

CS 744: MESOS

Shivaram Venkataraman Fall 2020

ADMINISTRIVIA

- Assignment I: How did it go? Fill out the poll !
- Assignment 2 out tonight \longrightarrow ML distributed
- Project details
 - Create project groups ~ 3 students
 - Bid for projects/Propose your own next week
 - Work on Introduction

L> 1-2 page

- Check-in - Final report and poster ression

COURSE FORMAT

Paper reviews

"Compare, contrast and evaluate research papers"

Discussion





BACKGROUND: OS SCHEDULING



CI USTER SCHEDULING Scale - large number of mailines La one scheduler? Fairness -> pinilarily petween OS & cluster space sharing MR space spark fault tolerance (Preferences (placement, constraint) Multi - core time ware NVMA juling MR amail Monc







CONSTRAINTS

Examples of constraints :

Data locality -s soft GPU machines -> hard

Constraints in Mesos:

La franeworks can reject offer La "Filters" - Boolean functions

DESIGN DETAILS





PLACEMENT PREFERENCES

Lo If you more proneworks with prep than machines available in the cluster What is the problem?

How do we do allocations?

La weighted lottery scheme make offers in size to the overall resources that proportional in size to the overall resources that a framework needs

CENTRAI IZED VS DECENTRAI IZED

Centralized

-> Scalability ~100% of frameworks each ~100% of apps Optimel solution

Decentralized

handle new frameworks i Complexity for framework developer

CENTRALIZED VS DECENTRALIZED

Framework complexity



Inter-dependent framework



COMPARISON: BORG - Google

Single centralized scheduler

Requests mem, cpu in cfg Priority per user / service

La Better packing

Support for quotas / reservations



SUMMARY

- Mesos: Scheduler to share cluster between Spark, MR, etc.
- Two-level scheduling with app-specific schedulers
- Provides scalable, decentralized scheduling
- Pluggable Policy ? Next class!

DISCUSSION

https://forms.gle/urHSeukfyipCKjue6

What are some problems that could come up if we scale from 10 frameworks to 1000 frameworks in Mesos?

- Fragmentation / starvation odds go up -> Master bottleneck? La it takes time to wait for frameworks to reply Mesos master hes soft state -> Pre-emption? Yes. -> Failure recovery takes longer? why? ~ unclear ?



List any one difference between an OS scheduler and Mesos La Motivation part of the lecture Data locality Lo Input (HDFS) Lo Memory (Executor) Spark on oversubscribed dusters Lo long running jobs Executor -> cache RPDs 1 -> "shuffle files" on local disk Shulffle (Eventor) output - gete pre-empted -> rache is bloom away A shuffle files long lived Deyond "gamanteed slave" Coarse Grained Executor Backend

resource offers I how does it perform better"? "ramp - up " (i) Time to schedule. - (ii) Time to completion optimal policy ¢ (ພິດ Borg Comparisons YARN, mth



NEXT STEPS

Pre-emption La redo "work"

Next class: Scheduling Policy Assignment 2 mill be released

Further reading

- https://www.umbrant.com/2015/05/27/mesos-omega-borg-a-survey/
- https://queue.acm.org/detail.cfm?id=3173558

