# Tetris

*Multi-Resource Packing for Cluster Schedulers*

Robert Grandl, Ganesh Ananthanarayanan,

Srikanth Kandula, Sriram Rao, Aditya Akella

Microsoft® Research

WISCONSIN
UNIVERSITY OF WISCONSIN–MADISON

# Performance of cluster schedulers

**We find that:**

- *Resources are fragmented i.e. machines run below capacity*
- *Even at 100% usage, goodput is smaller due to over-allocation*
- *Pareto-efficient multi-resource fair schemes do not lead to good avg. performance*

## Tetris
**Up to 40% improvement in makespan[1] and job completion time with near-perfect fairness**

# Findings from Bing and Facebook traces analysis

**Applications have (very) diverse resource needs**

- Tasks need *varying amounts of each resource*

- *Demands* for resources *are weakly correlated*

**Multiple resources become tight**

**This matters, because no single bottleneck resource in the cluster:**

- E.g., enough cross-rack network bandwidth to use all cores

## Upper bound on potential gains

- *Makespan reduces by ≈ 49%*
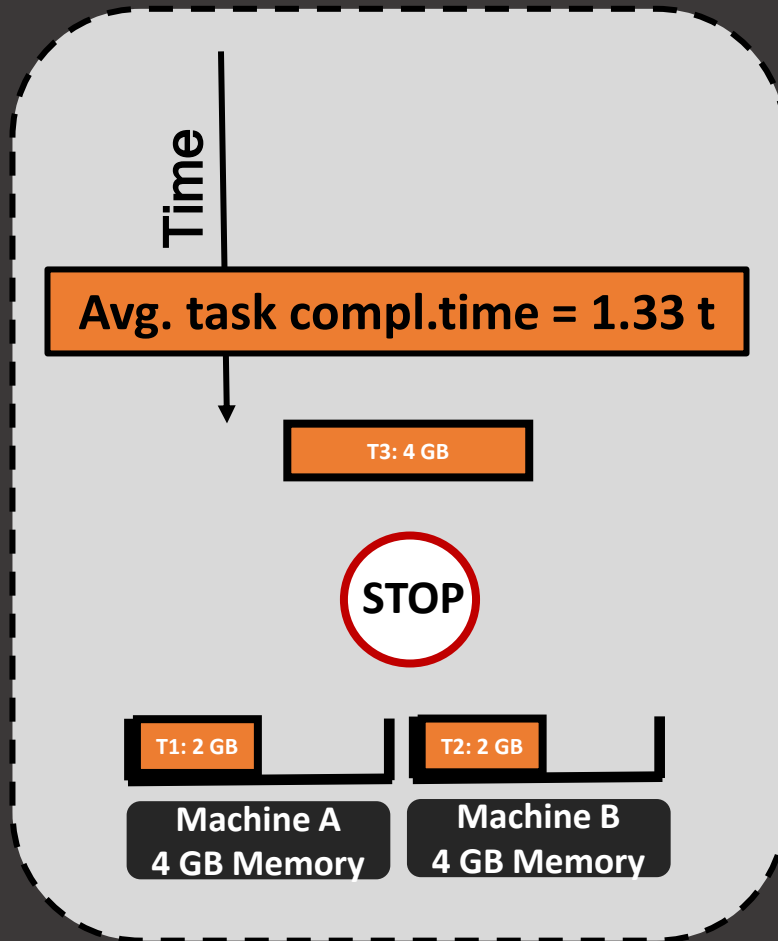
- *Avg. job completion time reduces by ≈ 46%*

**Production schedulers *neither pack tasks nor consider all* their relevant *resource demands***
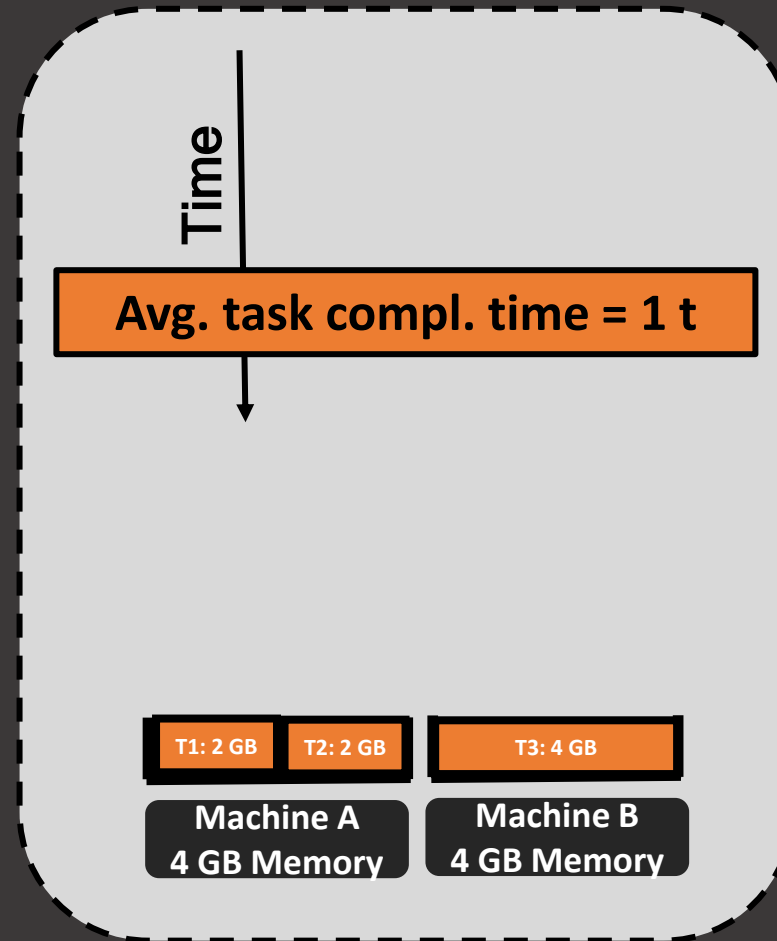
**#1 Resource Fragmentation**

**#2 Over-allocation**

# Resource Fragmentation (RF)



Current Schedulers

Current Schedulers

Time

Avg. task compl.time = 1.33 t

T3: 4 GB

STOP

T1: 2 GB     T2: 2 GB

Machine A
4 GB Memory     Machine B
4 GB Memory

"Packer" Scheduler

Time

Avg. task compl. time = 1 t

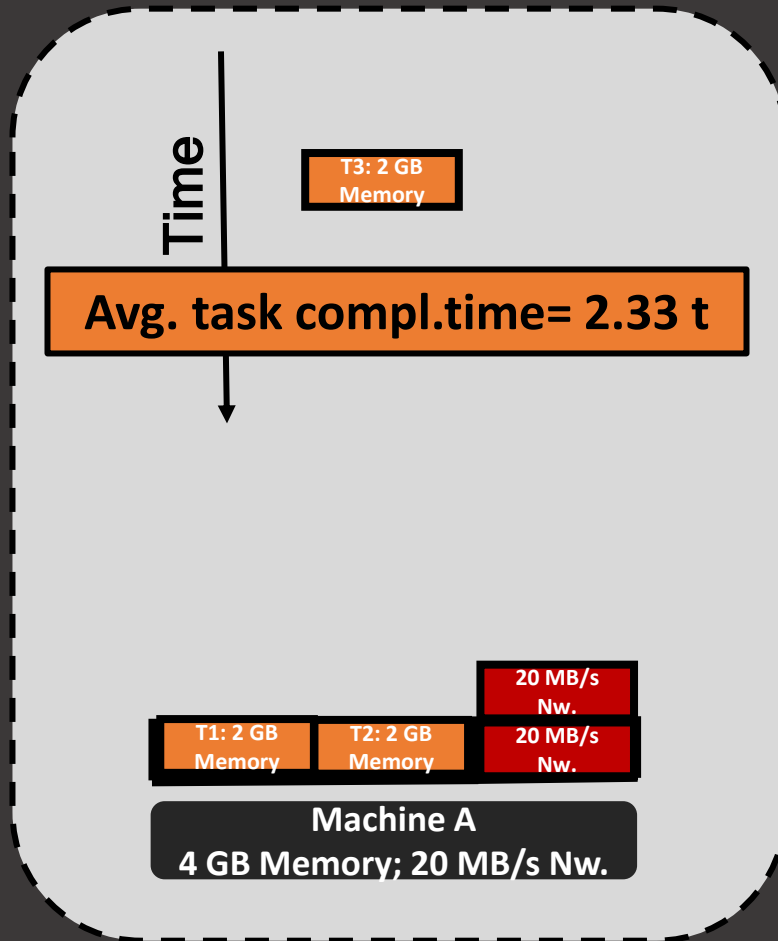T1: 2 GB   T2: 2 GB     T3: 4 GB

Machine A
4 GB Memory     Machine B
4 GB Memory

Are not explicit about packing.

Allocate resources per **slots, fairness.**

RF increase with the number of resources being allocated !

5

# Multi-resource Fairness Schemes do not solve the problem

**Example in paper**
**Packer vs. DRF: makespan and avg. completion time improve by over 30%**

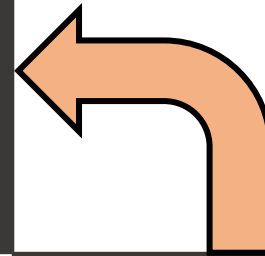Work Conserving  != no fragmentation, over-allocation

Pareto[1] efficient != performant

- **Treat cluster as a big bag of resources**
  - ❏ *Hides the impact of* **resource fragmentation**

- **Assume job has a fixed resource profile**
  - ❏ *Different tasks in the same job have* different demands
  - ❏ *How the job is* scheduled impacts *jobs' current resource* profiles
  - ❏ *Can schedule to create* complementarity

[1]*no job can increase its share without decreasing the share of another*

# Current Schedulers

1. **Resource Fragmentation**

2. **Over-Allocation**

3. **Fair allocations sacrifice performance**

# Competing objectives

**Cluster efficiency**

**vs.**

**Job completion time**

**vs.**

**Fairness**

# Tetris

# # 1

Pack tasks along multiple resources to improve cluster efficiency and reduce makespan

# Theory

## Multi-Resource Packing of Tasks
### similar to
## Multi-Dimensional Bin Packing

**APX-Hard[1]**

*Avoiding fragmentation looks like*:

❑ Tight bin packing

❑ Reduce # of bins → *reduce makespan*

*[1]APX-Hard is a strict subset of NP-hard*
*Balls could be tasks*
*Bin could be machine, time*

# Practice

Existing heuristics do not directly apply:

❑ Assume balls of a fixed size

❑ Assume balls are known apriori

- vary with time / machine placed
- elastic

- cope with online arrival of jobs, dependencies, cluster activity

# Tetris

## A packing heuristic

Fit

- Tasks resources demand vector  **<**  ▪ Machine resource vector

**Alignment score (A)**

## "A" works because:

**1.** Check for fit to ensure **no over-allocation**   ✓ Over-Allocation

**2.** Bigger balls *get bigger scores*

✓ Resource Fragmentation

**3.** Abundant resources *used first*

# Tetris

# # 2

Faster average job completion time

# Tetris

**CHALLENGE**

**Q**: What is the shortest "~~remaining time~~" ?

"**remaining work**" =  remaining # tasks
&
tasks' resource demands
&
tasks' durations

**Job Completion Time Heuristic**

## A job completion time heuristic

- Gives a score **P** to every job

- Extended SRTF to incorporate multiple resources

**Shortest Remaining Time First[1] (SRTF)**

*schedules jobs in ascending order of their remaining time*

13

[1]SRTF – M. Harchol-Balter et al. Connection Scheduling in Web Servers [USITS'99]

**Tetris**

Job Completion Time Heuristic

**CHALLENGE**

A: delays job completion time

?

Packing Efficiency

Completion Time

P: loss in packing efficiency

**Combine A and P scores !**

```
1:   among J runnable jobs
2:     score (j) = A(t, R)+ε P(j)
3:            max task t in j, demand(t) ≤ R (resources free)
4:  pick j*, t* = argmax score(j)
```

# Tetris

# # 3

Achieve performance and fairness

# Tetris

**Performance and fairness do not mix well in general**

**But** ….
We can get "perfect fairness" and much better performance

- **Packer** says: "*task T should go next to improve packing efficiency*"

- **SRTF** says: "schedule *job J to improve avg. completion time*"

- **Fairness** says: "**this set of *jobs** should be scheduled next*"

**Possible to satisfy all three**
**In fact, happens often in practice**

# 3

Fairness
Heuristic

# Tetris

**Fairness is not a tight constraint**

- **Lose a bit of fairness for a lot of gains in performance**
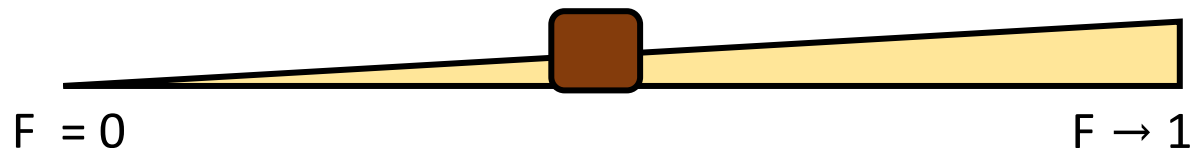- **Long term fairness not short term fairness**

**Heuristic**

- **Fairness Knob, F $\in [0, 1)$**

**Fairness Heuristic**

Pick the ***best-for-perf.*** task from among $\lceil 1-F \rceil$ fraction of jobs furthest from fair share

☐ Most unfair
☐ Most efficient scheduling

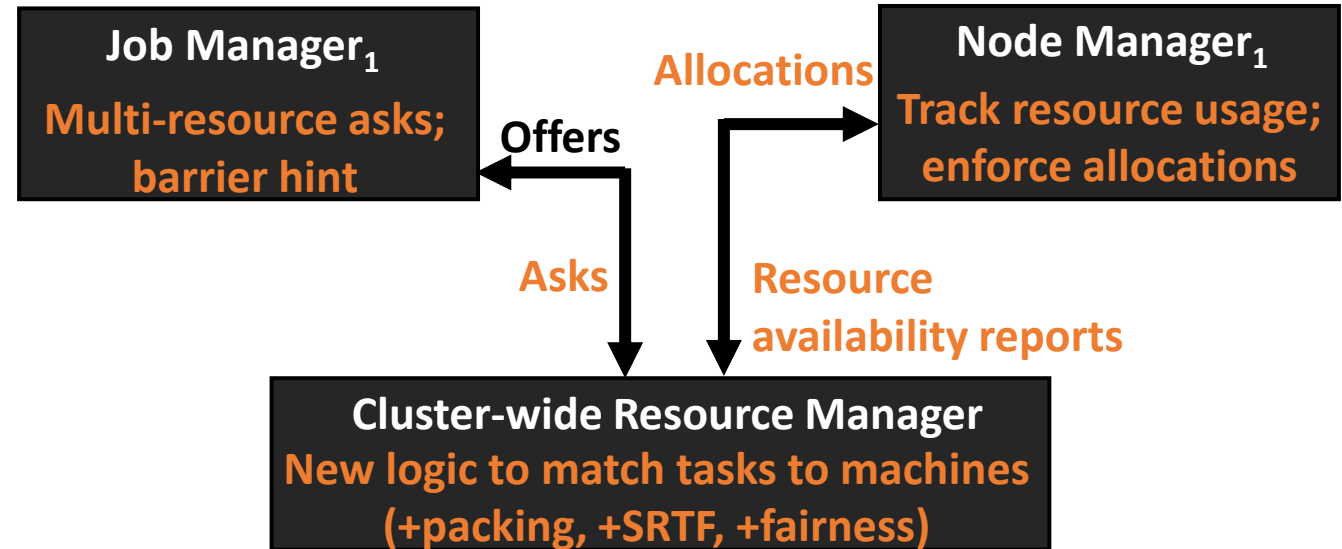Close to perfect fairness

F = 0                    F → 1

# Putting it all together

**We saw:**

- **Packing efficiency**
- **Prefer small remaining work**
- **Fairness knob**

**Other things in the paper:**

- **Estimate task demands**
- **Deal with inaccuracies, barriers**
- **Other cluster activities**

**Job Manager$_1$**
**Multi-resource asks;
barrier hint**

**Offers**

**Allocations**

**Node Manager$_1$**
**Track resource usage;
enforce allocations**

**Asks**

**Resource
availability reports**

**Cluster-wide Resource Manager**
**New logic to match tasks to machines
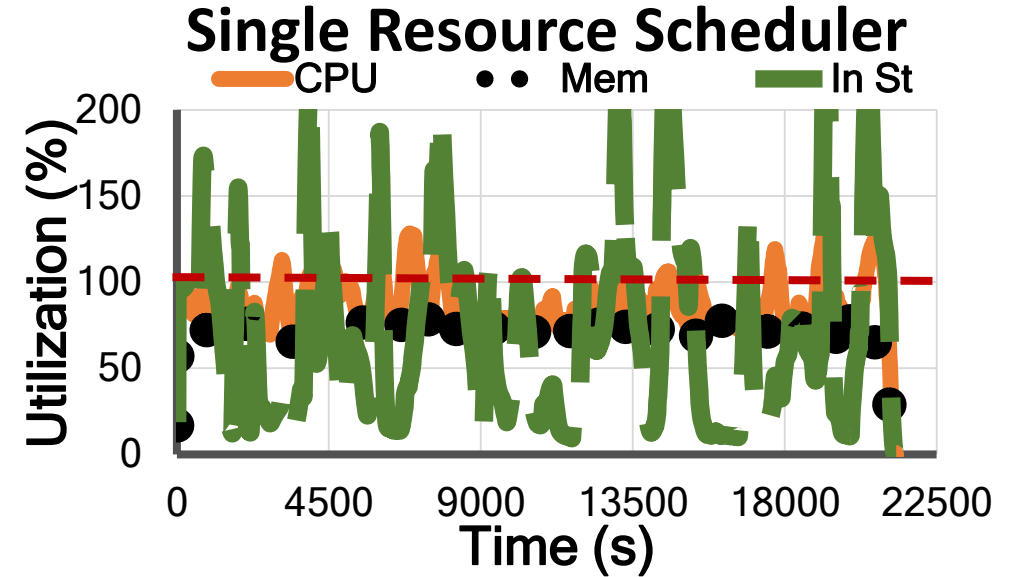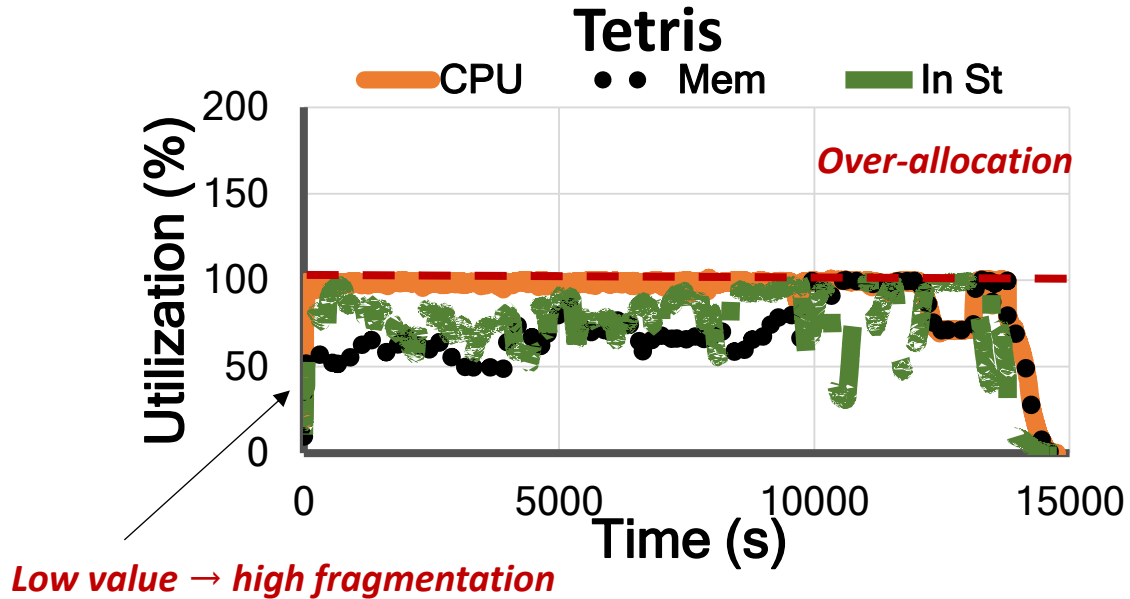(+packing, +SRTF, +fairness)**

## Yarn architecture

Changes to add Tetris(shown in orange)

# Evaluation

- Implemented in Yarn 2.4

- 250 machine cluster deployment

- Bing and Facebook workload

**Tetris**

CPU ●● Mem ▬ In St

Over-allocation

Utilization (%) — Time (s)

Low value → high fragmentation

**Single Resource Scheduler**

CPU ●● Mem ▬ In St

Utilization (%) — Time (s)

| Tetris vs. | Makespan | Avg. Job Compl. Time |
|---|---|---|
| Single Resource Scheduler | 29 % | 30 % |
| Multi-resource Scheduler | 28 % | 35% |

**Gains from**

- *avoiding fragmentation*
- *avoiding over-allocation*

# Fairness

- *quantifies the extent to which Tetris adheres to fair allocation*

|  | Makespan | Job Compl. Time | Avg. Slowdown [over impacted jobs] |
|---|---|---|---|
| **No Fairness F = 0** | 50 % | 40 % | 25 % |
| **Full Fairness F → 1** | 10 % | 23 % | 2 % |
| **F = 0.25** | 25 % | 35 % | 5 % |

# Tetris

**Pack efficiently along multiple resources**

**Prefer jobs with less "remaining work"**

**Incorporate Fairness**

- Combine heuristics that *improve packing efficiency* with those that *lower* average *job completion time*

- Achieving desired amounts of *fairness can coexist* with improving cluster performance

- Imp...
sho...

*We are working towards a Yarn check-in*

http://research.microsoft.com/en-us/UM/redmond/projects/tetris/