# CS 744: GAVEL

Shivaram Venkataraman
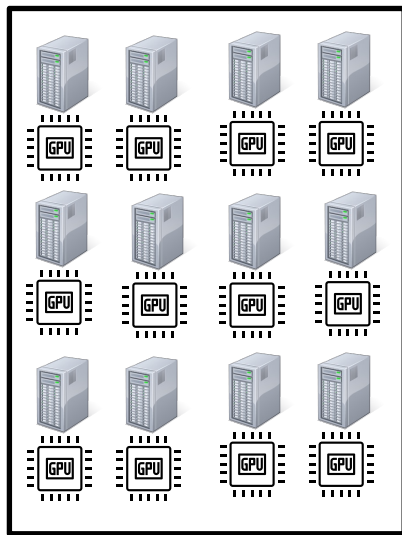
Spring 2024

# ADMINISTRIVIA

- Course project assignments
  - Emails will go out end of this week (March 1)
  - Introductions due March 8th

- Midterm Exam
  - In class on March 14th
  - Includes everything from beginning to the end of scheduling (including INFaaS)

DC as a
Computer

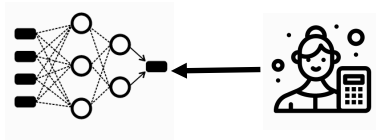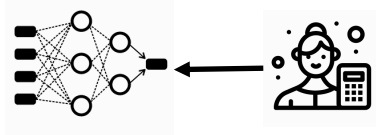# MACHINE LEARNING: TRAINING

Microsoft
University

PyTorch / Pipedream

Resnet
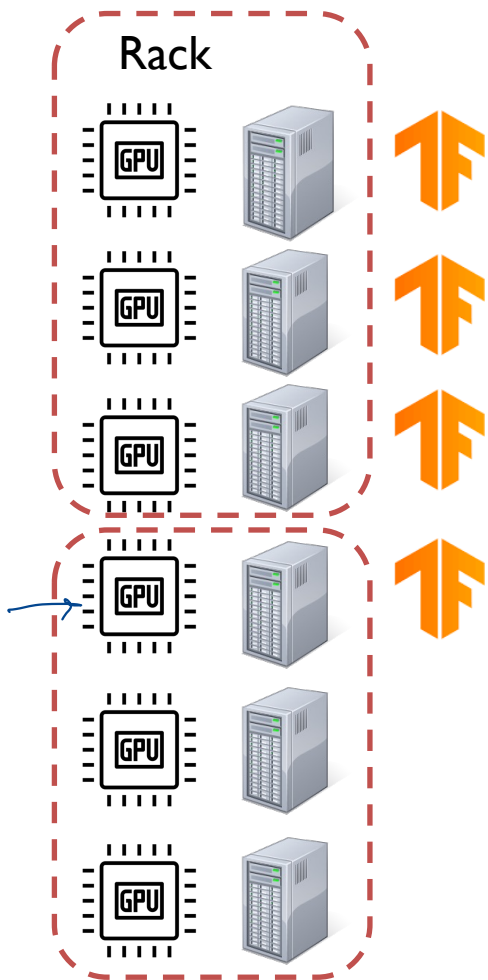
BERT

GPU devices
accelerators

Setup is

more focused

on a specific

workload

# WORKLOAD CHARACTERISTICS

Rack

hours / days ← **Long running tasks**

**Gang scheduling** → run together at the same time

✗ →

**Heterogeneity?** → Hardware generations

Task runs until training completes

k80, V100, A100

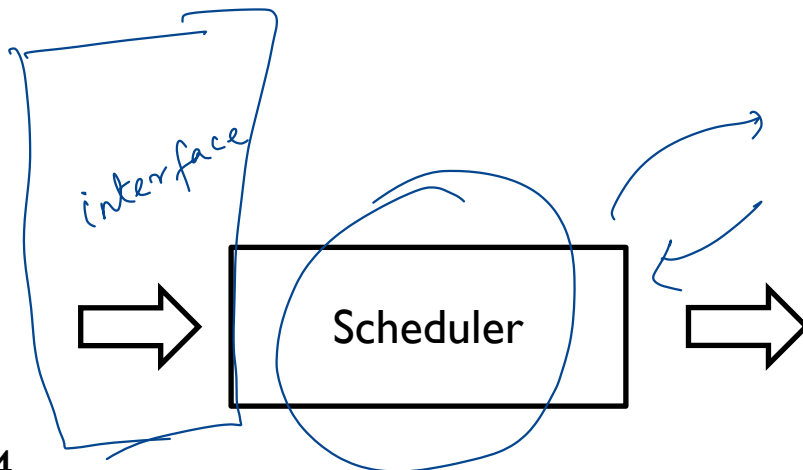# DL SCHEDULER INTERFACE

Job

Run job Resnet18
With BatchSize = 64
on Num GPUs = 4

interface

Scheduler

Goals:
→ Maximize throughput
→ Fairness
  Minimize JCT
  …

allocate these 4 GPUs to job 0

**(a)** Throughput.

**(b)** Dollar-normalized.

speed up with a newer GPU is diff for diff models

MOTIVATION: HETEROGENEITY

P100 or V100

lower than 9.6

worse to use V100 compared to K80

# ADDITIONAL GOALS

utilization

- Support a wide range of objectives

   Minimize makespan → last job finishes

responsiveness

   Average JCT

   Fairness (Sharing incentive) → DRF

   …

General framework

→ Pluggable Policy

→

locality MR    task close to input

- Placement sensitivity/Co-location

All Reduce

close to each
other → same mc

same rack

# GAVEL: SYSTEM DESIGN



Handwritten annotations:
- Simplicity → submits a job
- Accuracy of profile affect the scheduler?
- How fast does each model run on each hardware?
- profiling step input to the scheduler
- What is the overhead?
- Throughput augmented feedback
- profiles using loop

Diagram labels:
- If measurements provided by user
- User objective
- Throughput Estimator
- Throughput tensor
- Policy
- Allocation
- Scheduling Mechanism
- Per-round placement
- Training jobs written in existing frameworks
- PyTorch
- TensorFlow
- Throughput measurements from runs fed back into throughput estimator
- ... V100
- ... P100
- ...

# SCHEDULING POLICY: OPTIMIZATION PROBLEM

$$\text{Maximize}_X \sum_{m \in \text{jobs}} \text{throughput}(m, X)$$

obj: maximize total tput ←

$$\text{throughput}(m, X) = \sum_{\substack{j \in \\ \text{accelerator types}}} T_{mj} \cdot X_{mj}$$

→ weighted by acc. types in cluster

$$0 \leq X_{mj} \leq 1 \qquad \forall (m, j) \quad (1)$$

$$\sum_j X_{mj} \leq 1 \qquad \forall m \qquad (2)$$

$$\sum_m X_{mj} \cdot \text{scale\_factor}_m \leq \text{num\_workers}_j \qquad \forall j \qquad (3)$$

total is ≤ 1

$$X^{\text{example}} = \begin{pmatrix} V100 & P100 & K80 \\ 0.6 & 0.4 & 0.0 \\ 0.2 & 0.6 & 0.2 \\ 0.2 & 0.0 & 0.8 \end{pmatrix} \begin{matrix} \text{job 0} \\ \text{job 1} \\ \text{job 2} \end{matrix}$$

solve this opt problem

get back X ⟶ allocation that achieves this objective

# POLICY: MAX-MIN FAIRNESS

Classic: Weighted max-min fairness based on accelerator hours consumed

$$\text{Maximize}_X \min_m \frac{1}{w_m} X_m$$

$X_0 = 0.33$

$X_1 = 0.33$ → equal share

$X_2 = 0.33$ if all GPUs are same

Gavel: Use weighted normalized effective throughputs

$$\text{Maximize}_X \min_m \frac{1}{w_m} \frac{\boxed{\text{throughput}(m, X)}}{\text{throughput}(m, X_m^{\text{equal}})}$$

$$\text{throughput}(m, X) = \sum_{\substack{j \in \\ \text{accelerator types}}} T_{mj} \cdot X_{mj}$$

weighted sum per acc. type

# EXAMPLE

Profiler or tput
estimator

$$T = \begin{array}{cc} V\,100 & K\,80 \end{array} \begin{pmatrix} 40.0 & 10.0 \\ 12.0 & 4.0 \\ 100.0 & 50.0 \end{pmatrix} \begin{array}{l} \text{job 0} \\ \text{job 1} \\ \text{job 2} \end{array}$$

homogeneous

$$X^{\text{hom.}} = \begin{bmatrix} 0.33 & 0.33 \\ 0.33 & 0.33 \\ 0.33 & 0.33 \end{bmatrix} \begin{array}{l} J_0 \\ J_1 \\ J_2 \end{array}$$

$$X^{\text{het.}} = \begin{array}{cc} V\,100 & K\,80 \end{array} \begin{pmatrix} 0.45 & 0.0 \\ 0.45 & 0.09 \\ 0.09 & 0.91 \end{pmatrix} \begin{array}{l} \text{job 0} \\ \text{job 1} \\ \text{job 2} \end{array} \text{ higher}$$

tput

Eff tput

$$J_0 = 40 \times 0.33 + 10 \times 0.33 = 16.5$$

$$J_1 = 12 \times 0.33 + 4 \times 0.33 = 5.28$$

$$J_2 = 49.5$$

$$J_0 = 0.45 \times 40 = \boxed{18}$$

$$J_1 = 0.45 \times 12 + 0.09 \times 4$$
$$= 5.76$$

$$J_2 = 0.09 \times 100 + 0.91 \times 50$$
$$= 54.5$$

# HIERARCHICAL POLICIES

Single resource
GPUs
"Synergy"

Share physical cluster among sub-organizations
Different policies at levels of hierarchy

**Weighted fairness**  $0.8$  **Organization**  $0.2$

$w_1$   $w_2$

**Product Team**   **Research Team**

**Fairness**

FIFO

Job 1   Job 2   Job 3   Job 4   Job 5

$0.33$   $0.33$   $0.33$

$w_4 = 1.0$   $w_5 = 0$

needs to run first

Solve an LP problem across the organization
Weights constrained by policy within entity
  (e.g., w4 = 1 and w5 = 0)

↳ only start after job 4 finishes

Use water-filling to remove bottlenecked jobs

↳ Prior work in max min fairness

# MECHANISM: ROUND-BASED SCHEDULING

Schedule in "rounds" – every round is ~6 mins → round length
↳ very long running tasks

In every round:

Consider a list of schedulable jobs and $X^{opt}$ (from policy)

Solution opt problem

↳ end of a round, pause stop tasks

Decide which jobs are chosen to run in this round
Track time spent by job m on accelerator type j

Give high priority to jobs which are farthest from $X^{opt}$

→ compute next round

Greedy policy that converges across rounds

→ start jobs for next round

# MECHANISM: PRIORITIES

$$X^{example} = \begin{pmatrix} V100 & P100 & K80 \\ 0.6 & 0.4 & 0.0 \\ 0.2 & 0.6 & 0.2 \\ 0.2 & 0.0 & 0.8 \end{pmatrix} \begin{matrix} \text{job 0} \\ \text{job 1} \\ \text{job 2} \end{matrix}$$

$\rightarrow$ $X^{OPT}$   target allocation

V100 | P100 | K80
$$\begin{pmatrix} 3 & 1 & 0 \\ 1 & 3 & 0 \\ 0 & 0 & 4 \end{pmatrix} \begin{matrix} \text{job 0} \\ \text{job 1} \\ \text{job 2} \end{matrix}$$

rounds_received$_n$

$\longrightarrow$

V100 | P100 | K80
$$\begin{pmatrix} 0.2 & \mathbf{0.4} & 0 \\ 0.2 & 0.2 & \infty \\ \infty & 0 & 0.2 \end{pmatrix} \begin{matrix} \text{job 0} \\ \text{job 1} \\ \text{job 2} \end{matrix}$$

priorities$_n$

Priority $= \dfrac{x_{jm}}{r\,r_{jm}}$   element wise division

V100 | P100 | K80
$$\begin{pmatrix} 3 & \mathbf{2} & 0 \\ 1 & 3 & \mathbf{1} \\ \mathbf{1} & 0 & 4 \end{pmatrix} \begin{matrix} \text{job 0} \\ \text{job 1} \\ \text{job 2} \end{matrix}$$

rounds_received$_{n+1}$

Jobs placed on resources where they have high priority (marked in **red**)

Job 1   K80

Job 2   V100

Job 0   P100

- Switch frequently
- Doesn't work jobs shorter than 1 round

# SUMMARY

DL training workloads properties

Clusters with mix of accelerators

Gavel: Framework to capture many scheduling goals

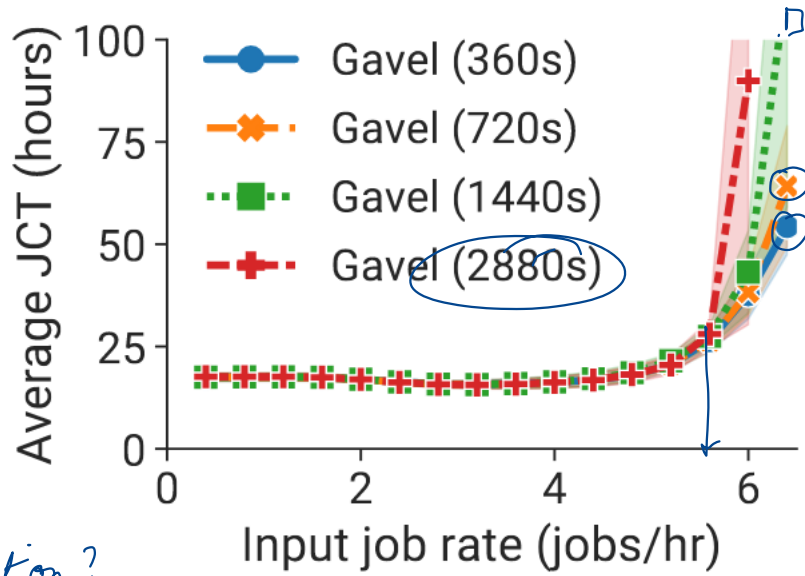Mechanism based on round-based assignments

# DISCUSSION

https://forms.gle/pYnFErGi54HEHcuj7

What are some similarities or differences between Mesos/DRF and DL schedulers like Gavel?

round length



Small round length

→ overhead of
computing allocation?

→ Pre-emption
overhead ( < 10%? )

→ Increase seems
to be linear!

round length is
high JCT

is high for
> 5.5 jobs/hr

→ round length is large

→ fragmentation
jobs finish early
in the round

# NEXT STEPS

Next Class: INFaaS

Course Project Introductions!