



Profiling and optimizing a Spark job with Babar

BENOIT HANOTTE, BERLIN BUZZWORDS 2018

criteo.

“Connecting shoppers to the things they need and love.”





500TB

INGESTED DAILY

15PB

READ DAILY

4000+

HADOOP NODES

criteo.



...

criteo.

Context

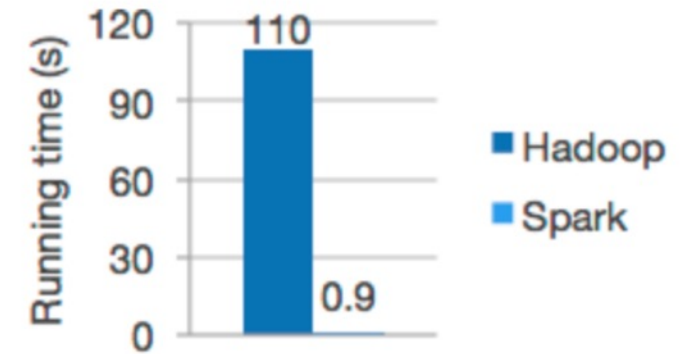
Teams moving from MapReduce-based frameworks to newer alternatives

Spark very popular, advertises faster & more efficient processing

Speed

Run workloads 100x faster.

Apache Spark achieves high performance for both batch and streaming data, using a state-of-the-art DAG scheduler, a query optimizer, and a physical execution engine.

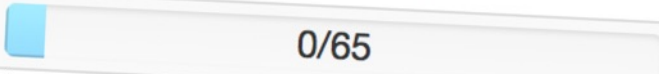


Logistic regression in Hadoop and Spark

Source: <http://spark.apache.org> (June 2018)

Context

But this is not what we experienced...

Duration ▾	Tasks: Succeeded/Total
40.5 h	 0/65

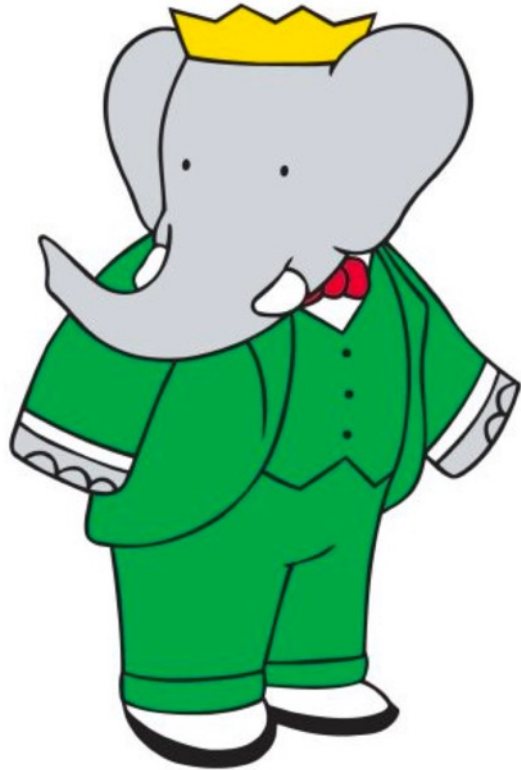
ExecutorLostFailure (executor 21 exited caused by one of the running tasks) Reason: Executor heartbeat timed out after 150977 ms

java.lang.OutOfMemoryError: GC overhead limit exceeded

ExecutorLostFailure (executor 19 exited caused by one of the running tasks) Reason: Container killed by YARN for exceeding memory limits. 17.1 GB of 17 GB physical memory used. Consider boosting spark.yarn.executor.memoryOverhead.

Jobs were unstable, teams were solving it by allocating too much resources to their jobs.

Babar



Profiler for **distributed applications** on **Hadoop**

Works with any JVM framework



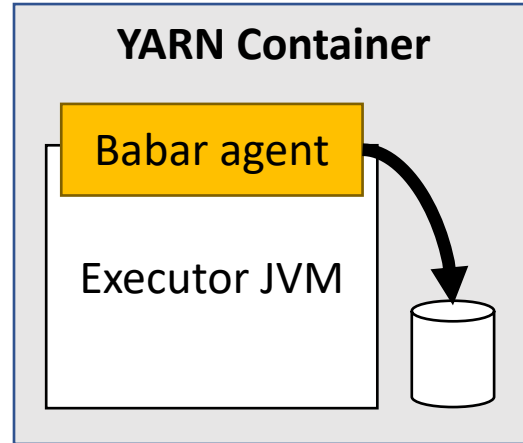
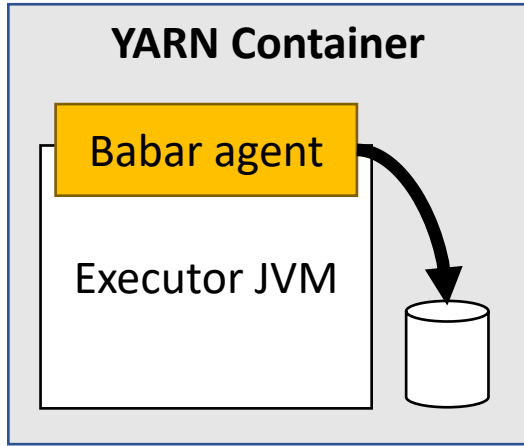
...

- Easy to get started, no infrastructure required
- Made for distributed applications
- Works on 10000+ containers apps
- Exports ready-to-use graphs as HTML file



<https://github.com/criteo/babar>

Babar



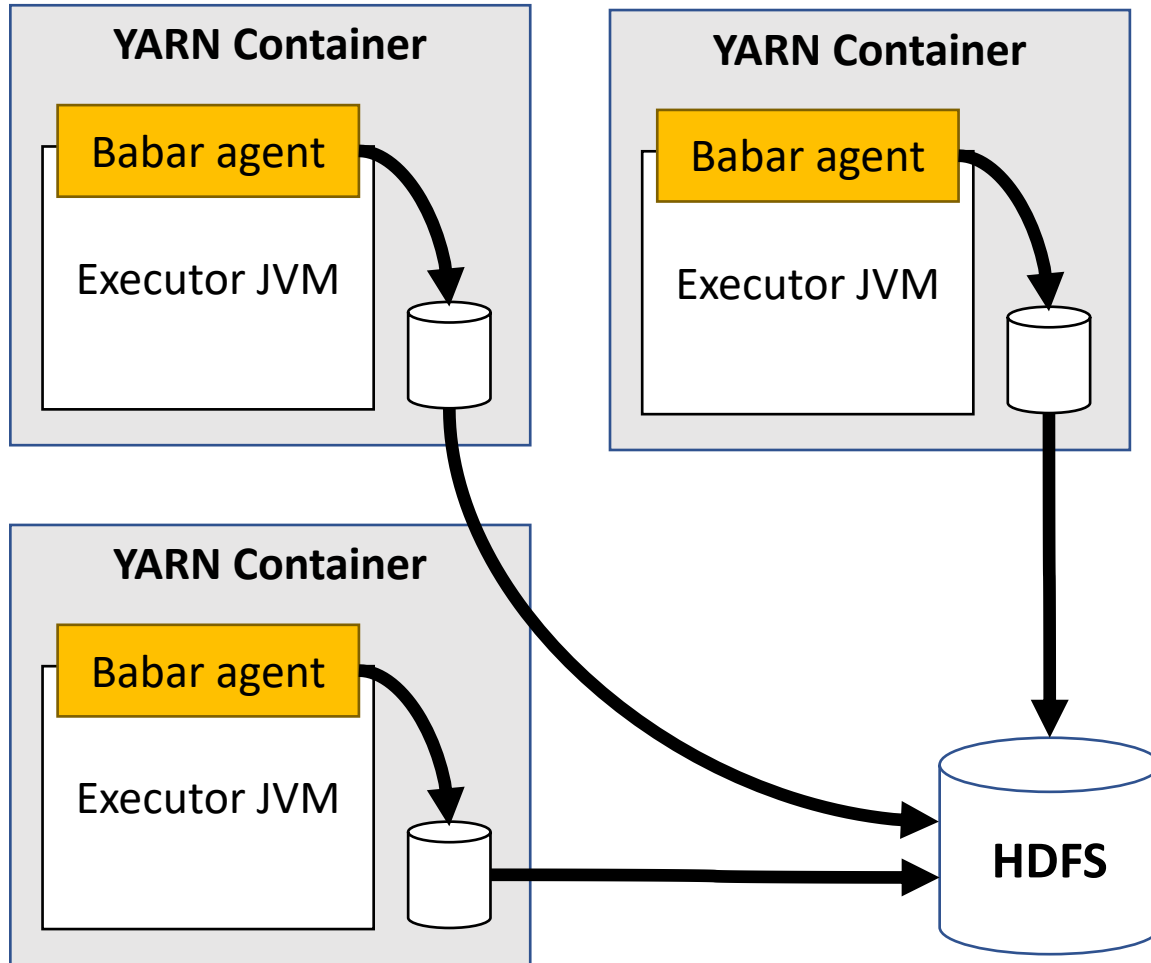
Babar agent instruments the executor JVM, logs metrics to local FS.

Agent can be distributed with spark-submit (no installation).

Metrics from:

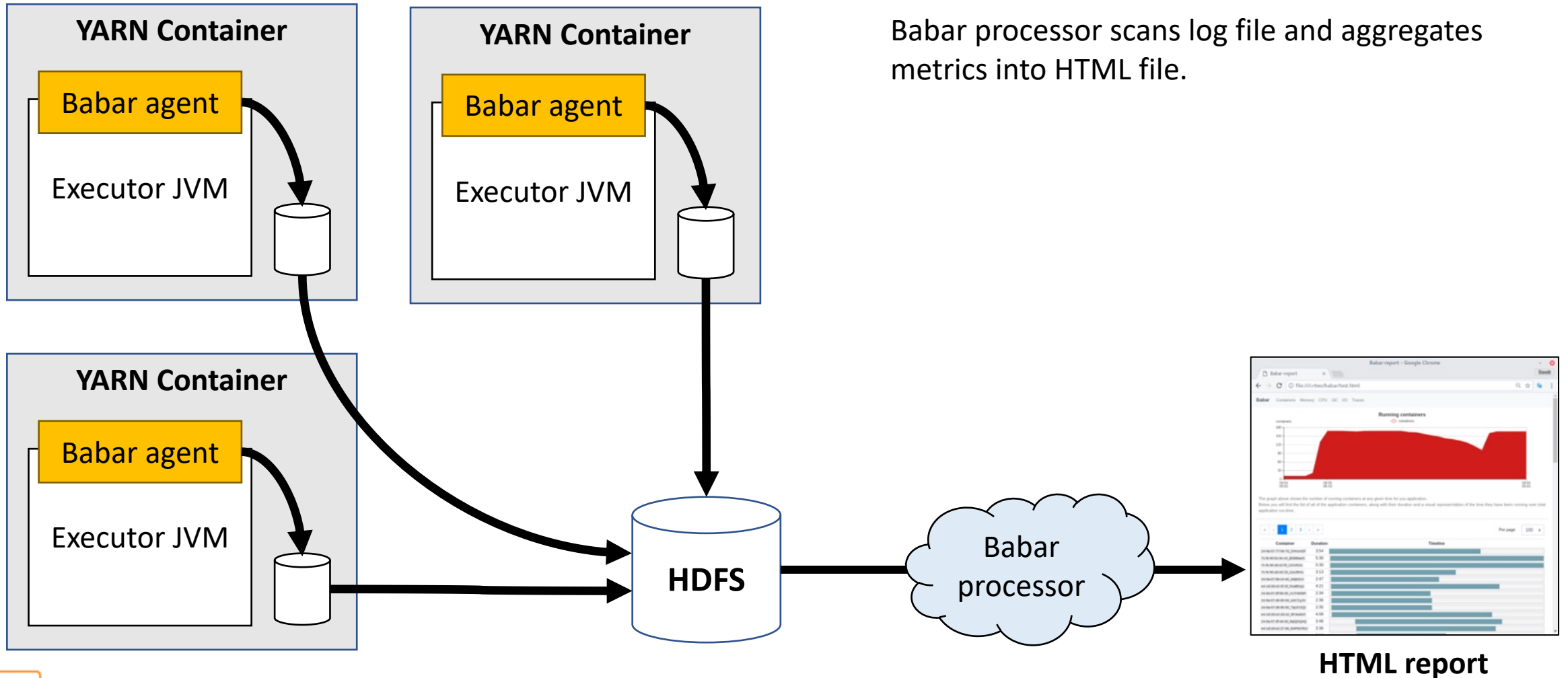
- JVM
- Stack traces
- /proc/

Babar



On job complete, YARN aggregates logs from all executors into single file on HDFS.

Babar



Running containers

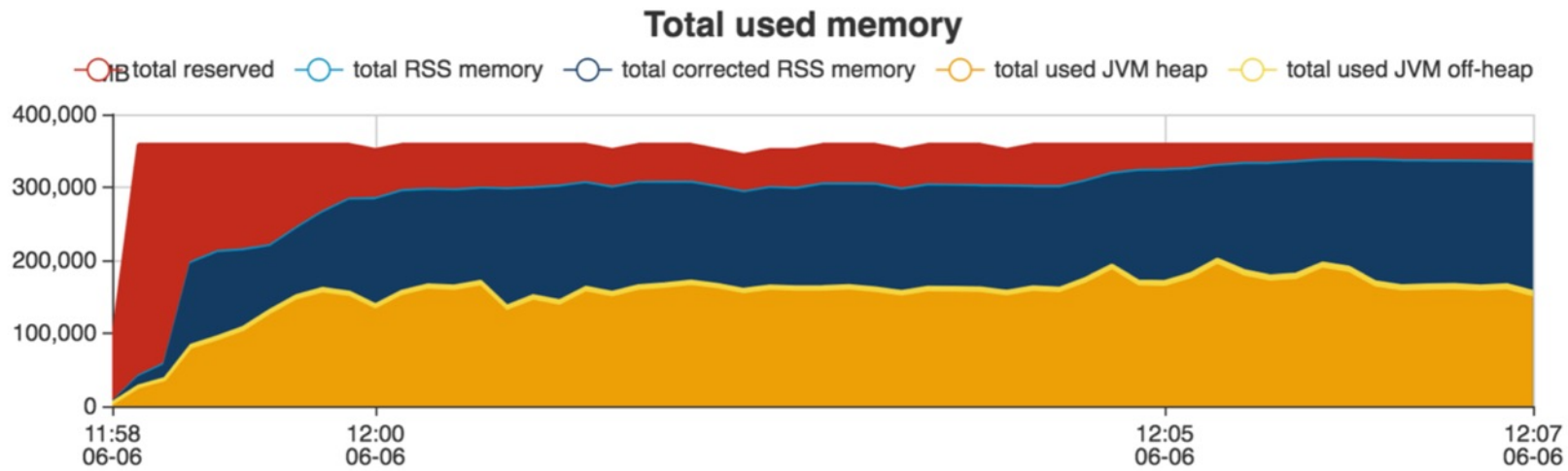


The graph above shows the number of running containers at any given time for your application. Below you will find the list of all of the application containers, along with their duration and a visual representation of the time they have been running over total application run-time.

Container	Duration	Timeline
24-8a-07-77-67-e0_IN2FxvP	8:53	
24-8a-07-bb-b6-90_mRGM3QrD	8:52	
24-8a-07-bb-e2-90_D9jqc9xa	8:52	
24-8a-07-bb-e2-00_DoyBHK3q	8:51	

Total memory shows the sum of the memory usage over all containers
Max memory shows the maximum memory usage for any container

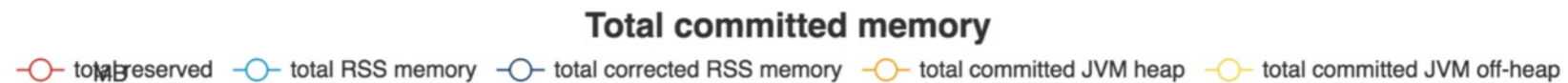
Total memory Max memory



This graph shows the total memory used by all of your application's running containers at any given time.

Reserved memory is the amount of memory reserved on the infrastructure.

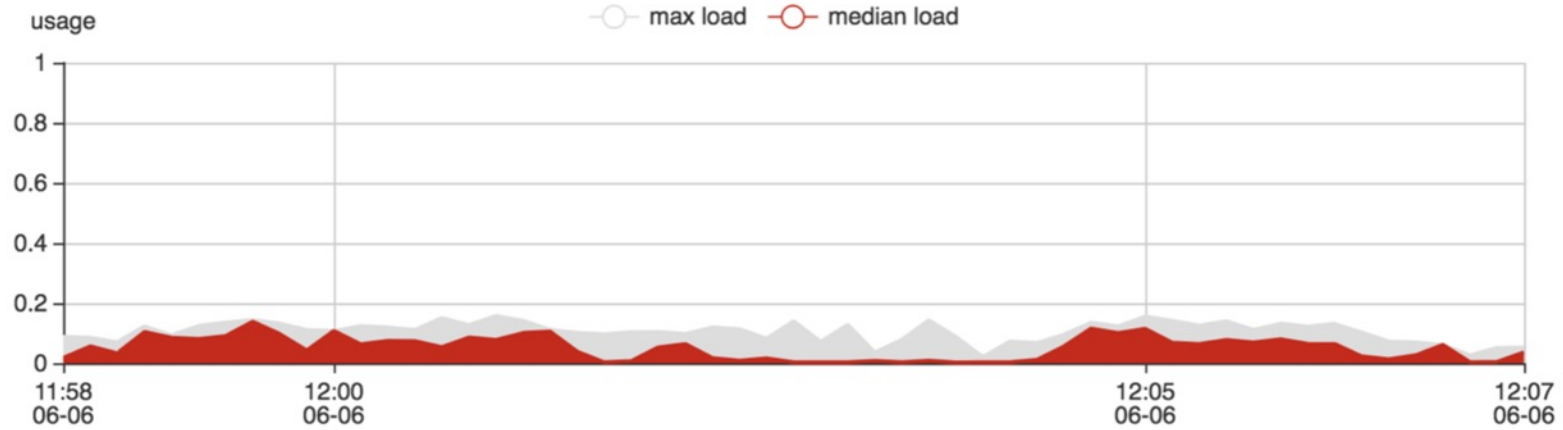
RSS memory is the size of the memory pages that are loaded in the physical memory (RAM) for your containers process-tree (including non JVM programs).



JVM shows the memory usage as reported by the JVM instrumentation using the *JVMProfiler*.
ProcFS shows the memory usage as reported by the `/proc` filesystem using the *ProcFSProfiler*.

JVM ProcFS

JVM CPU usage

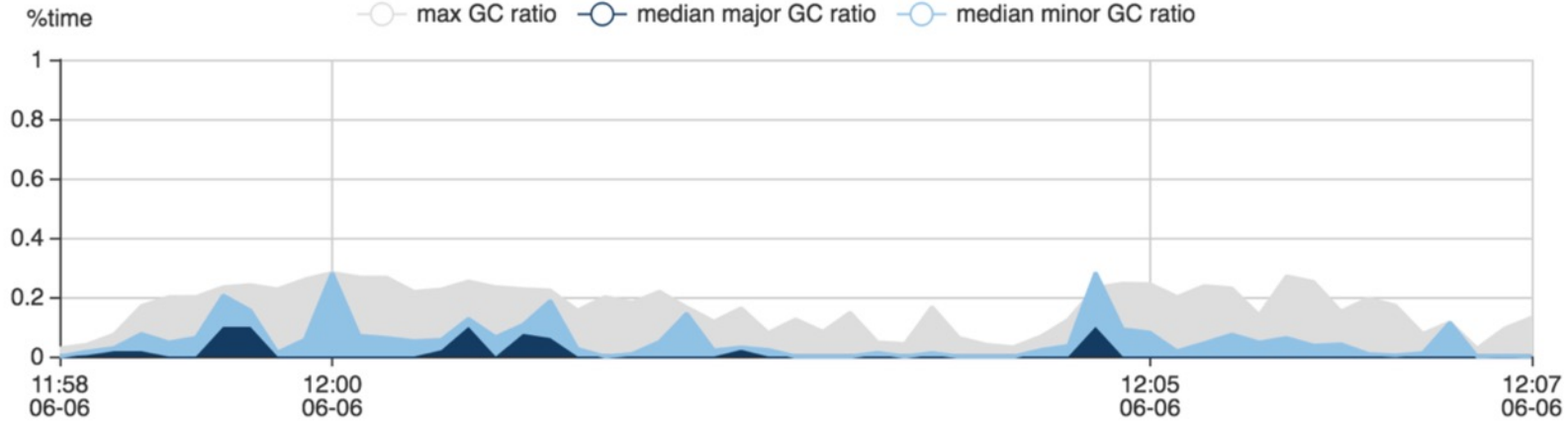


This graph shows the median and max **CPU usage as reported by the JVM instrumentation** over all the containers. It only reports the JVM CPU usage and will ignore children processes spawned by the java application.

Host CPU usage



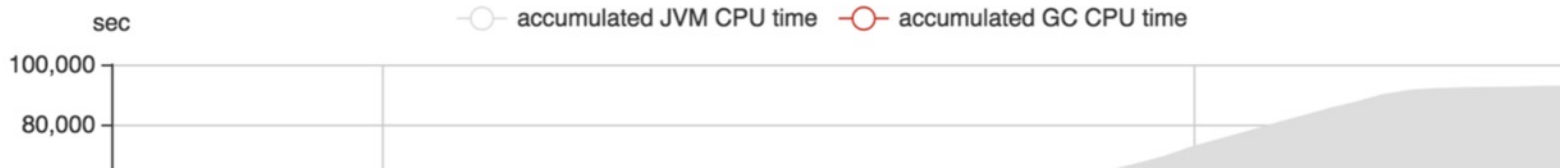
Minor & Major median GC ratio



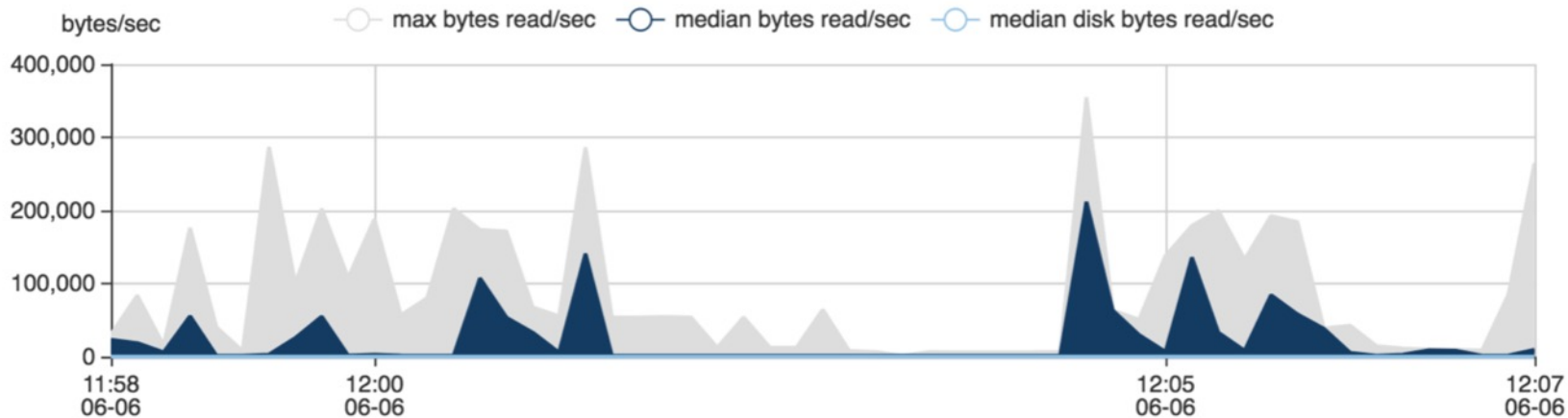
This graph shows the median ratio of wall-clock time spent doing minor and major garbage collections in the JVMs over all containers. **Minor GC** only clean the young generation, while **Major GC** cleans both the young and old ones. The major GC should be much less frequent that the minor one, otherwise it could indicate that too many short lived-objects are promototed to the old generation. If this is the case, you may want to resize the generations and make sure that no humongous object uses most of the old generation (which could trigger frequent major GC).

You can tune the size of the generations either by specifying the `-XX:NewRatio` (integer value only) parameter or with the `-XX:NewSize` and `-XX:MaxNewSize` parameters for finer-grained tuning.

Accumulated JVM CPU time and GC CPU time



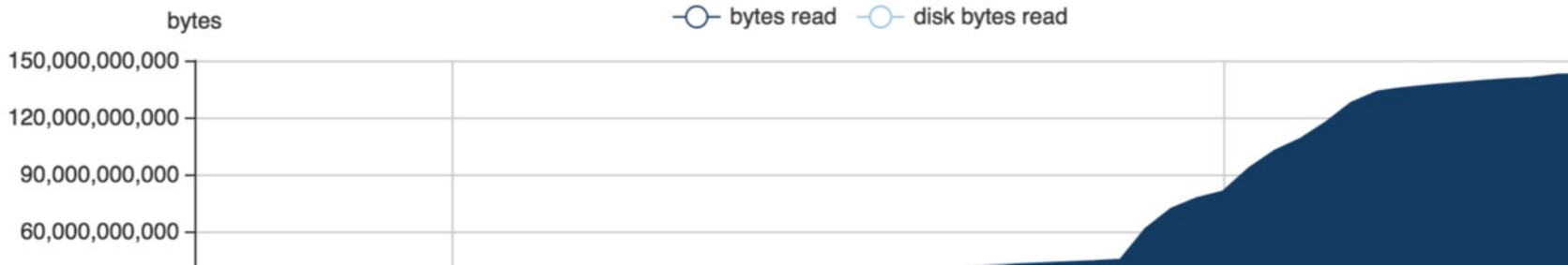
Process tree peak bytes read / sec



This graph shows the peak bandwidth used to read data from disk I/O by the entire process-tree as reported by the `/prod/[pid]/io` file.

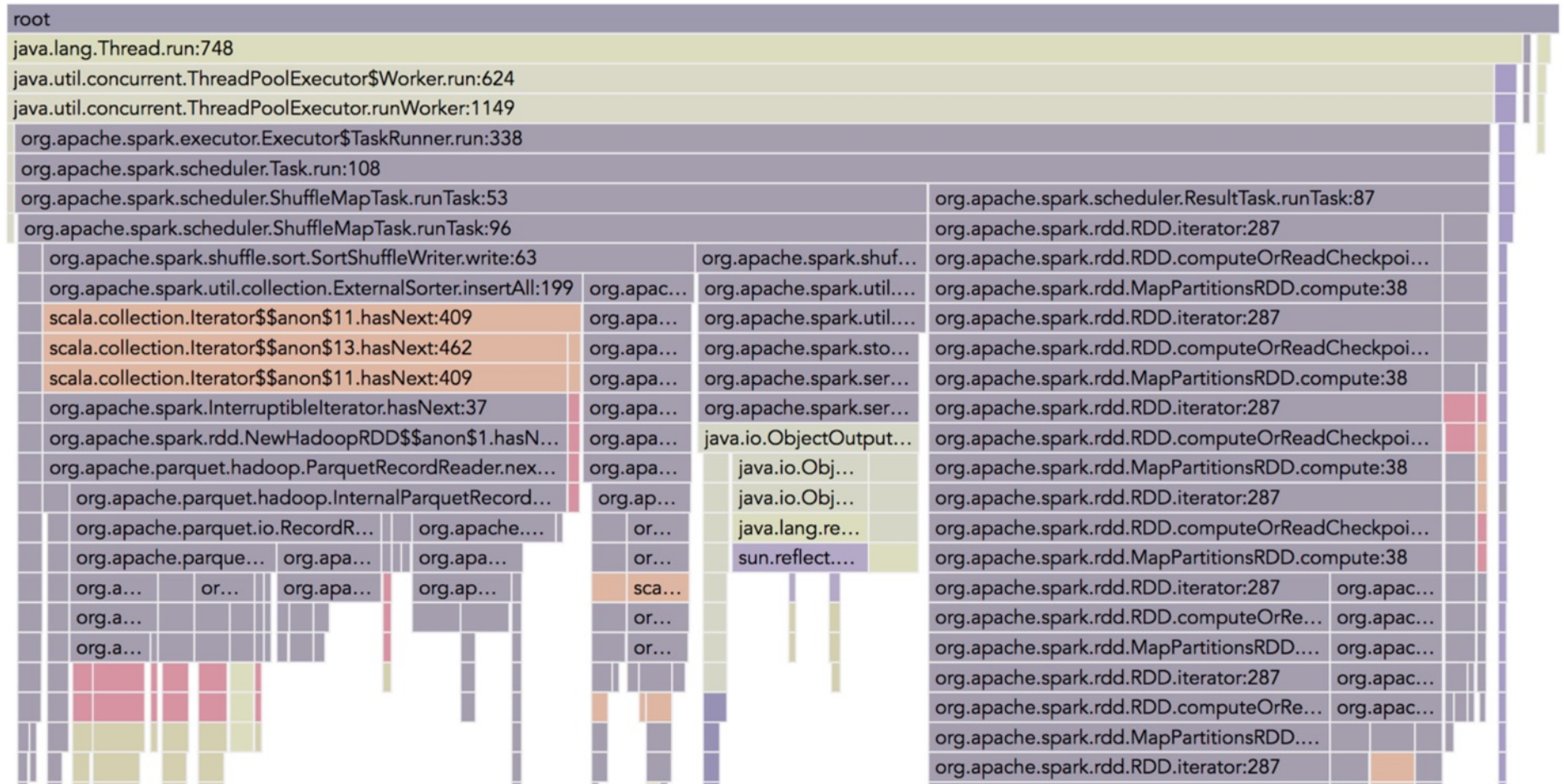
Bytes read/sec show the bandwidth for all read operations, while **disk bytes read/sec** only shows the bandwidth for I/O operations that actually read data from the block storage layer (i.e. excluding pagecache).

Process tree Accumulated bytes read



Search for method-prefixes (commas-separated)

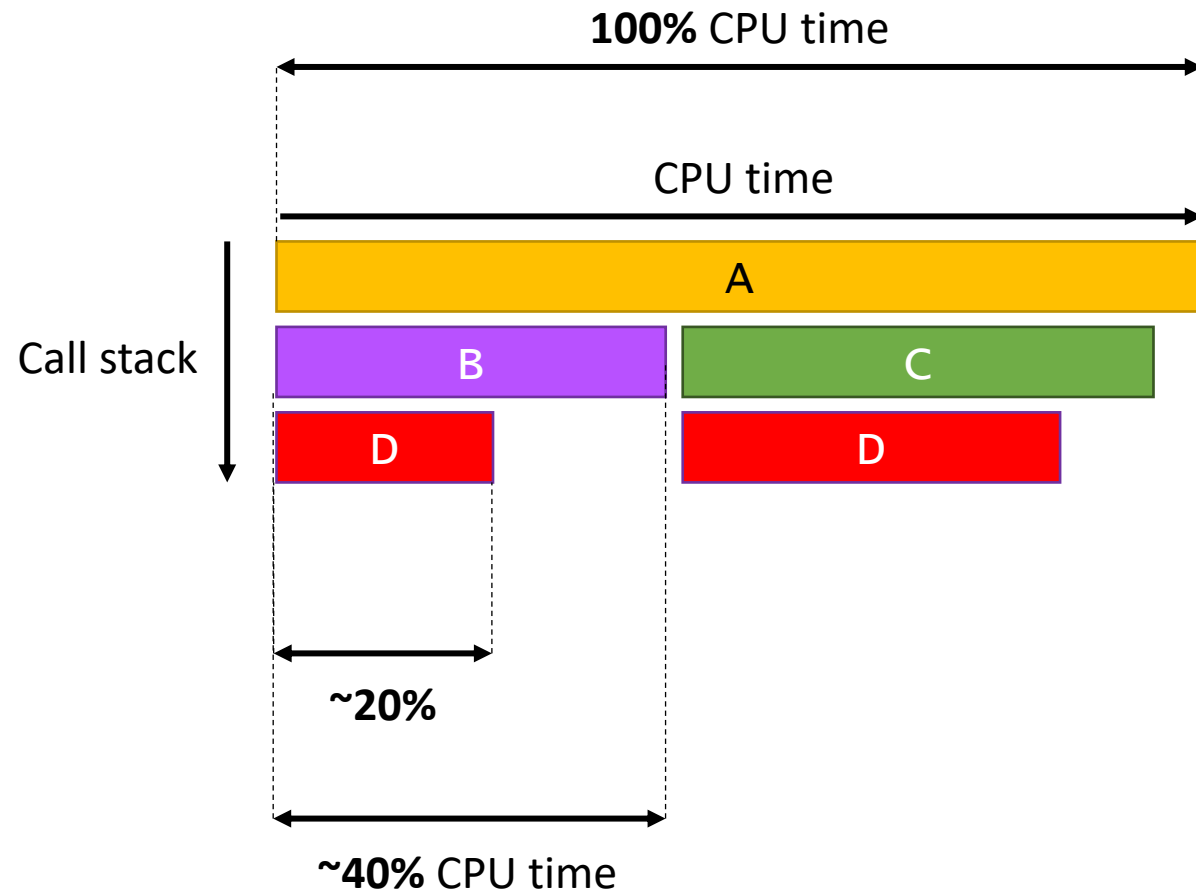
eg: org.project.MyClass.myMethod

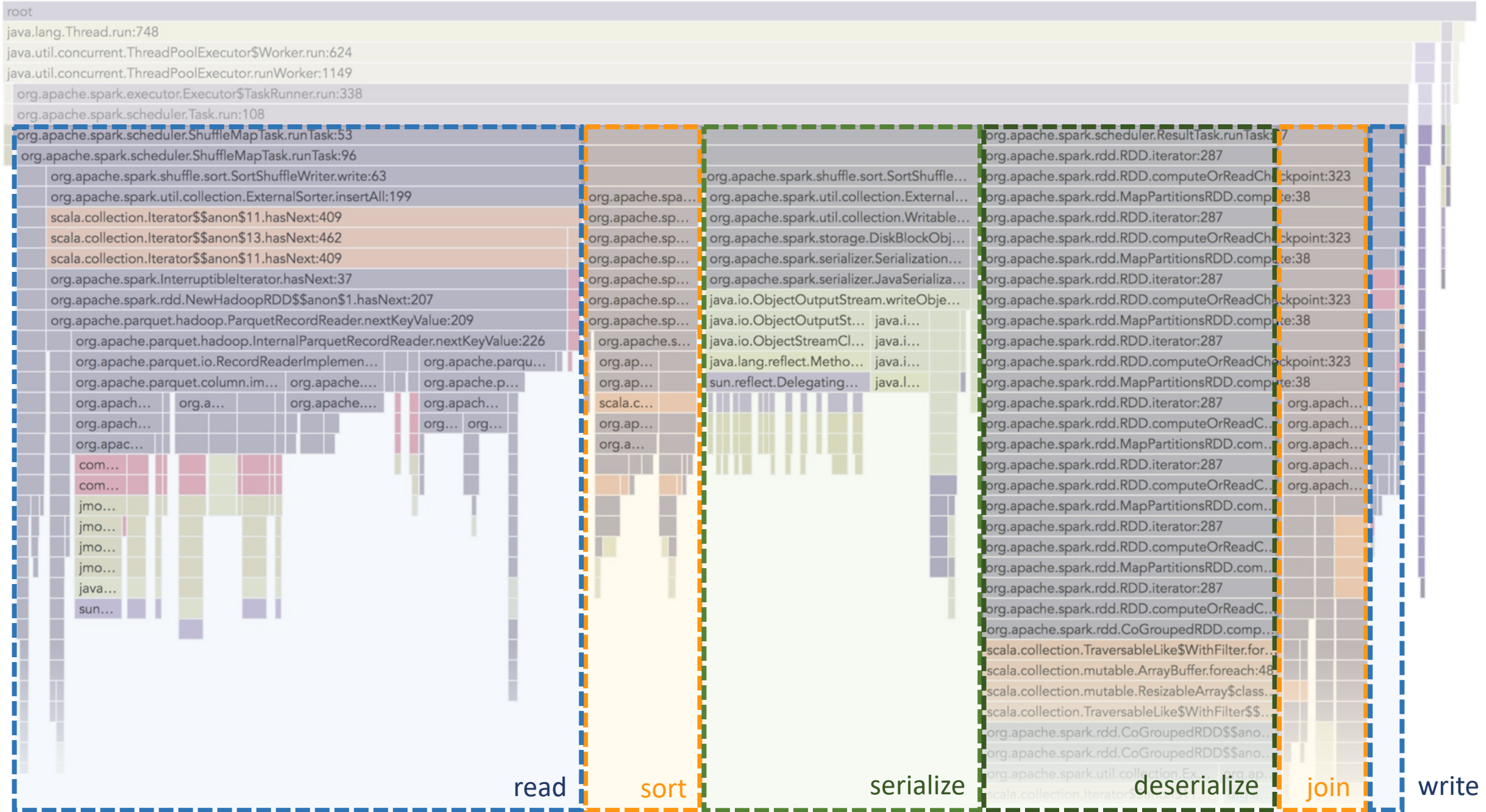


In practice: Flame Graphs & CPU time

Visualize expensive code paths from sampled stack traces

```
A() {  
  B() {  
    D() {  
    }  
  }  
  C() {  
    D() {  
    }  
  }  
}
```

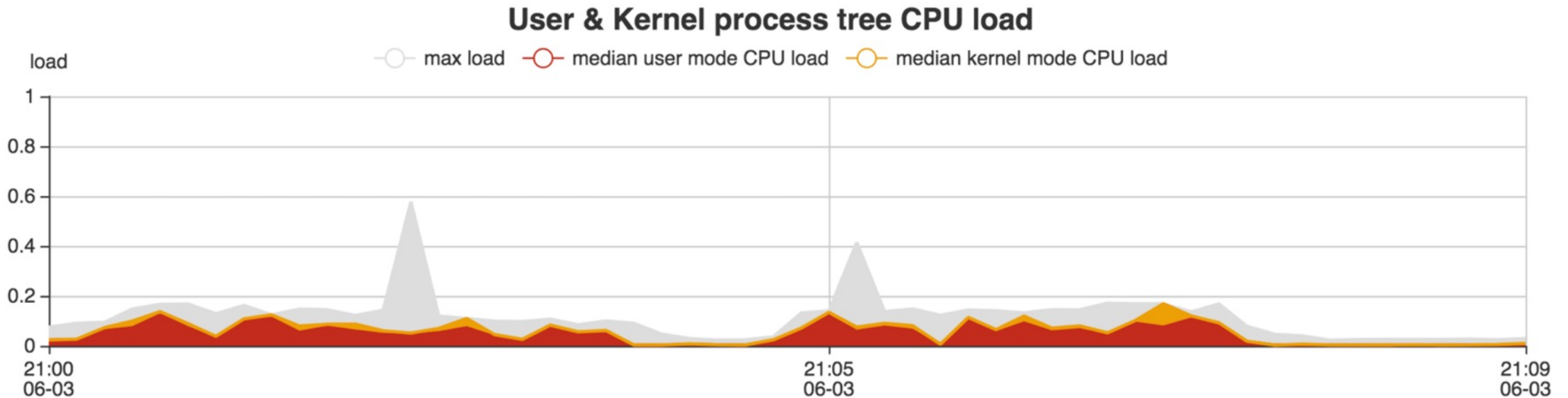




In practice: Flame Graphs & CPU time

Shuffle is expensive

Not network, not IO, but serialization!



In practice: Flame Graphs & CPU time

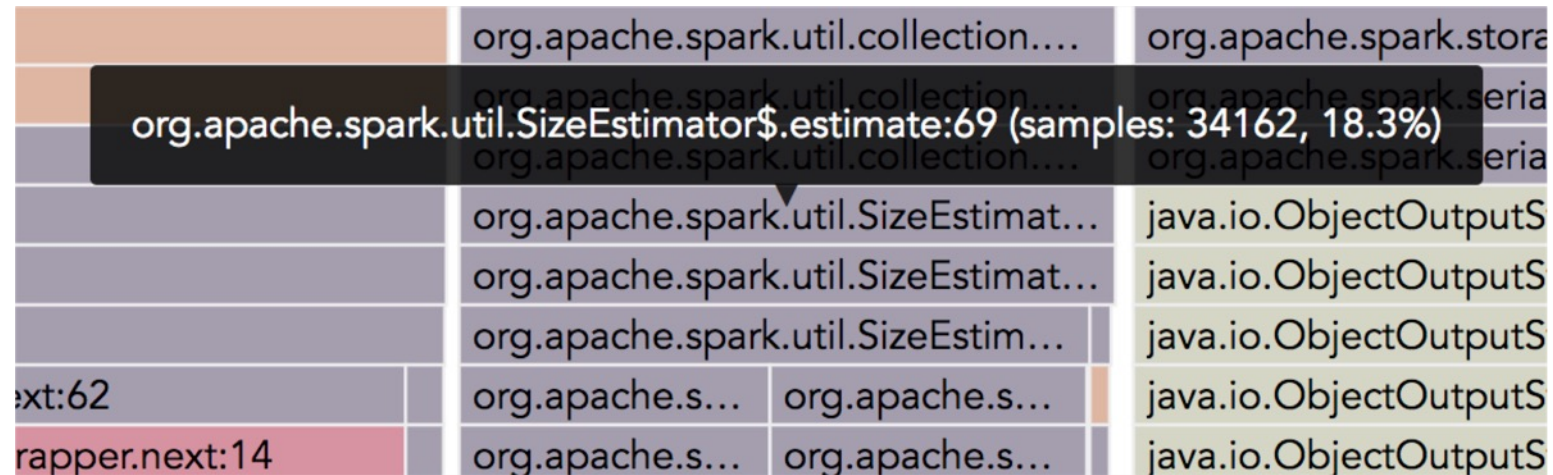
Pick models accordingly

- Serialization/deserialization efficiency often more important than memory footprint
- **Use specialized serializers** (can allow further optimizations by Spark)

In this example join: **-40% CPU time!**

Models also impact cost of **size estimation**

Can be very expensive (we have seen up to 30% CPU time)



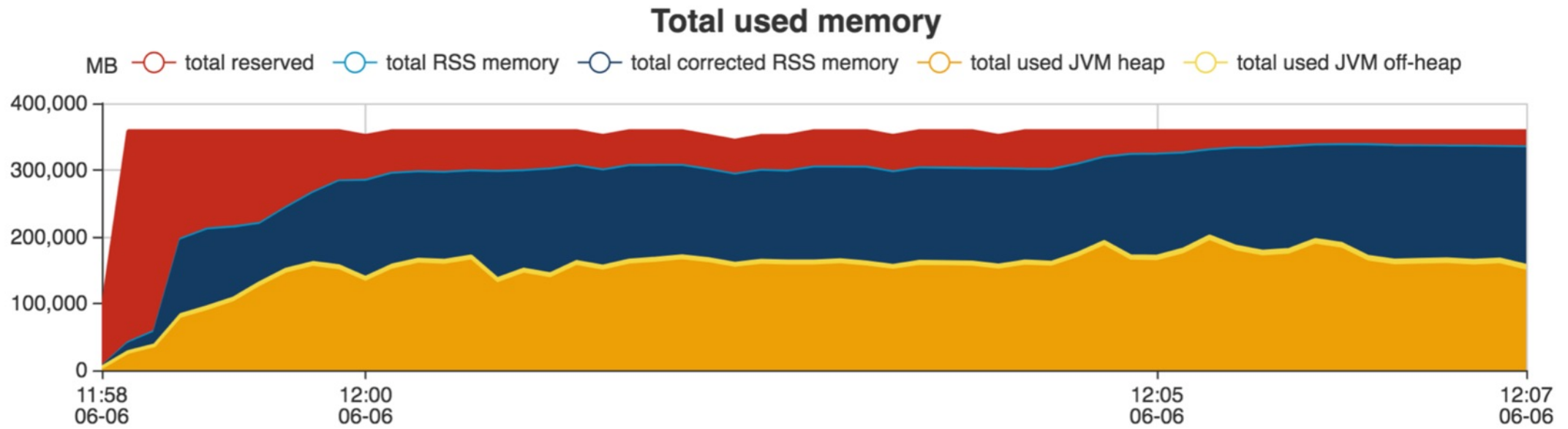
In practice: Memory

Understanding the memory usage

- JVM heap & off-heap memory used

- Physical memory used (RSS memory)

- Reserved memory

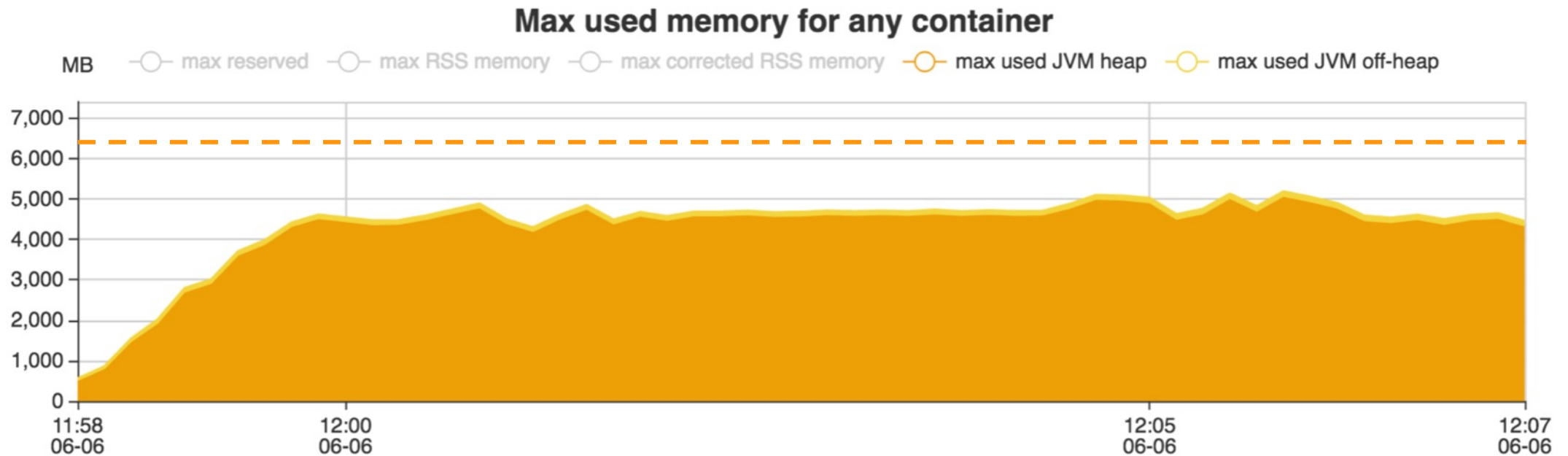


In practice: Memory

Dimensioning and tuning memory

Setting executor & container memory to accommodate the largest executor:

- Avoiding OutOfMemoryErrors (`spark.executor.memory`)

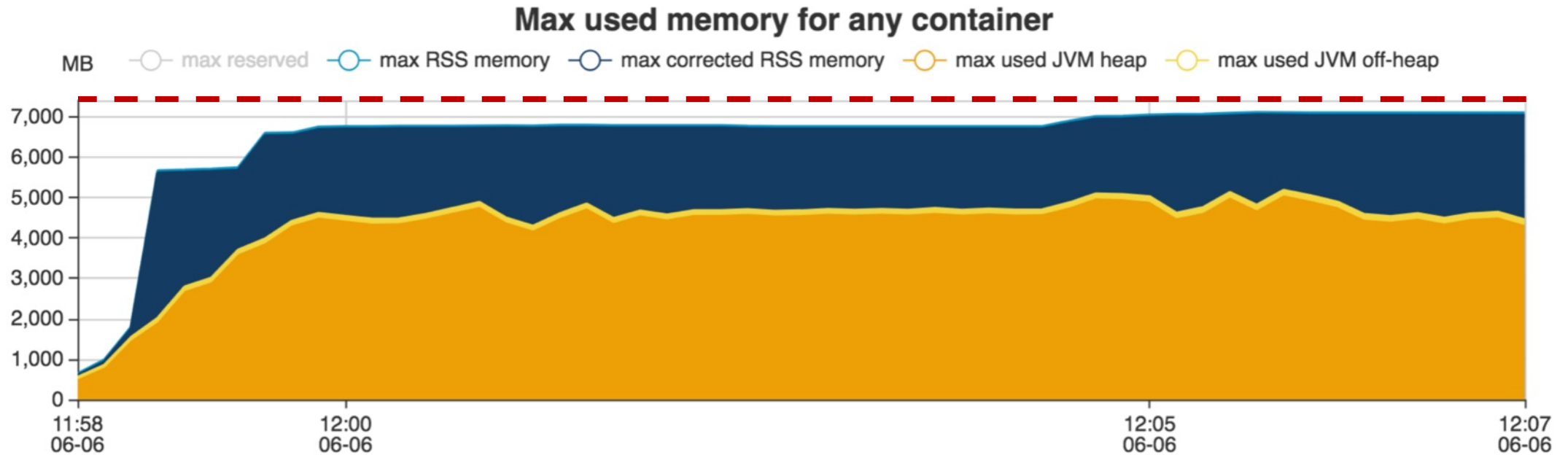


In practice: Memory

Dimensioning and tuning memory

Setting executor & container memory to accommodate the largest executor:

- Avoiding OutOfMemoryErrors (spark.executor.memory)

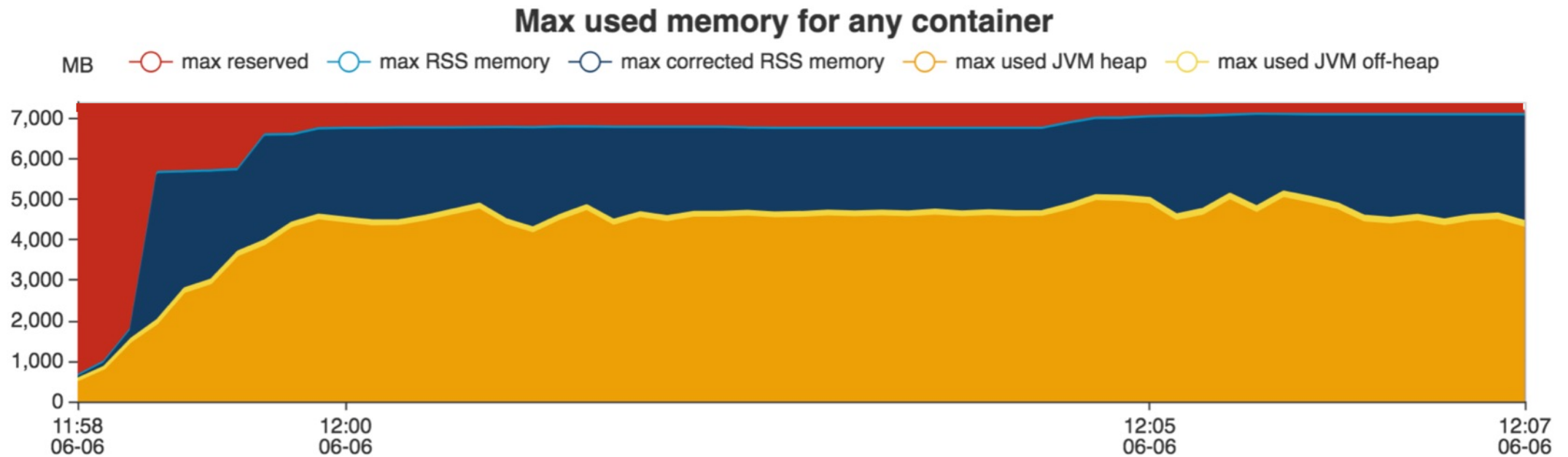


In practice: Memory

Dimensioning and tuning memory

Setting executor & container memory to accommodate the largest executor:

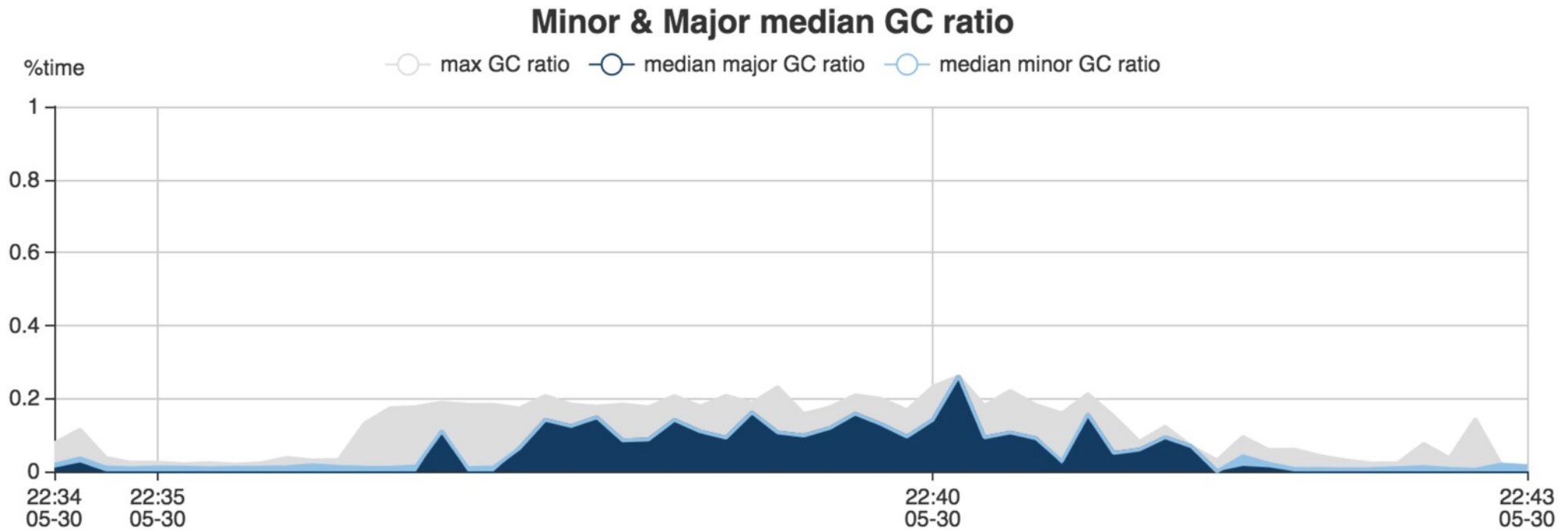
- Avoiding OutOfMemoryErrors (spark.executor.memory)
- Avoiding YARN killing containers for exceeding reserved memory (spark.yarn.executor.memoryOverhead)



In practice: Memory

Dimensioning and tuning memory

Keeping Garbage Collection under control



Recap

Babar helped us a lot, we hope it can help you too!



<https://github.com/criteo/babar>

THANK YOU!