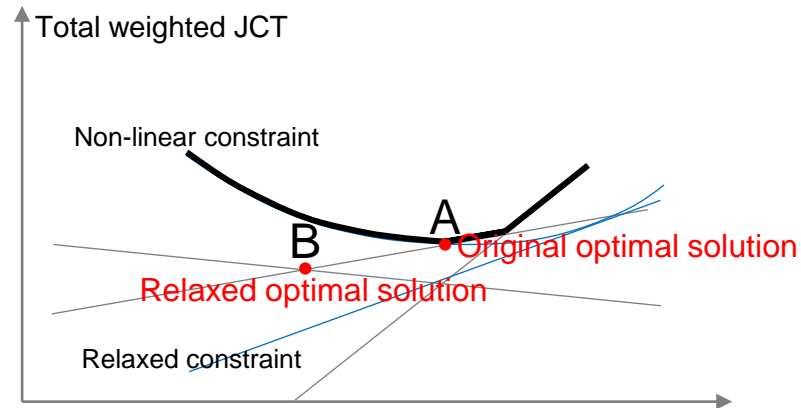- **Solution:** the start running time of each task

💥**Theorem:** The above problem is NP-hard, which cannot be solved within a polynomial time.

💥**Proof:** Reduce the well-known SS13 problem [11].

[11] Garey, Michael R., and David S. Johnson. Computers and intractability. Vol. 174. San Francisco: freeman, 1979.

# Scheduling Algorithm: Algorithm Design

Step 1: Problem relaxation [12]

Total weighted JCT

Non-linear constraint

A

B

Original optimal solution

Relaxed optimal solution

Relaxed constraint

Step 2: Task scheduling
- Decide the scheduling ordering of tasks according to the solution of relaxed problem
- Greedily assign tasks to GPUs with the earliest available time

✹**Theorem:** Our scheduling algorithm is $\alpha(2 + \alpha)$-approximation.

✹**Proof:** Please refer to our paper.

$\alpha = \max_{task}\{\frac{T^{c,max}}{T^{c,min}}, \frac{T^{s,max}}{T^{s,min}}\}$
$T^c$: task training time
$T^s$: task synchronization time

[12] Queyranne, Maurice. "Structure of a simple scheduling polyhedron." Mathematical Programming 58.1 (1993): 263-285.

# Evaluation: Experimental settings

**Testbed:** <span style="color:red">15</span> heterogeneous GPUs (V100 + T4 + K80 + M60)

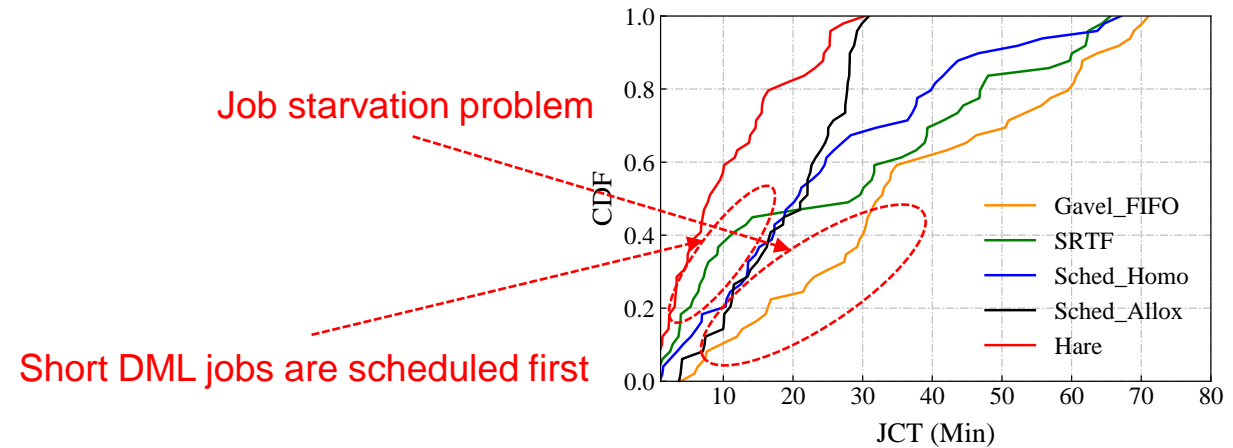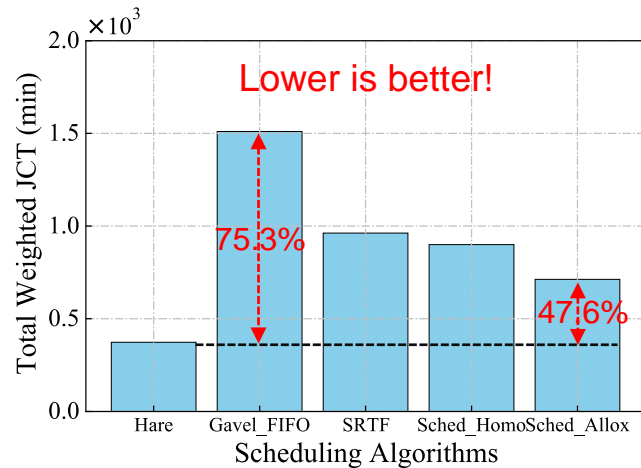**Simulator:** up to <span style="color:red">200</span> GPUs and <span style="color:red">300</span> DML jobs

**Workload:** <span style="color:red">8</span> popular models across domains of CV, NLP, Speech, and Recognition

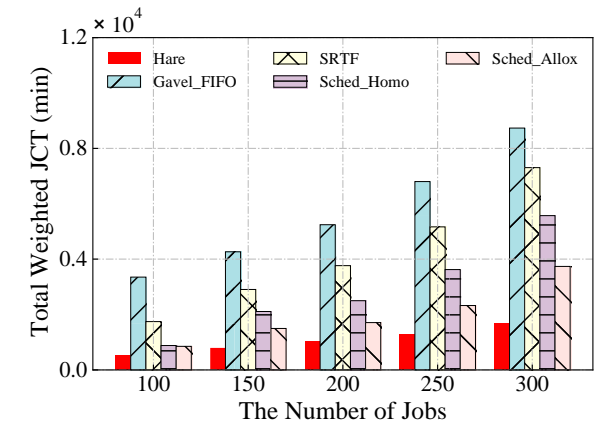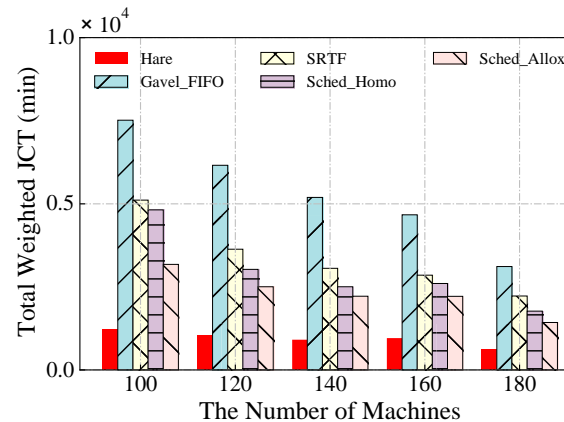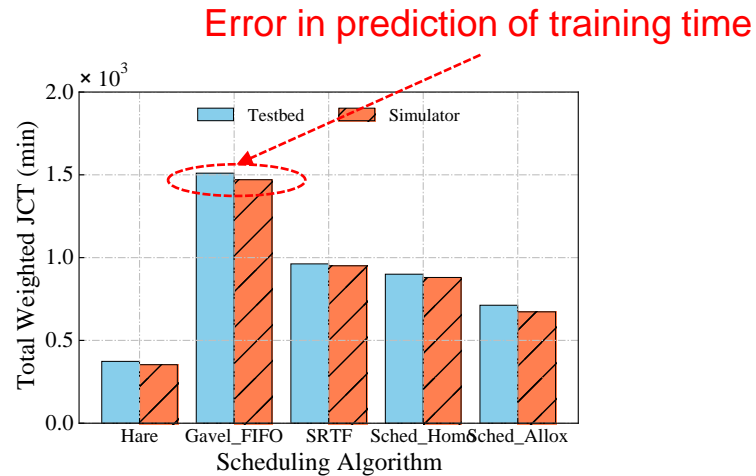**Baseline:** <span style="color:red">4</span> popular scheduling schemes.

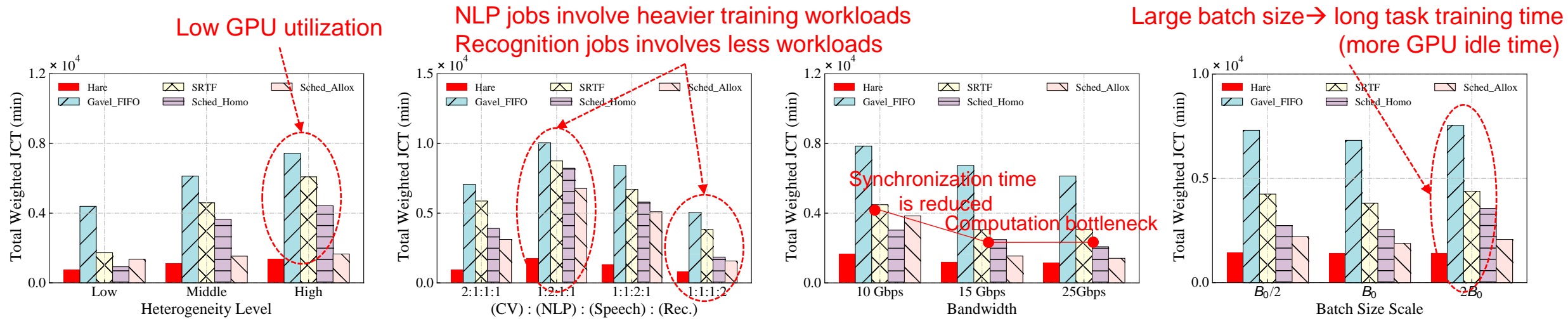(*Please refer to our paper for more details on experimental settings.)

# Evaluation: Testbed



- Our scheduling algorithm reduces the total weighted JCT by 47.6% to 75.3% by others.

- About 90.5% of jobs can complete within 25 minutes based on our scheduling algorithm.

# Evaluation: Simulation



- The maximum performance gap between the testbed and simulator is only 5%.

- Our scheduling algorithm outperforms others under a scaling of GPUs and jobs.

# Evaluation: Simulation



- Our scheduling algorithm outperforms others under various settings.

# Conclusion

- **Hare:** An DML jobs training system with the objectives of ***high training efficiency***, ***high GPU utilization***, and ***starvation-free***

- **Task Scheduling Algorithm:**
  - Minimize total weighted JCT
  - Reduce total weighted JCT by **47.6% to 75.3%** over other schemes

- **Constraint:**
  - Jobs with changed setting may incur high overhead when scheduling
  - Hare is short in handling dynamic jobs

# *Efficient Design Space Exploration for Sparse Mixed Precision Neural Architectures*

**Krishna Teja Chitty-Venkata**[1], Arun Somani[1]
Murali Emani[2], Venkatram Vishwanath[2]

[1]Iowa State University, Ames, IA, USA
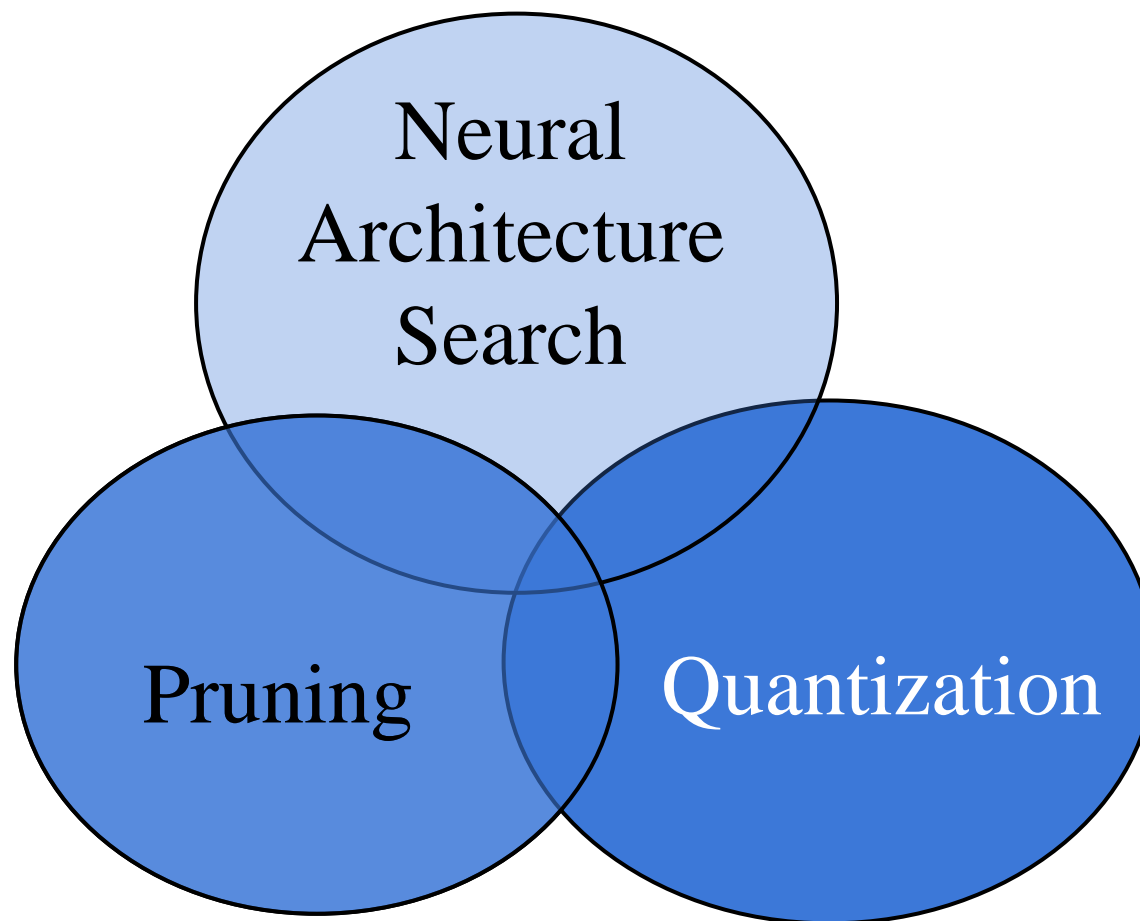[2]Argonne National Laboratory, Lemont, IL, USA
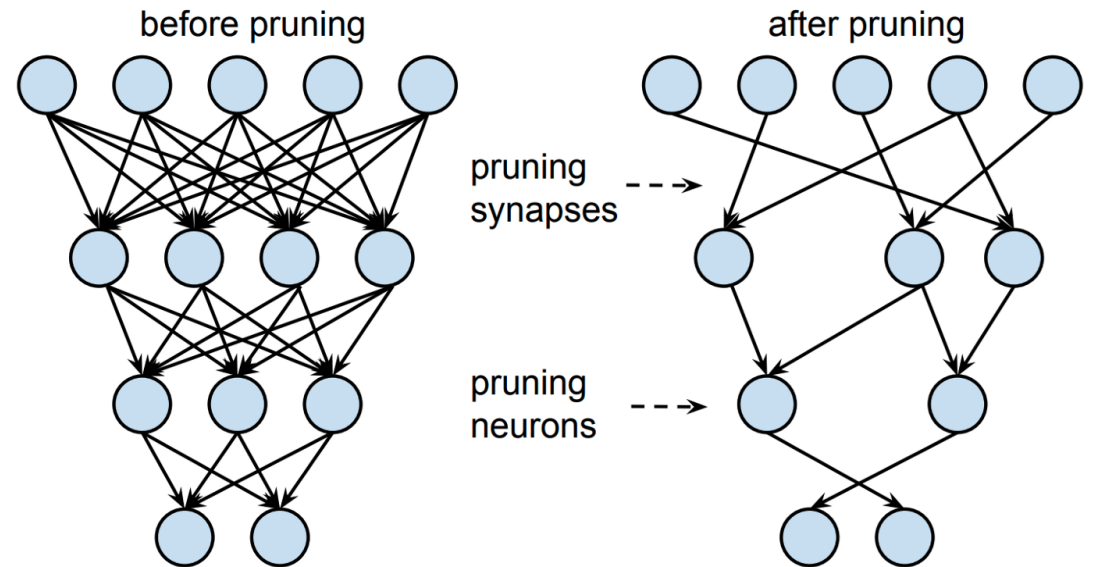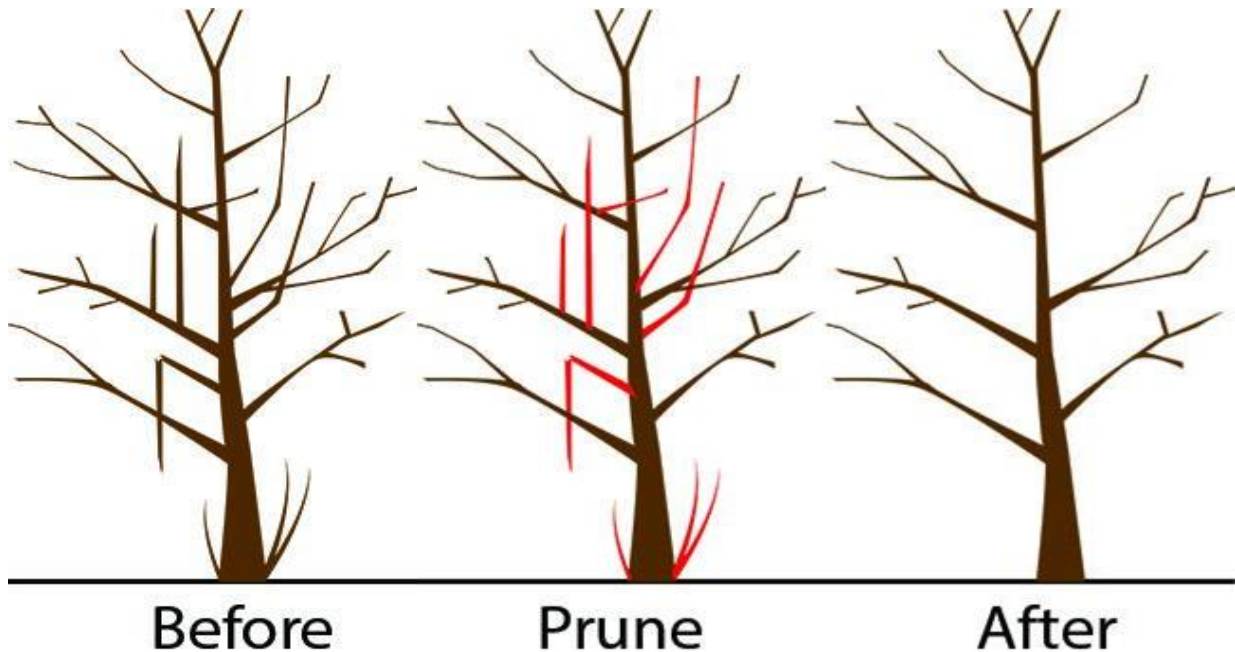
ACM HPDC 2022

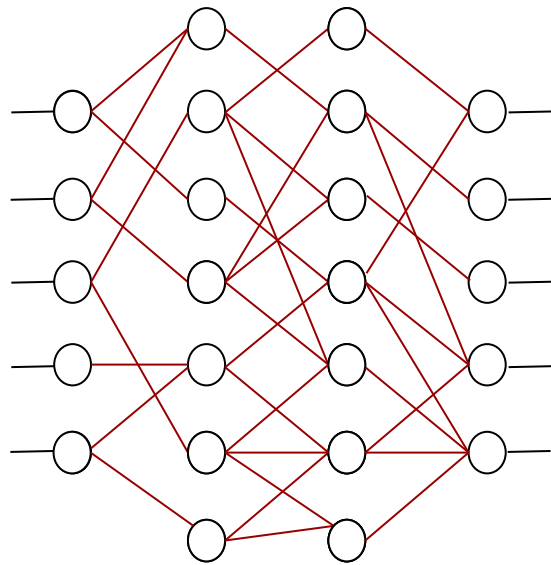The 31st International Symposium on High-Performance Parallel and Distributed Computing

# *Pruning*

- Pruning in general refers to cutting down branches/leaves

- Neural Network Pruning refers to removing weights/connection without compromising the accuracy

- The parameters which do not contribute to the final accuracy are removed to save memory

- Example: Magnitude-based



Song Han et.al. Deep Compression    44

# Neural Network Pruning

## Weight/Irregular Pruning



### Weight Matrix

| 2 | 0 | 0 | 1 | 0 | 6 |
|---|---|---|---|---|---|
| 0 | 0 | 2 | 0 | 0 | 0 |
| 0 | 0 | 0 | 6 | 0 | 0 |
| 1 | 0 | 0 | 0 | 7 | 0 |
| 0 | 0 | 3 | 0 | 0 | 0 |
| 7 | 0 | 0 | 8 | 0 | 9 |

- Creates Sparse Matrices
- Requires sparse decoding
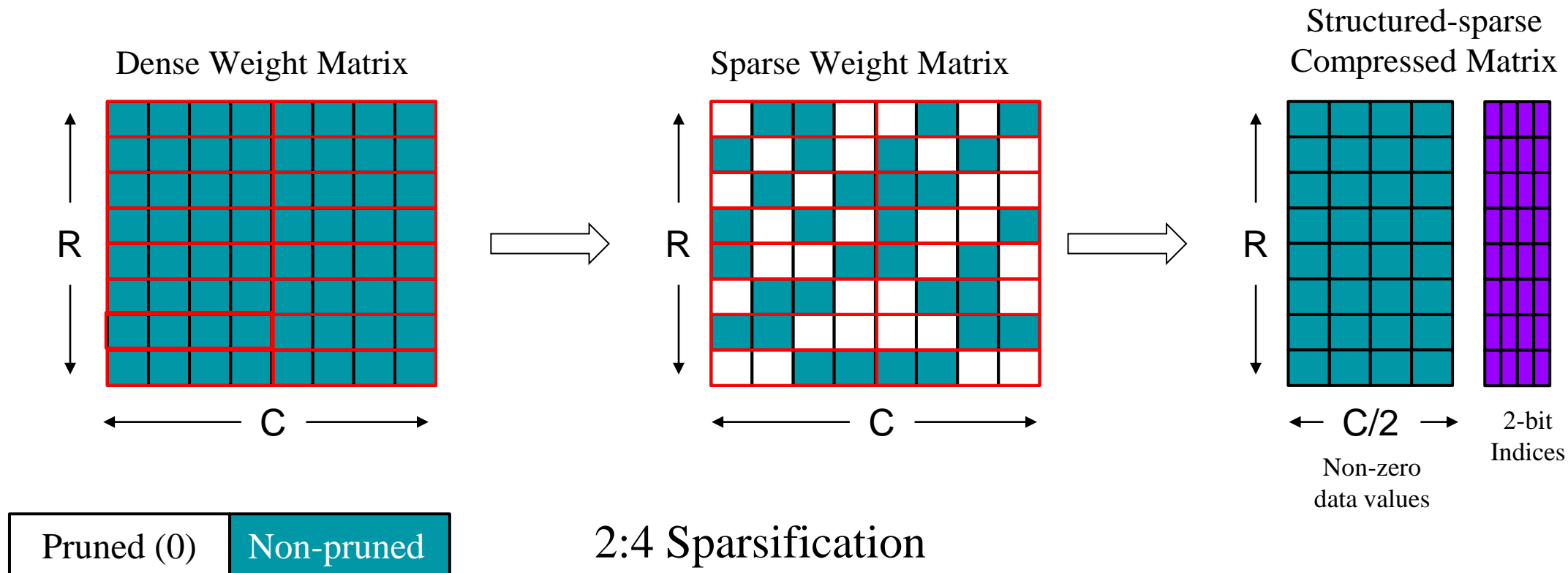- Irregular memory accesses if not stored contigously

## Node/Structured Pruning



### Weight Matrix

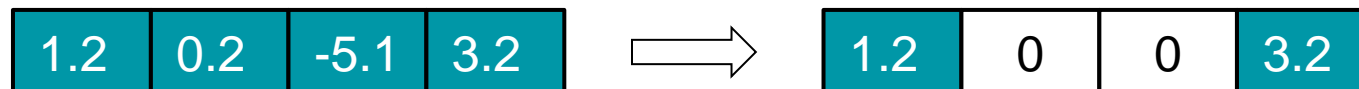| 0 | 1 | 9 | 0 | 2 | 0 |
|---|---|---|---|---|---|
| 0 | 8 | 2 | 0 | 9 | 0 |
| 0 | 1 | 5 | 0 | 1 | 0 |
| 0 | 3 | 8 | 0 | 7 | 0 |
| 0 | 5 | 3 | 0 | 7 | 0 |
| 0 | 2 | 7 | 0 | 1 | 0 |

- Regular Dense Matrices
- Regular DNN implementation
- Regular memory accesses as easily stored in contiguous memory

45

# 2:4 (two-to-four) Sparsity Pattern by Nvidia A100 GPU Tensor Cores:

- The 2:4 pattern mandates that for each group of 4 values, at least 2 must be zero
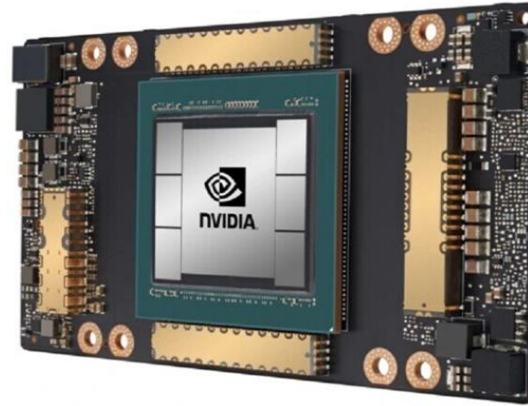- The pattern leads to 50% sparsity while maintaining accuracy



Dense Weight Matrix

Sparse Weight Matrix

Structured-sparse
Compressed Matrix

R

C

R

C

R

C/2

Non-zero
data values

2-bit
Indices

| Pruned (0) | Non-pruned |
|---|---|

2:4 Sparsification

Example Pruning Strategy:

| 1.2 | 0.2 | -5.1 | 3.2 |
|---|---|---|---|

⟹

| 1.2 | 0 | 0 | 3.2 |
|---|---|---|---|

46

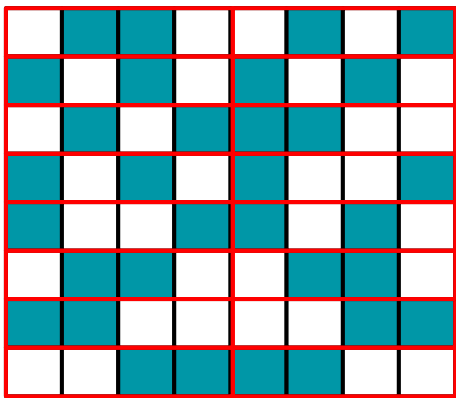Types of Multiplications Supported by Nvidia A100 GPU
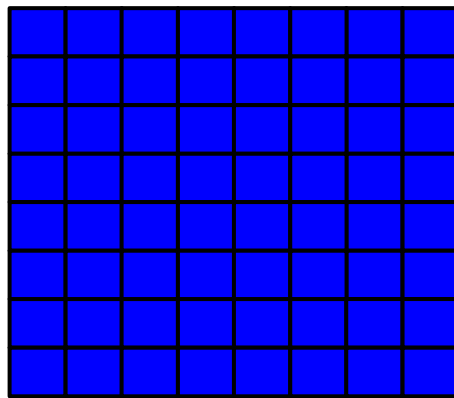
Theoretical Speedup

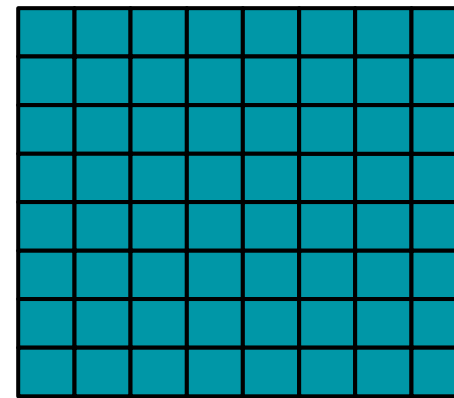Pruned (0) | Non-pruned | Non-pruned
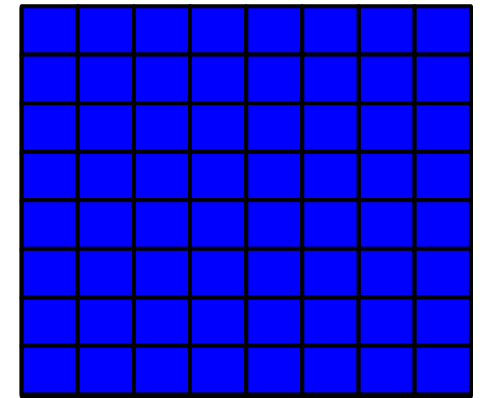
## 2:4 Sparse Multiplication

Weight Matrix * Activation Matrix

## Dense Multiplication

Weight Matrix * Activation Matrix
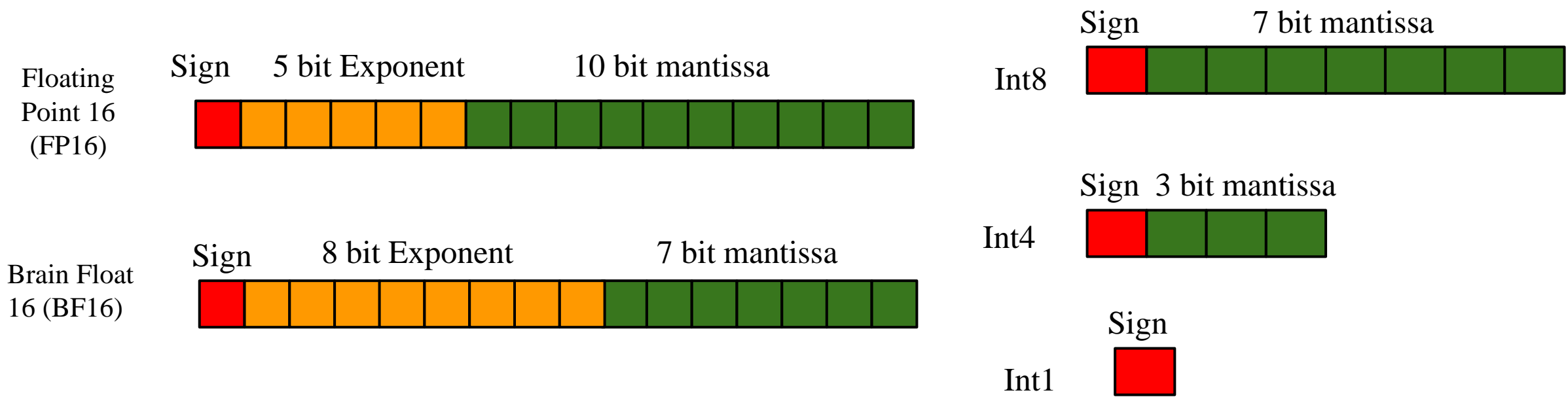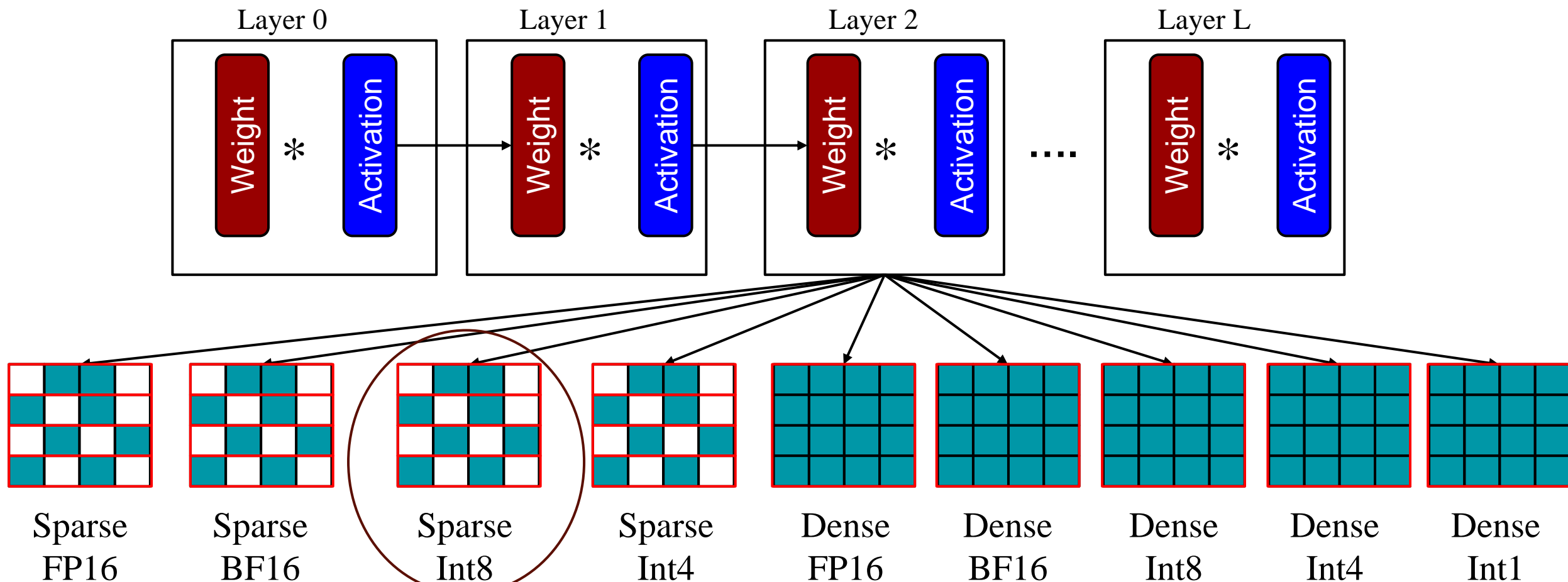
# Types of Precisions Supported by Nvidia A100 Tensor Cores



Summary of Sparse-Precisions Supported by Nvidia A100 GPU

| Multiplication | FP16 | BF16 | Int8 | Int4 | Int1 |
|---|---|---|---|---|---|
| 2:4 Sparse | Yes | Yes | Yes | Yes | No |
| Dense | Yes | Yes | Yes | Yes | Yes |

# Summary of Sparse-Precisions Supported by Nvidia A100 GPU Tensor Cores



Nvidia State-of-the-art
- 4x Latency
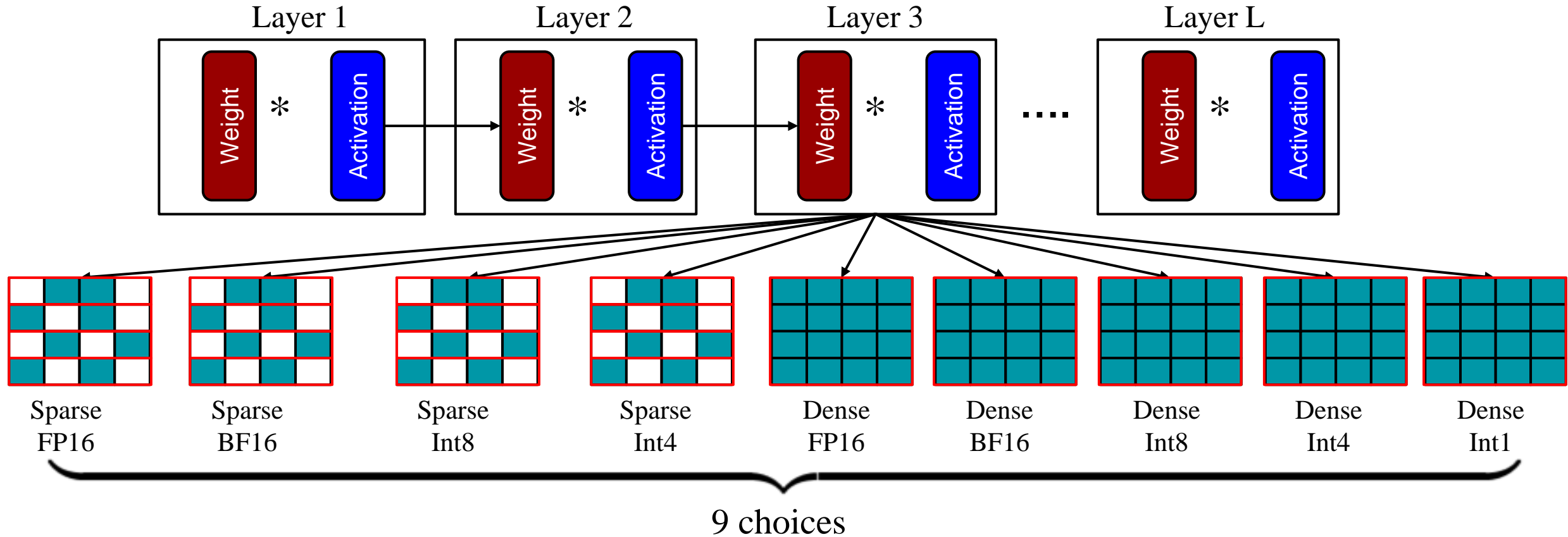- ~ Accuracy

Mishra et al. "Accelerating Sparse Deep Neural Networks"

# *Latency vs Accuracy*

- Example:
  **ResNet50** on **Imagenet Dataset**

- Latency: Time taken for all
  the matrix multiplications
  in the Network

- Accuracy: Percentage of
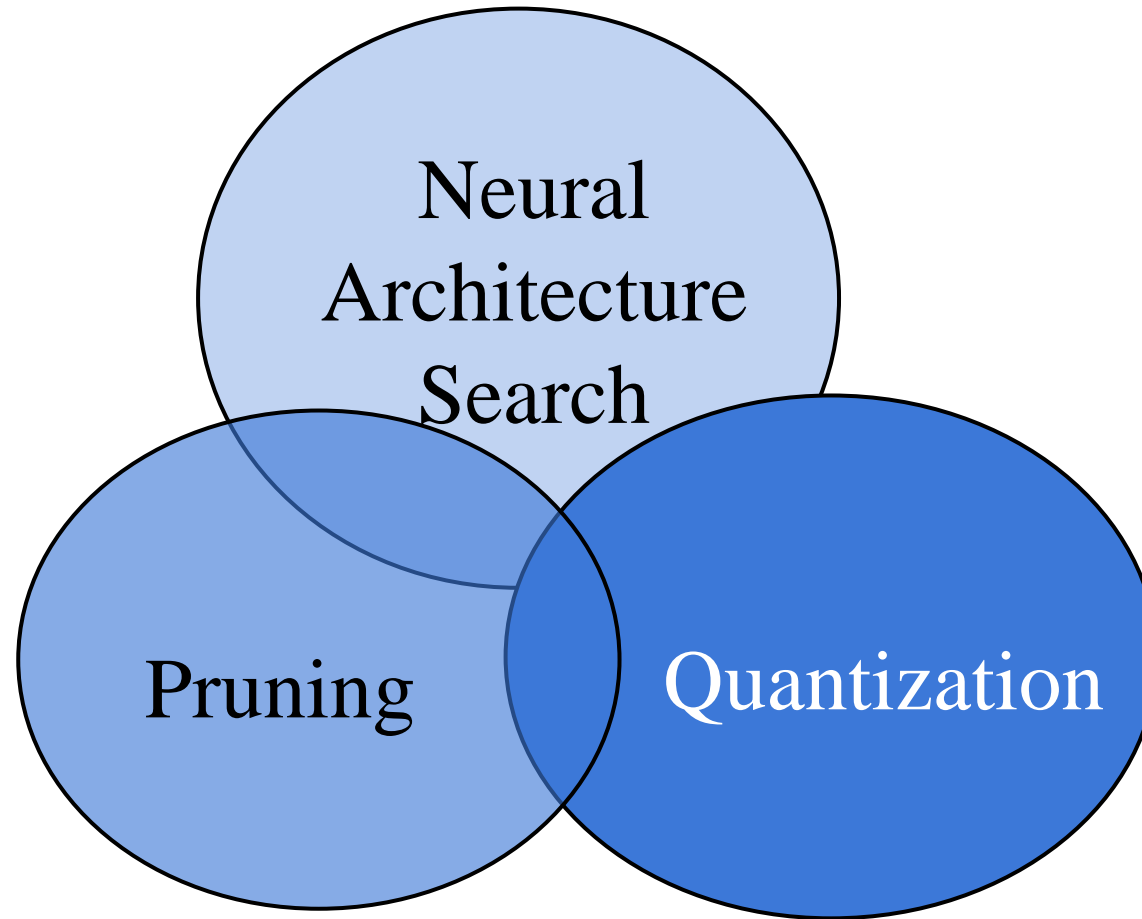  images correctly classified



**Latency**

| | Dense FP16 | Sparse BF16 | Dense Int8 | Dense BF16 | Sparse FP16 | Dense Int4 | Sparse Int8 | Sparse Int4 | Dense Int1 |
|---|---|---|---|---|---|---|---|---|---|
| Latency (ms) | 95.68 | 88.04 | 73.11 | 70.8 | 67.48 | 53.02 | 52.8 | 48.78 | 44.43 |

**Accuarcy**

| | Dense FP16 | Sparse BF16 | Dense Int8 | Dense BF16 | Sparse FP16 | Dense Int4 | Sparse Int8 | Sparse Int4 | Dense Int1 |
|---|---|---|---|---|---|---|---|---|---|
| Accuarcy (%) | 76.1 | 76.05 | 76.05 | 76.05 | 76.2 | 72.5 | 76.2 | 72.6 | 30.5 |

# Problem: *Need for Automated Search Method*
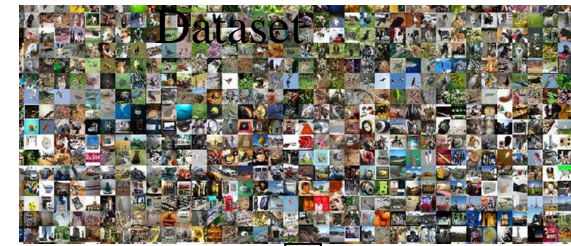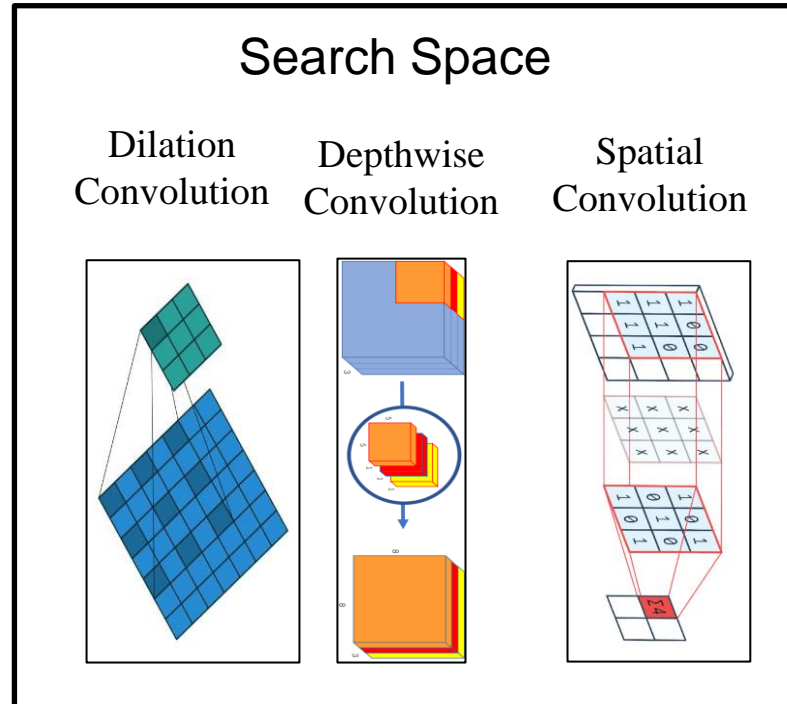


9 choices

- Total Number of distinct combinations for "N" layer $= 9^N$

- Example: For Resnet50, number of combinations $= 9^{50}$ (practically impossible to go through each combination)

- Hence, automated Neural Architecture Search method is required to find optimal combination
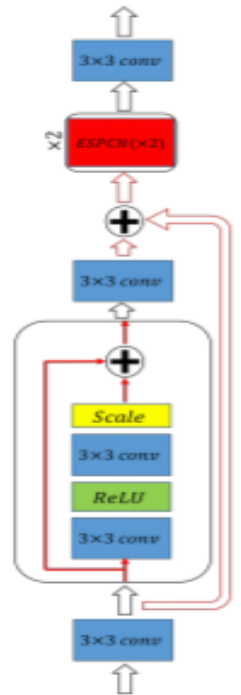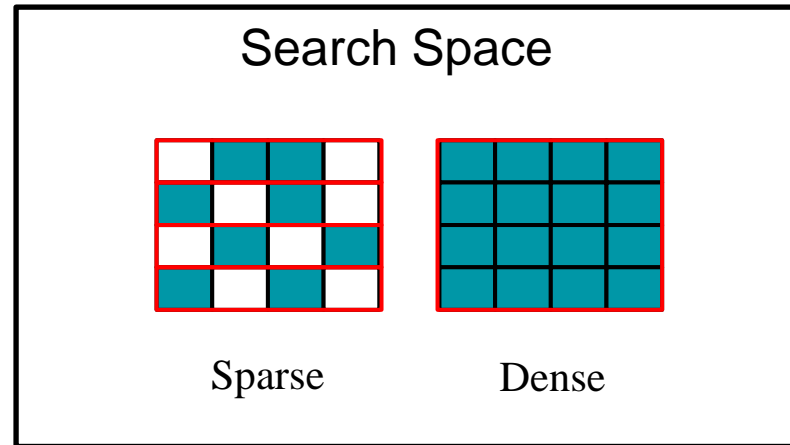
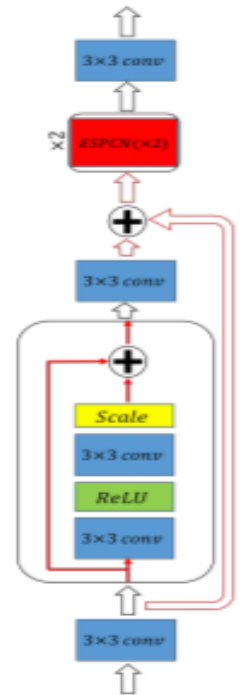*Deep Neural Network Optimization*

# *Neural Architecture Search (NAS)*

Dataset

## Search Space

| Dilation Convolution | Depthwise Convolution | Spatial Convolution |

Hardware-aware

Neural Architecture Search

# Neural Architecture Search (NAS)



Dataset

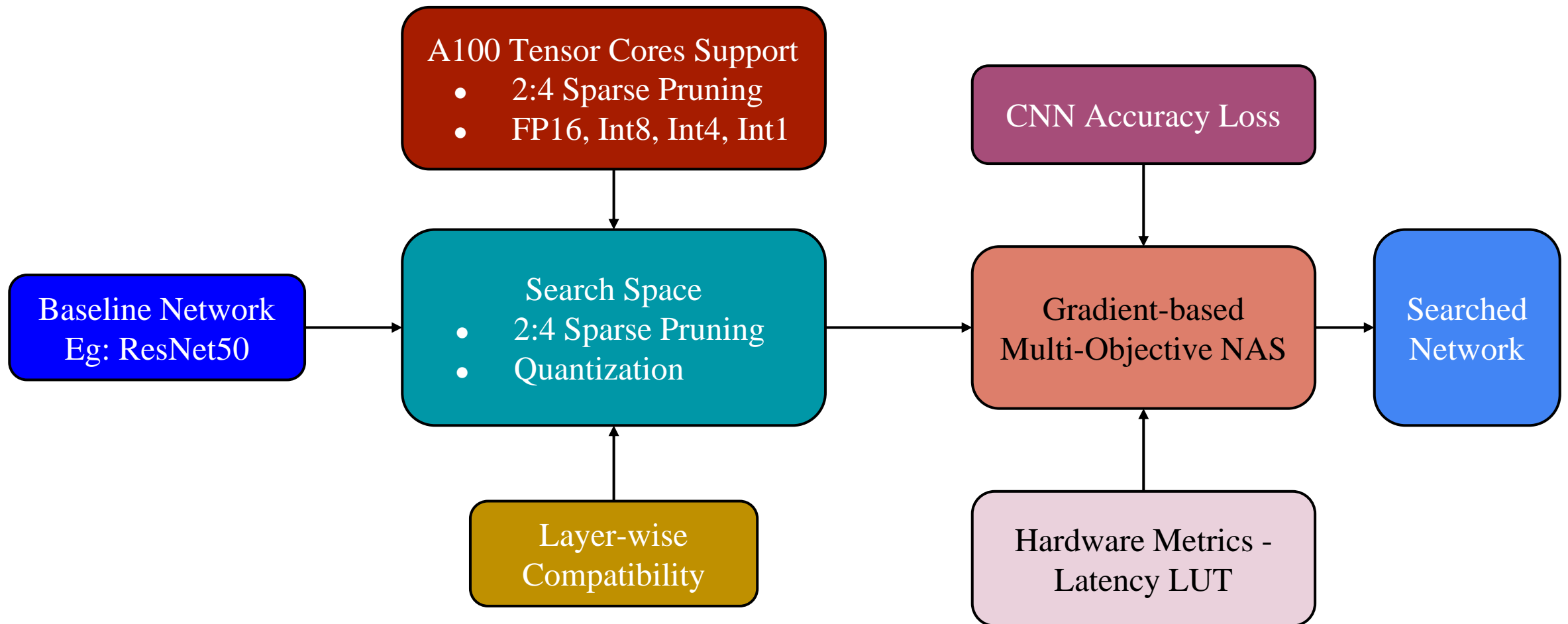Search Space

Sparse          Dense
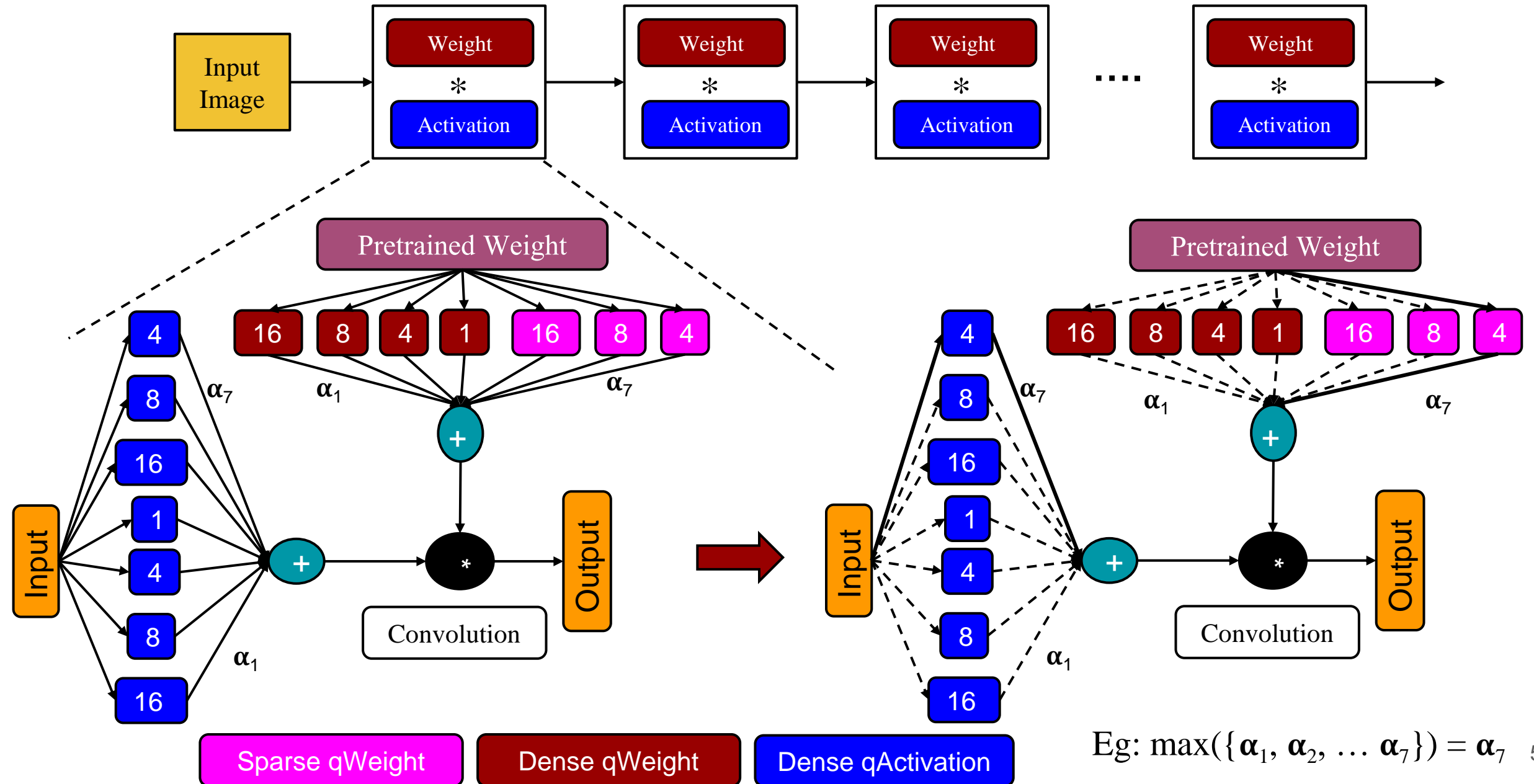
A100 Tensor Core-aware

Neural Architecture Search

55

# *Neural Architecture Search Schematic for A100 GPUs*
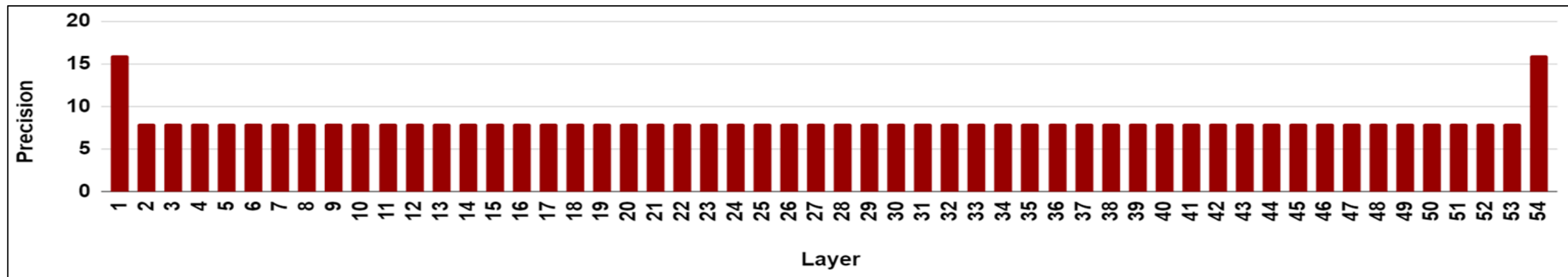
Mixed Sparse & Precision Search (MSPS) - *Sparse and Mixed Precision Quantized Supernetwork*

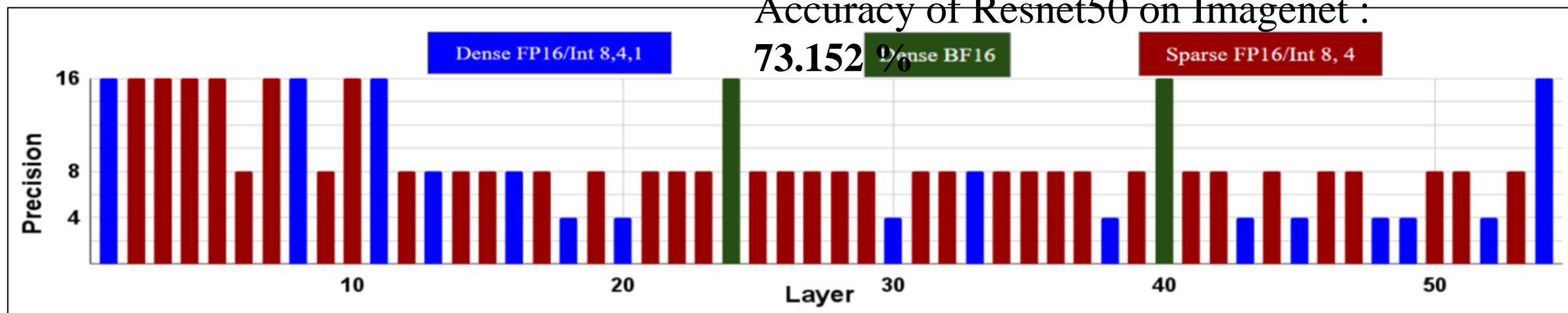Eg: $\max(\{\alpha_1, \alpha_2, \ldots \alpha_7\}) = \alpha_7$

57

# *Results: ResNet50 on Imagenet Dataset*



**Uniform Sparse Int8 Network on Resnet50**

Latency on Nvidia A100 GPU :
**52.8 ms**
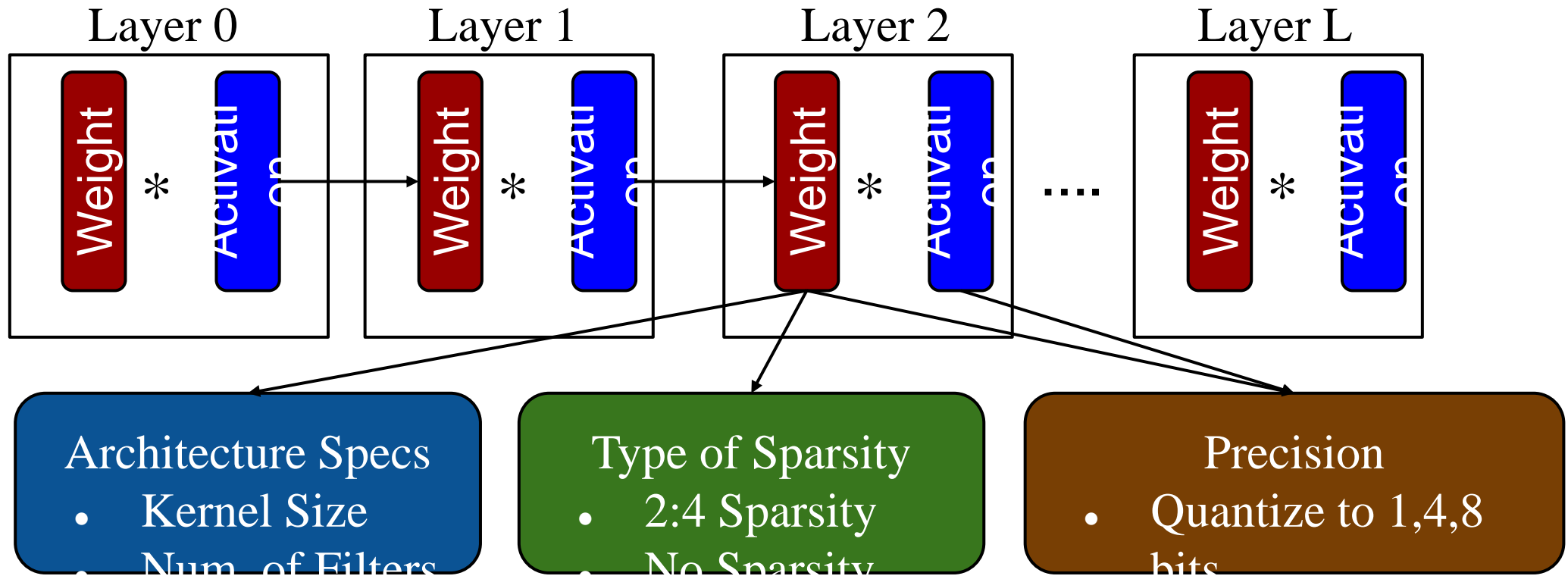Accuracy of Resnet50 on Imagenet :
**73.152 %**



Legend:
- Dense FP16/Int 8,4,1
- Dense BF16
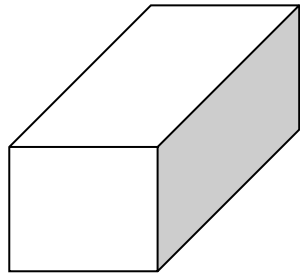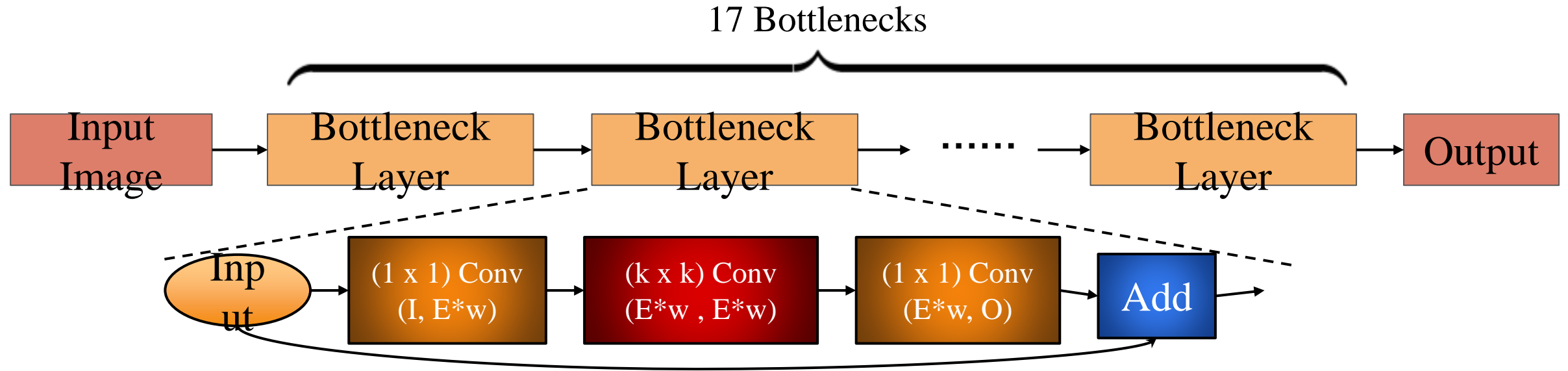- Sparse FP16/Int 8, 4

**Our Mixed Searched Network on Resnet50**

Latency on Nvidia A100 GPU : **52.37 ms**
Accuracy of Resnet50 on Imagenet : **73.32 %**

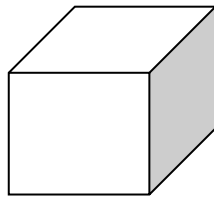# Neural Architecture, Sparsity and Precision Search (ASPS)

- Problem: Simultaneously search for the following dimensions for Ampere 100 Tensor Cores:
  - **Architecture Choices:** Kernel Size and Number of Filters
  - **Optimization:** Mixed Sparse and Precision Combination

*Benchmark: ResNet50 Network on Imagenet dataset*

17 Bottlenecks

Input Image → Bottleneck Layer → Bottleneck Layer → ...... → Bottleneck Layer → Output

Input → (1 x 1) Conv (I, E*w) → (k x k) Conv (E*w , E*w) → (1 x 1) Conv (E*w, O) → Add

kernel = 3
Width Multiplier = 1

kernel = 3
Width Multiplier = 0.5

kernel = 5
Width Multiplier = 1

kernel = 5
Width Multiplier = 0.5

# Architecture and Sparse-Precision Search Space on ResNet50 Network
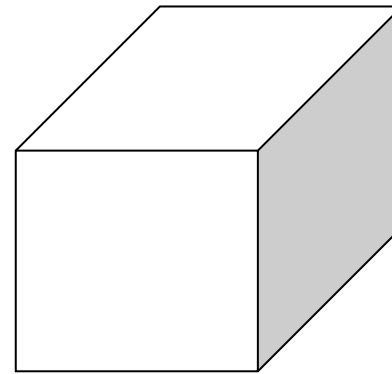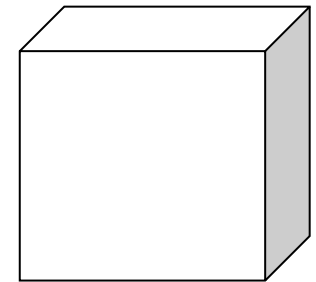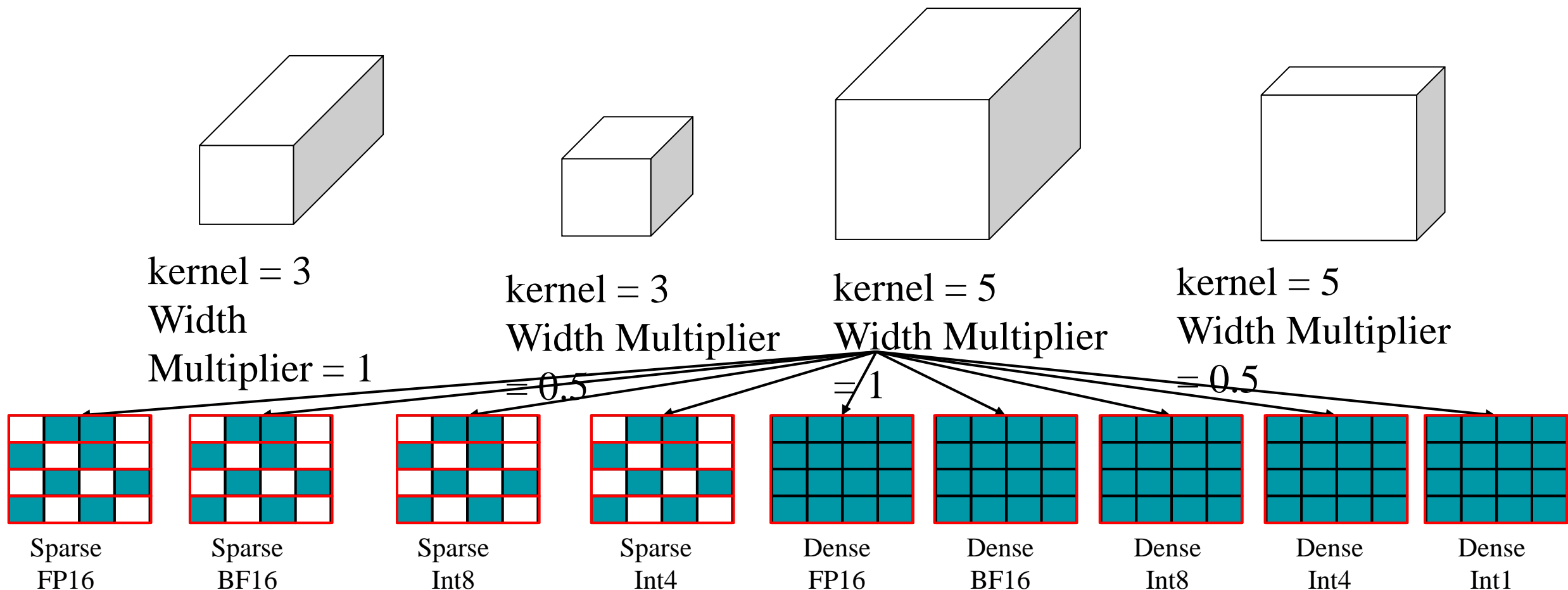


kernel = 3
Width
Multiplier = 1

kernel = 3
Width Multiplier
= 0.5

kernel = 5
Width Multiplier
= 1

kernel = 5
Width Multiplier
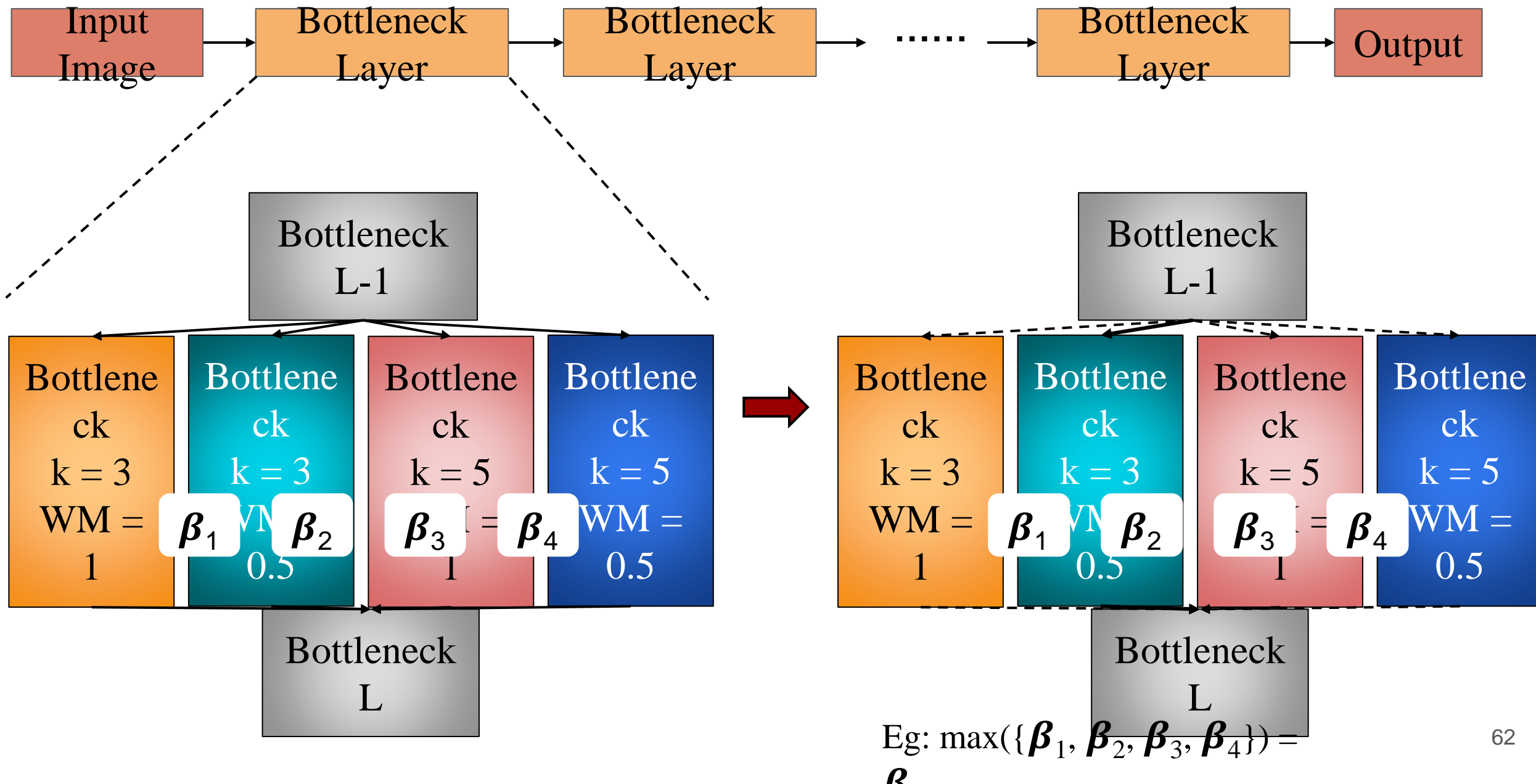= 0.5

| Sparse FP16 | Sparse BF16 | Sparse Int8 | Sparse Int4 | Dense FP16 | Dense BF16 | Dense Int8 | Dense Int4 | Dense Int1 |

Number of Distinct Combinations: $17^{36} * 9^{35}$

# Architecture Search Supernetwork



Eg: max($\{\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \boldsymbol{\beta}_3, \boldsymbol{\beta}_4\}$) = $\boldsymbol{\beta}$

63

*Sampling the best Architecture, Sparse and Precision Combination*

$\max(\{\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \boldsymbol{\beta}_3, \boldsymbol{\beta}_4\}) = \boldsymbol{\beta}_3$

$\max(\{\boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2, \dots \boldsymbol{\alpha}_7\}) = \boldsymbol{\alpha}_7$

*Searched Resnet50 Network on Imagenet dataset*

End-to-end Microarchitecture on Resnet50
Search Space

End-to-end Sparse and Mixed Precision Quantization Combination of the
corresponding Searched Network

# Architecture, Mixed Sparse & Precision Search (ASPS) - Performance Comparison between Uniform Sparse Int8 and our Searched Networks

*CIFAR dataset*

| Configuration | Lat. (ms) | CIFAR10 Acc. (%) | CIFAR100 Acc. (%) |
|---|---|---|---|
| Uniform Sparse Int8 | 5.4 | 94.32 | 74.29 |
| MSPS Most Efficient Model | 5.27 | **94.76** | 74.54 |
| ASPS best Model | **5.19** | 94.38 | **74.7** |

*Imagenet dataset*

| Configuration | Lat. (ms) | Acc. (%) |
|---|---|---|
| Uniform Sparse Int8 | 52.8 | 73.15 |
| MSPS Most Efficient Model | 49.44 | 72.86 |
| ASPS Model ($\lambda = 0$) | 52.8 | 73.42 |
| ASPS best Model | **49.36** | **73.72** |

*Conclusion*

We developed the following methods:

- ***Mixed Sparse and Precision Search (MSPS):*** to search for the optimal weight matrix type (sparse or dense) and precision (FP16, BF16, Int8, Int4, Int1) combination for every layer

- ***Architecture, Sparsity, and Precision Search (ASPS):*** a method to search for better hyperparameters (kernel and filter sizes) along with the matrix type and bit-width in the same loop

- Our searched models on Resnet50 Search Space outperforms the manually designed models on the ImageNet dataset

- For example, our best ASPS model is ~1.1x faster and 0.57% more accurate than the baseline sparse-only Integer 8 ResNet50 network