

# Profile-guided heterogeneous-aware scheduling for cloud workloads

Daniel Araújo de Medeiros

<sup>1</sup>Department of Computer Science  
Federal University of Bahia  
Salvador, Brazil

Status on  
September 2018

# Table of Contents

- 1 Hurry-up Proposal
  - Algorithm
  - Instrumentation
- 2 Conclusions

# Premises

The proposed scheme, Hurry-Up, is based on the following premises:

- **Heavy vs Light requests:** There is significant difference in computing demands in the user requests (heavy and light ones).
- **Hot functions:** The application has functions showing high computing demands (hot functions), accounting for the increase of the tail latency.

# Empirical Observations

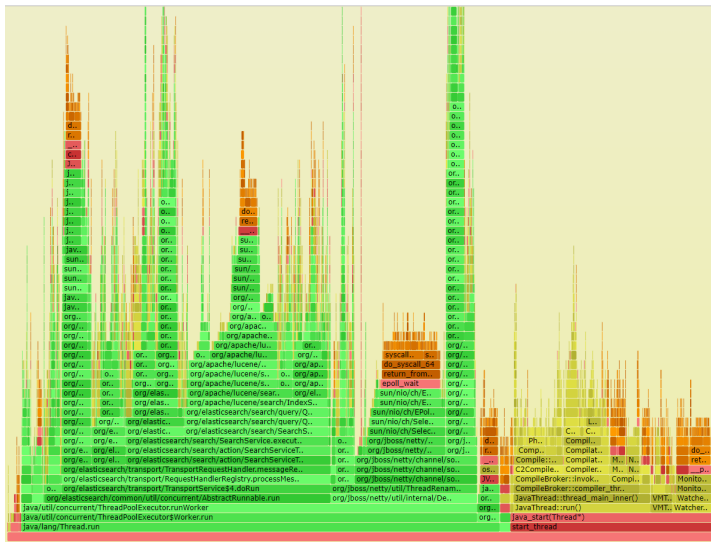
Empirical observations noticed in the algorithm design:

- **Threshold:** To minimize unnecessary up and down migrations, a threshold should be defined when monitoring hot function execution. This is used to better distinguish between heavy and light loads.
- **Priority:** When multiple requests compete for a limited number of cores, there's a need to manage in the task mapping decisions.

## Profiling phase - How to identify a hot function?

- Hot functions are determined empirically via standard software profiling.
- The chosen function should not have a very large number of calls because it may lead to many unnecessary migrations.
- The chosen function should be related to the tail latency experienced by the application because the algorithm is designed to improve the service time and efficiency.

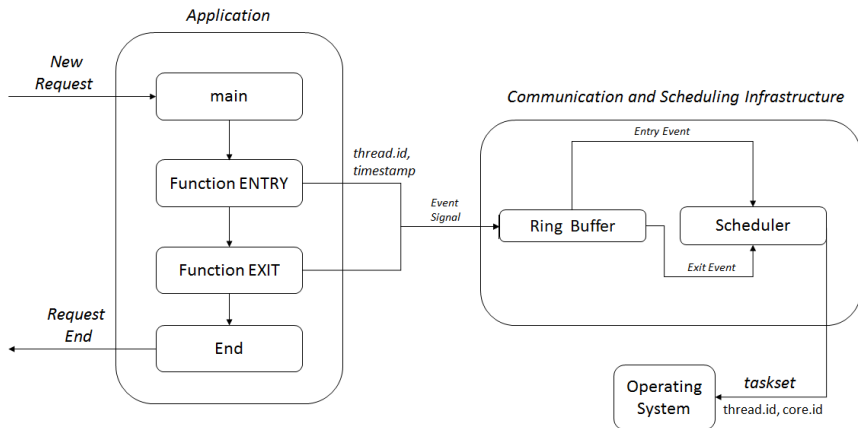
# Elasticsearch's Flame Graph (perf tool)



# Elasticsearch's call tree (YourKit profiler)

Call Tree			Time (ms)	Count
<All threads>			10,780,512	100%
java.lang.Thread.run()			9,575,828	89%
ThreadPoolExecutor.java:1149	org.elasticsearch.common.util.concurrent.AbstractRunnable.run()		8,362,682	78%
AbstractRunnable.java:37	org.elasticsearch.transport.TransportService\$4.doRun()		8,360,366	78%
TransportService.java:378	org.elasticsearch.transport.RequestHandlerRegistry.processMessageReceived(TransportRequest, TransportChannel)		8,360,366	78%
RequestHandlerRegistry.java:77	org.elasticsearch.transport.TransportRequestHandler.messageReceived(TransportRequest, TransportChannel, TransportRequestHandler)		8,360,273	78%
TransportRequestHandler.java:33	org.elasticsearch.search.action.SearchServiceTransportAction\$SearchQueryTransportHandler.messageReceived(TransportRequest, TransportChannel)		8,296,286	77%
SearchServiceTransportAction.java:365	org.elasticsearch.search.action.SearchServiceTransportAction\$SearchQueryTransportHandler.messageReceived(TransportRequest, TransportChannel)		8,296,286	77%
SearchServiceTransportAction.java:368	org.elasticsearch.search.SearchService.executeQueryPhase(ShardSearchRequest)		8,295,466	77%
SearchService.java:385	org.elasticsearch.search.SearchService.loadOrExecuteQueryPhase(ShardSearchRequest, SearchContext, QueryPhase)		8,292,433	77%
SearchService.java:372	org.elasticsearch.search.query.QueryPhase.execute(SearchContext)		8,292,421	77%
QueryPhase.java:113	org.elasticsearch.search.query.QueryPhase.execute(SearchContext, IndexSearcher)		8,292,415	77%
QueryPhase.java:384	org.apache.lucene.search.IndexSearcher.search(Query, Collector)		8,292,168	77%
IndexSearcher.java:535	org.apache.lucene.search.IndexSearcher.search(List, Weight, Collector)		8,228,665	76%
IndexSearcher.java:821	org.apache.lucene.search.BulkScorer.score(LeafCollector, Bits)		8,209,549	76%
BulkScorer.java:39	org.apache.lucene.search.BooleanScorer.score(LeafCollector, Bits, int, int)		8,106,592	75%
BooleanScorer.java:335	org.apache.lucene.search.BooleanScorer.scoreWindow(BooleanScorer\$BulkScorerAndDoc, LeafCollector, Bits, int, int)		8,095,698	75%
BooleanScorer.java:311	org.apache.lucene.search.BooleanScorer.scoreWindowMultipleScorers(LeafCollector, Bits, int, int)		8,057,975	75%
BooleanScorer.java:266	org.apache.lucene.search.BooleanScorer.scoreWindowIntoBitSetAndReplay(LeafCollector, Bits, int, int)		8,043,423	75%
BooleanScorer.java:219	org.apache.lucene.search.BooleanScorer\$BulkScorerAndDoc.score(LeafCollector, Bits, int, int)		7,012,186	65%
BooleanScorer.java:61	org.apache.lucene.search.BooleanScorer.score(LeafCollector, Bits, int, int)		4,980,101	46%
BooleanScorer.java:335	org.apache.lucene.search.BooleanScorer.scoreWindowMultipleScorers(LeafCollector, Bits, int, int)		4,978,243	46%
BooleanScorer.java:311	org.apache.lucene.search.BooleanScorer.scoreWindowMultipleScorers(LeafCollector, Bits, int, int)		4,974,738	46%
BooleanScorer.java:266	org.apache.lucene.search.BooleanScorer.scoreWindowIntoBitSetAndReplay(LeafCollector, Bits, int, int)		4,970,490	46%
BooleanScorer.java:219	org.apache.lucene.search.BooleanScorer\$BulkScorerAndDoc.score(LeafCollector, Bits, int, int)		4,065,679	38%
BooleanScorer.java:61	org.apache.lucene.search.BooleanScorer.score(LeafCollector, Bits, int, int)		4,064,482	38%
Weight.java:183	org.apache.lucene.search.Weight\$DefaultBulkScorer.score(LeafCollector, Bits, int, int)		4,061,427	38%
Weight.java:196	org.apache.lucene.search.BooleanScorer\$OrCollector.score(LeafCollector, Bits, int, int)		3,125,437	29%
BooleanScorer.java:143	org.apache.lucene.search.TermScorer.score(LeafCollector, Bits, int, int)		2,561,114	24%
TermScorer.java:66	org.apache.lucene.search.similarities.TFIDFSimilarity.score(LeafCollector, Bits, int, int)		2,010,255	19%
TFIDFSimilarity.java:715	org.apache.lucene.codecs.lucene53.Lucene53NormsProducer.score(LeafCollector, Bits, int, int)		879,723	8%
Lucene53NormsProducer.java:139	org.apache.lucene.store.ByteBufferStore.score(LeafCollector, Bits, int, int)		159,392	1%
Lucene53NormsProducer.java:139	org.apache.lucene.store.ByteBufferStore.score(LeafCollector, Bits, int, int)		150,249	1%
TFIDFSimilarity.java:713	org.apache.lucene.search.similarities.ClassicSimilarity.score(LeafCollector, Bits, int, int)		277,372	3%
Weight.java:198	org.apache.lucene.codecs.lucene50.Lucene50PostingList.score(LeafCollector, Bits, int, int)		360,752	3%
Weight.java:171	org.apache.lucene.search.TermScorer.docID()		567	0%
Weight.java:175	org.apache.lucene.search.TermScorer.docID()		463	0%
BooleanScorer.java:222	org.apache.lucene.search.BooleanScorer.scoreMatch(LeafCollector, Bits, int, int)		903,610	8%
BooleanScorer.java:223	org.apache.lucene.util.PriorityQueue.fill(long[], long)		≥ 2	≥ 3,849
BooleanScorer.java:271	org.apache.lucene.util.PriorityQueue.insertWithOverflow(long, long)		2,712	0%
BooleanScorer.java:264	org.apache.lucene.util.PriorityQueue.clear()		333	0%

# Overview of the Scheduler Design





# Illustration of the algorithm design

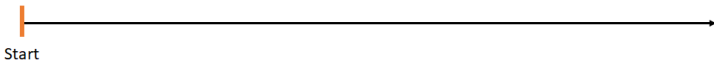
Number of physical cores;  
In case of Juno, 6 cores.

Long integer type

0: Not in Function  
1: Executing hot function  
2: 1 + Over threshold

	<i>core.id</i>	<i>thread.id</i>	<i>timestamp</i>	<i>thread priority</i>
Big cores	0			
	1			
Little cores	2			
	3			
	4			
	5			

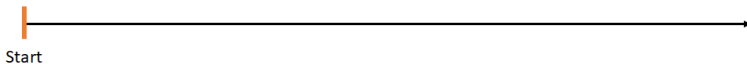
## EVENT TIMELINE



# Case 1: No Competition and Below Threshold

<i>core.id</i>	<i>thread.id</i>	<i>timestamp</i>	<i>thread priority</i>
0			
1			
2			
3			
4			
5			

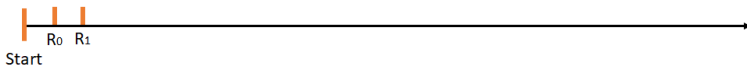
EVENT TIMELINE



# Case 1: No Competition and Below Threshold

<i>core.id</i>	<i>thread.id</i>	<i>timestamp</i>	<i>thread priority</i>
0	900	2	1
1	901	3	1
2			
3			
4			
5			

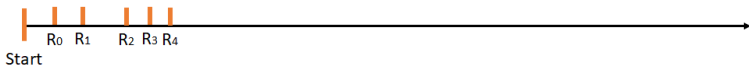
## EVENT TIMELINE



# Case 1: No Competition and Below Threshold

<i>core.id</i>	<i>thread.id</i>	<i>timestamp</i>	<i>thread priority</i>
0	900	2	1
1	901	3	1
2	902	5	1
3	903	6	1
4	904	7	1
5			

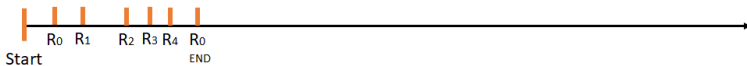
## EVENT TIMELINE



# Case 1: No Competition and Below Threshold

<i>core.id</i>	<i>thread.id</i>	<i>timestamp</i>	<i>thread priority</i>
0	900	2	0
1	901	3	1
2	902	5	1
3	903	6	1
4	904	7	1
5			

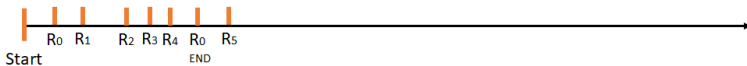
## EVENT TIMELINE



# Case 1: No Competition and Below Threshold

<i>core.id</i>	<i>thread.id</i>	<i>timestamp</i>	<i>thread priority</i>
0	900	2	0
1	901	3	1
2	902	5	1
3	903	6	1
4	904	7	1
5	905	9	1

## EVENT TIMELINE



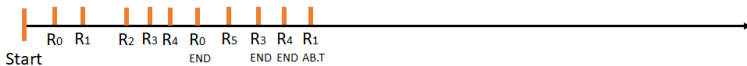
# Notes

- In the initialization, the threads are allocated in a round-robin fashion.
- Up migration to a big core in the initialization is not performed, even when the big core is available, so as to reduce the number of unnecessary migrations.
- If the thread is not above the threshold in a little core (priority = 2 and core id between 2 and 5), a migration is not performed (the request was light and little core was already able to handle it).

## Case 2: No Competition and Above Threshold

<i>core.id</i>	<i>thread.id</i>	<i>timestamp</i>	<i>thread priority</i>
0	900	2	0
1	901	3	2
2	902	5	1
3	903	6	0
4	904	7	0
5	905	9	1

### EVENT TIMELINE

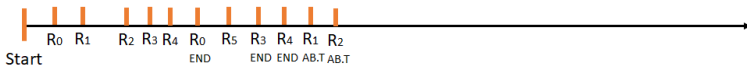




## Case 2: No Competition and Above Threshold

<i>core.id</i>	<i>thread.id</i>	<i>timestamp</i>	<i>thread priority</i>
0	902	2	2
1	901	3	2
2	900	5	0
3	903	6	0
4	904	7	0
5	905	9	1

### EVENT TIMELINE



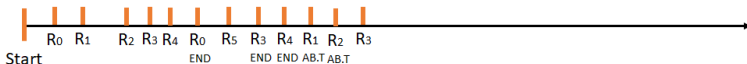
# Notes

- If the request is already on big core and goes to priority 2, we can't do anything about it - already on fastest core.
- In case a heavy request (above threshold) is on a little core and there's a slot available on a big core (priority = 0 or 1), swap it.

# Case 3: Competition and Above Threshold

<i>core.id</i>	<i>thread.id</i>	<i>timestamp</i>	<i>thread priority</i>
0	902	2	2
1	901	3	2
2	900	5	0
3	903	20	1
4	904	7	0
5	905	9	1

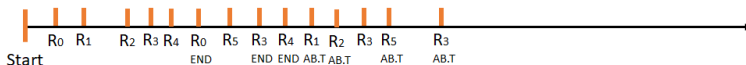
## EVENT TIMELINE



# Case 3: Competition and Above Threshold

<i>core.id</i>	<i>thread.id</i>	<i>timestamp</i>	<i>thread priority</i>
0	902	2	2
1	901	3	2
2	900	5	0
3	903	20	2
4	904	7	0
5	905	9	2

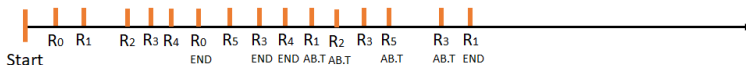
## EVENT TIMELINE



# Case 3: Competition and Above Threshold

<i>core.id</i>	<i>thread.id</i>	<i>timestamp</i>	<i>thread priority</i>
0	902	2	2
1	905	9	2
2	900	5	0
3	903	20	2
4	904	7	0
5	901	3	0

## EVENT TIMELINE



# Case 3: Competition and Above Threshold

<i>core.id</i>	<i>thread.id</i>	<i>timestamp</i>	<i>thread priority</i>
0	903	20	2
1	905	9	2
2	900	5	0
3	902	2	0
4	904	7	0
5	901	3	0

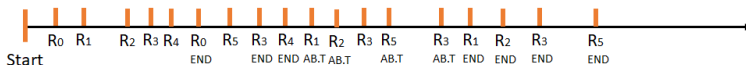
## EVENT TIMELINE



# Case 3: Competition and Above Threshold

<i>core.id</i>	<i>thread.id</i>	<i>timestamp</i>	<i>thread priority</i>
0	903	20	0
1	905	9	0
2	900	5	0
3	902	2	0
4	904	7	0
5	901	3	0

## EVENT TIMELINE

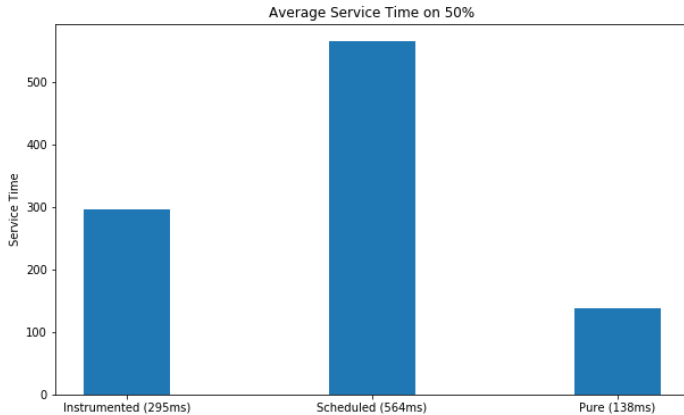


# Notes

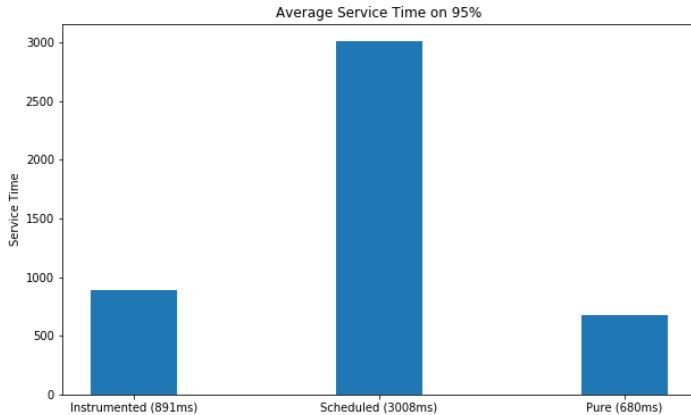
- Scheduler runs every  $t$  milliseconds in order to check for new events and perform thread management.
- Threshold should be determined empirically
- Every time a new request arrives, the timestamp on matrix is updated.
- In case there are more requests with priority = 2 than big cores, the oldest ones has priority of execution (based on the timestamp information)



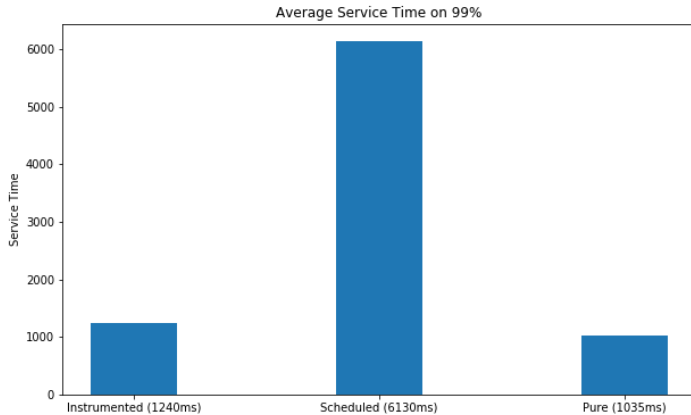
# Results



# Results



# Results



# Future steps

## Challenges ahead:

- Algorithm is designed for the situation where the number of cores is equal or less to number of threads.
- Code instrumentation can lead to excessive overhead.
- Initial results show that further refining is necessary in the algorithm implementation.
- Apply the algorithm to other cloud workloads, such as Cassandra.

# List of Contributors

- Daniel Mossé (Pitt): Advising Committee;
- Denilson Amorim (UFBA): Instrumentation infrastructure;
- George Lima (UFBA): Advising Committee;
- Paul Carpenter (BSC): Constructive feedback and access to the Juno board at BSC;
- Rajiv Nishtala (NTNU): Constructive feedback and help with experiments on Juno;
- Vinicius Petrucci (UFBA): Advising Committee, Chair.

Thank you!