

PFLOCK Report

Andres Calderon

University of California, Riverside

July 10, 2020

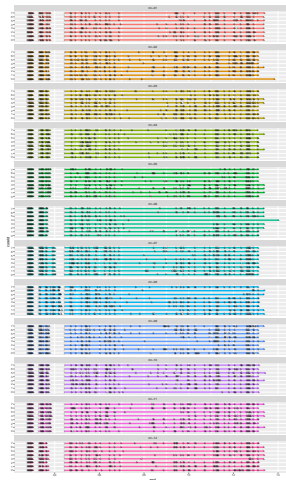
Task ID and Core/Thread mapping

- ▶ Technically a mapping between taskId and core/thread is not possible¹...
 - ▶ Apache Spark will run a Task as a Runnable process managed by the JVM.
 - ▶ JVM interacts directly with the OS to deal with those processes.
 - ▶ There is no guarantee an OS process run in only one core.
- ▶ However, based on the lifetime of a Task (launchTime and finishTime) we can explore how concurrent are the Tasks...

¹More info at <https://tinyurl.com/ycjvxtne>

Task histogram

Task distribution²



² click to enlarge...

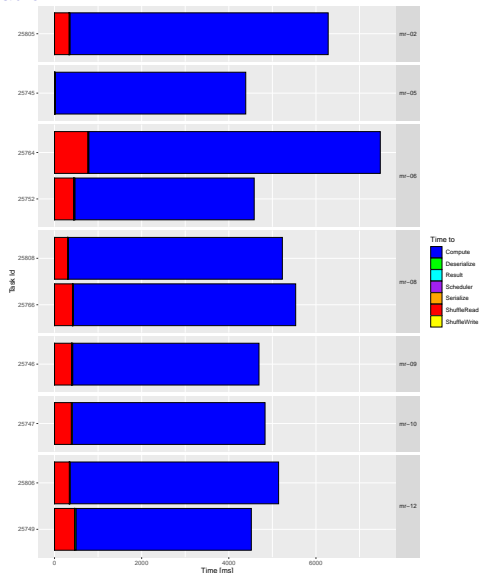
Node usage

	n
Min.	1.00
1st Qu.	8.00
Median	8.00
Mean	7.95
3rd Qu.	8.00
Max.	13.00

1	2	3	4	5	6	7	8	9	10	11	12
310	315	489	788	374	697	198	121763	4721	616	77	11

Task histogram

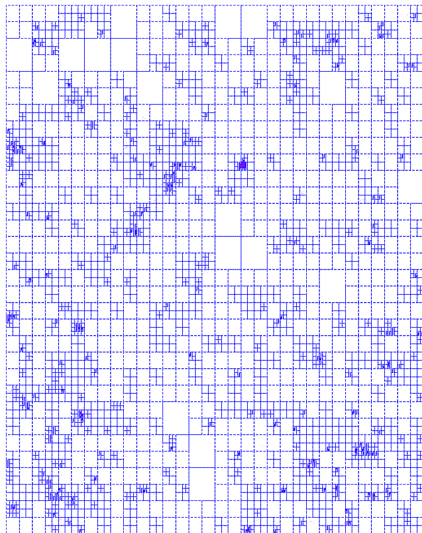
Top 10 tasks by duration



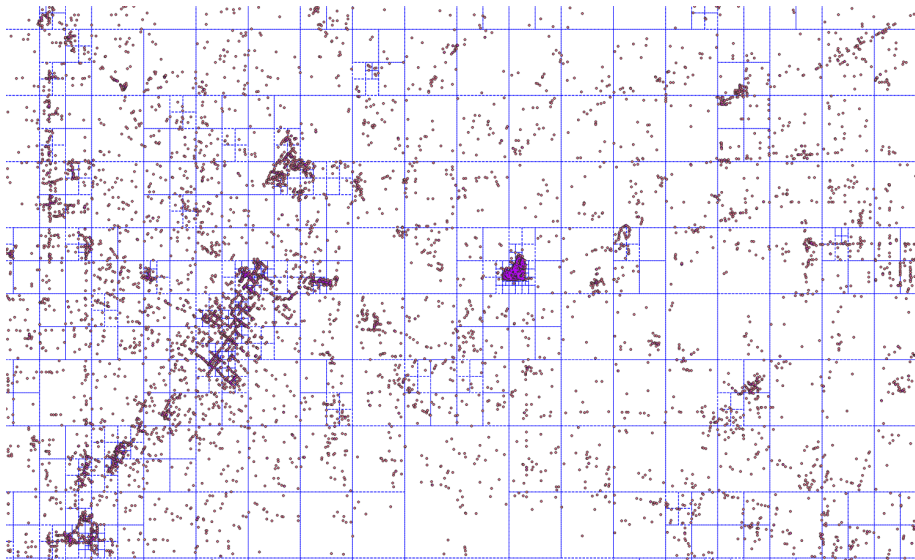
Task-Data Mapping

- ▶ Class `TaskContext` allows to collect information about the task execution on runtime.
- ▶ Methods `partitionId()` and `taskAttemptId()` give us the ID of the RDD partition and the task it was assigned.

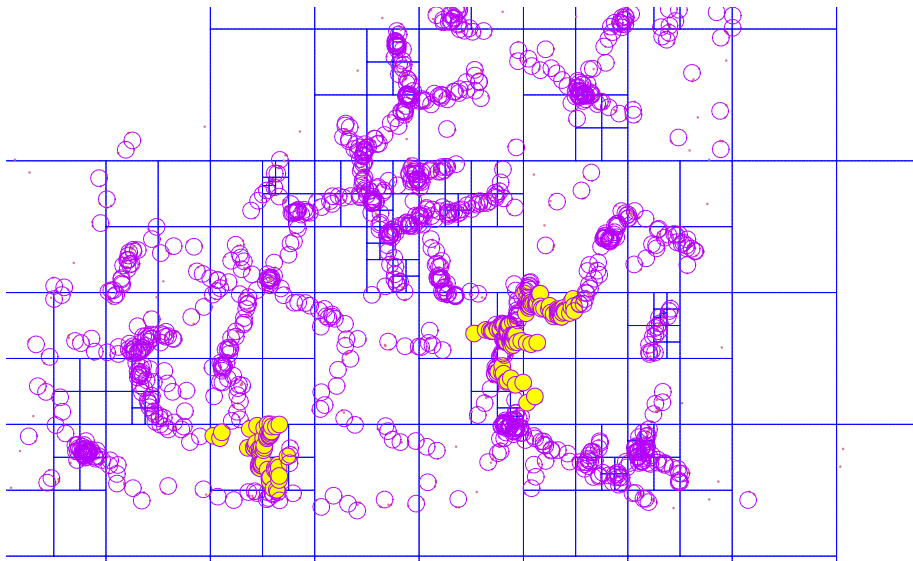
Task-Data Mapping



Task-Data Mapping



Task-Data Mapping



What's next

- ▶ Measure the impact of very small partitions.
- ▶ Still dealing with shuffle cost.
- ▶ Give a try with a different dataset.