



(1)



(2)



(3)

Runtime Object Lifetime Profiler for Latency Sensitive Big Data Applications

Rodrigo Bruno^{(1)*}, Duarte Patrício⁽²⁾, José Simão⁽²⁾, Luís Veiga⁽²⁾, Paulo Ferreira^(2,3)

EuroSys 2019 @ Dresden, 25-28 March

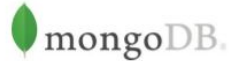
* Work done while at (2)

Big Data Application Stack

Big Data Application



Big Data Platform



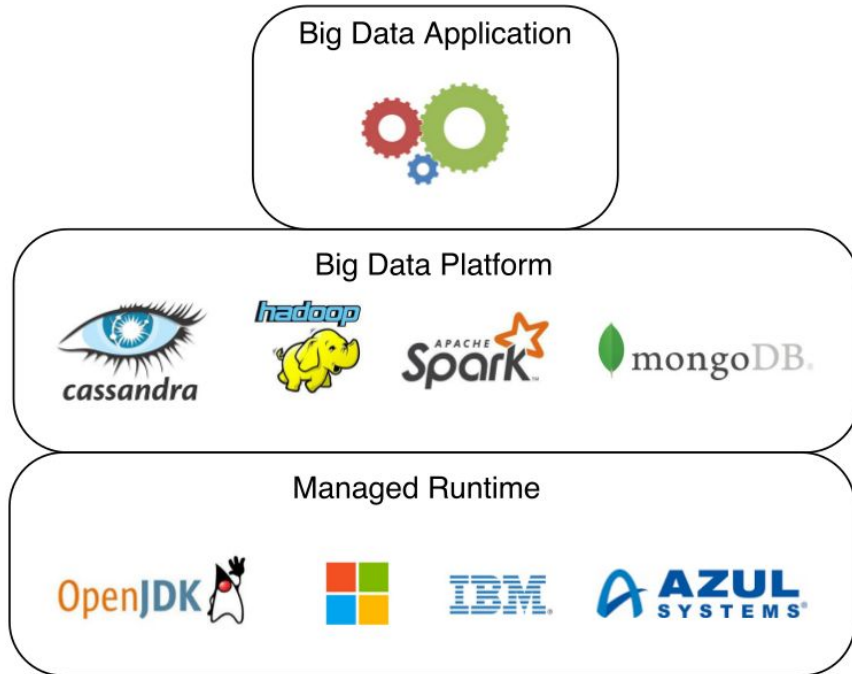
Managed Runtime



Big Data Application Stack

Examples of Latency Sensitive apps

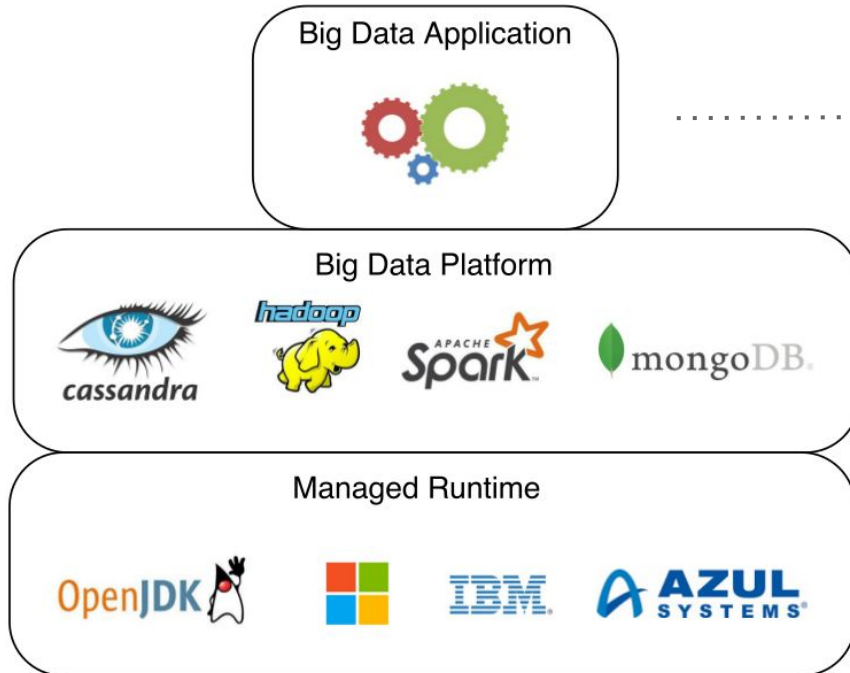
- Banking applications and services
- Context-Aware Ad services
- Games
- ...



Big Data Application Stack

Examples of Latency Sensitive apps

- Banking applications and services
- Context-Aware Ad services
- Games
- ...

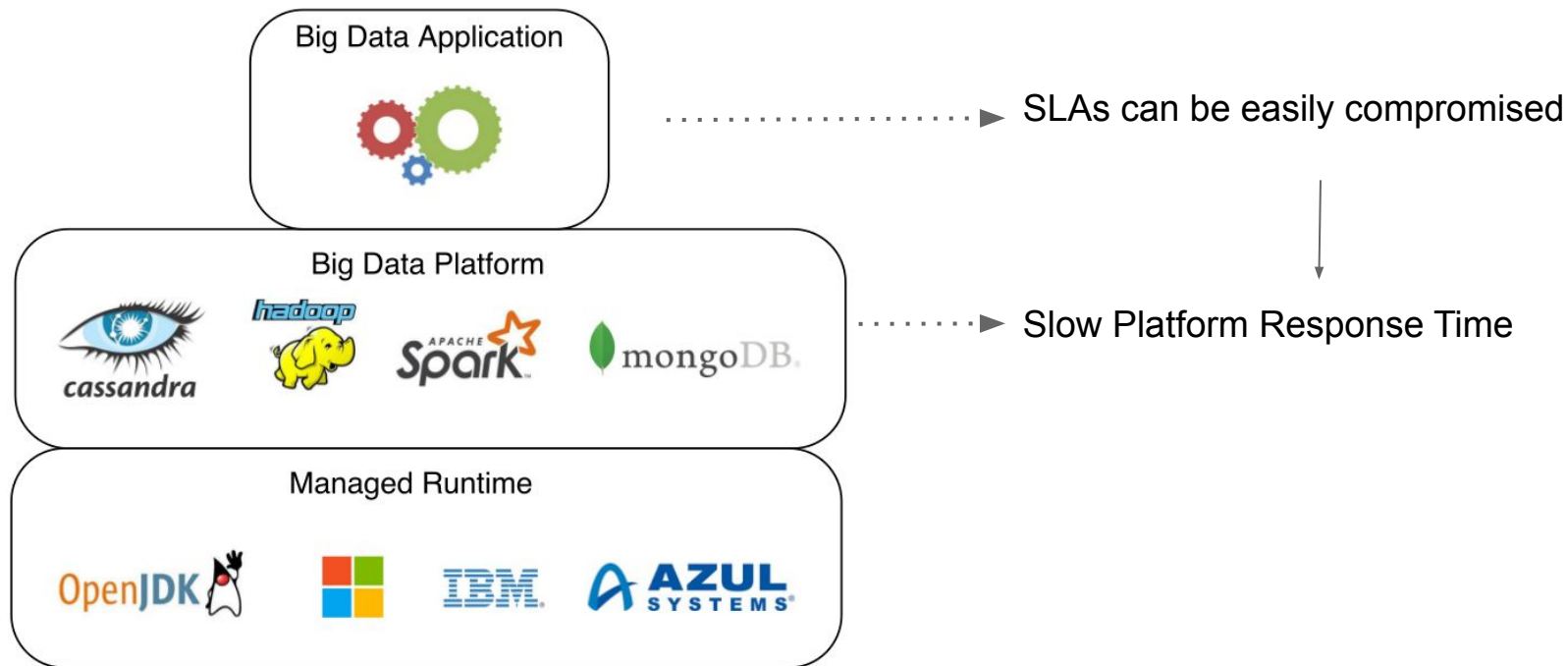


SLAs can be easily compromised

Big Data Application Stack

Examples of Latency Sensitive apps

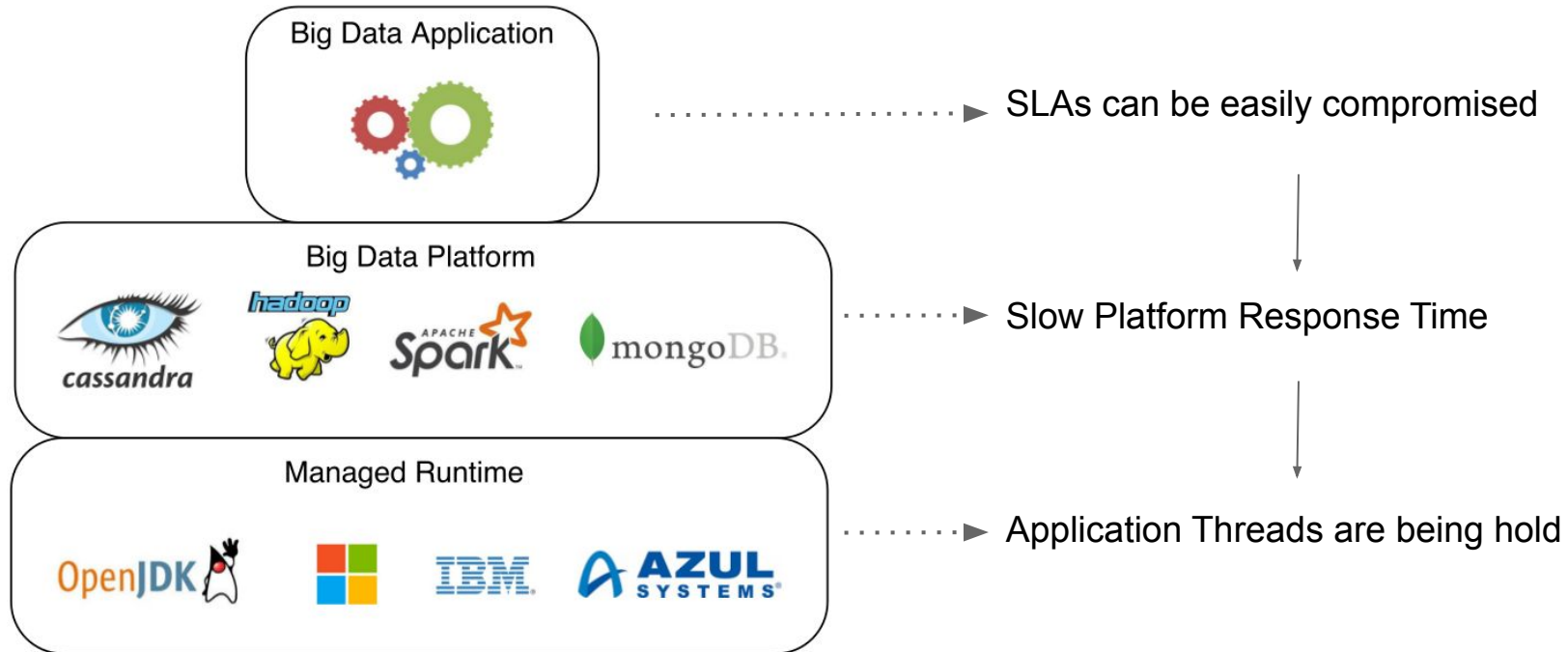
- Banking applications and services
- Context-Aware Ad services
- Games
- ...



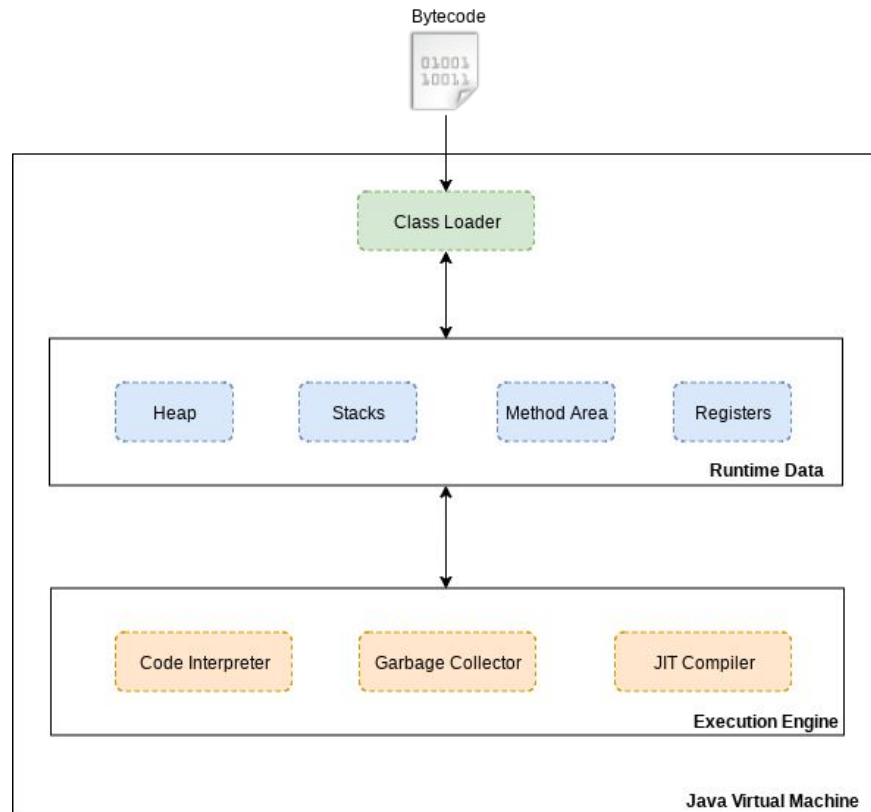
Big Data Application Stack

Examples of Latency Sensitive apps

- Banking applications and services
- Context-Aware Ad services
- Games
- ...



GC-induced Application Latency

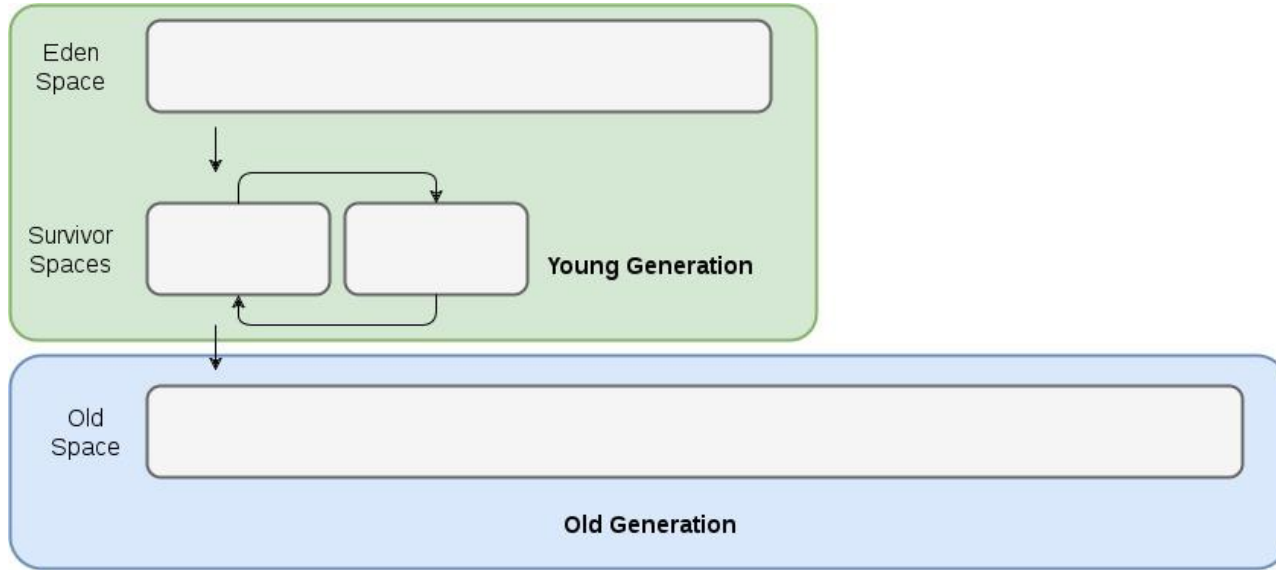


GC is known to have difficulties scaling to high number of cores and memory, mainly w.r.t.

Latencies:

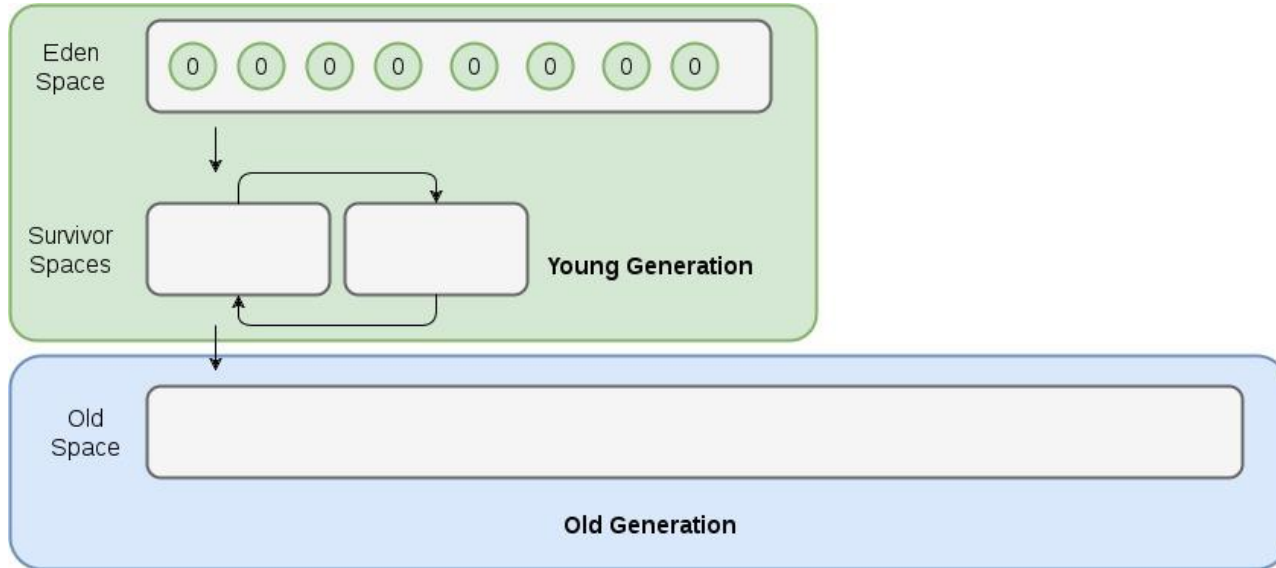
- [ACM CSUR 2018]
- [DSN 2018]
- [ISMM 2017]
- [ISMM 2015]
- [ASPLOS 2013]
- ...

OpenJDK HotSpot Generational GCs (PS, CMS, G1)



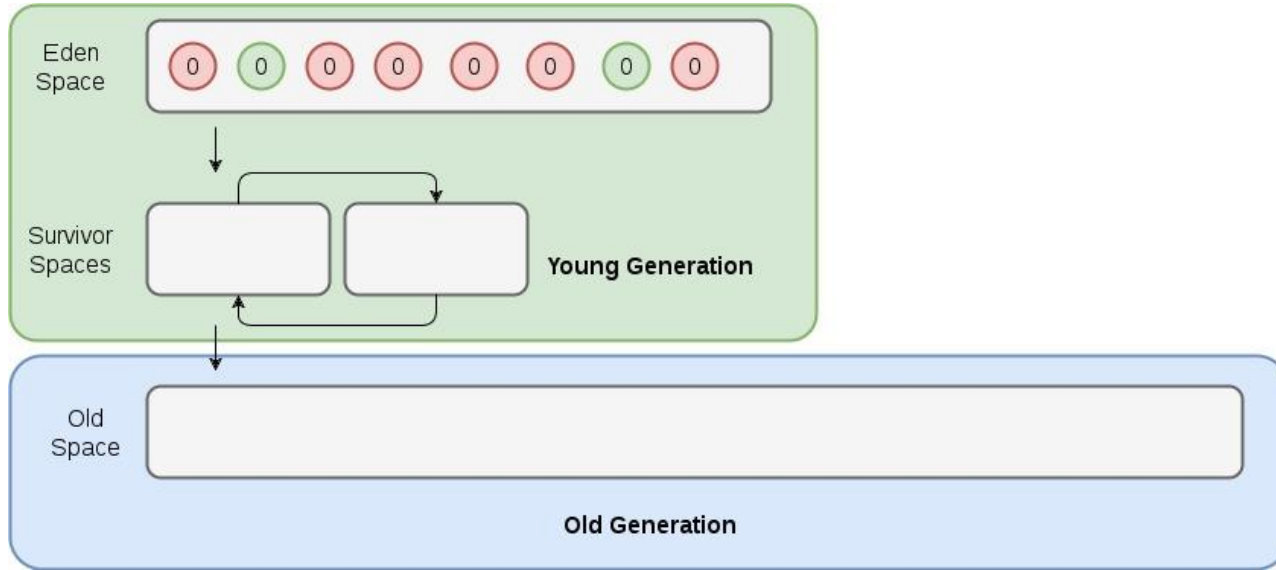
- Two generations:
 - **Young** and **Old**
- Surviving objects are copied to
 - **Survivor** spaces and then to
 - the **Old** generation.

OpenJDK HotSpot Generational GCs (PS, CMS, G1)



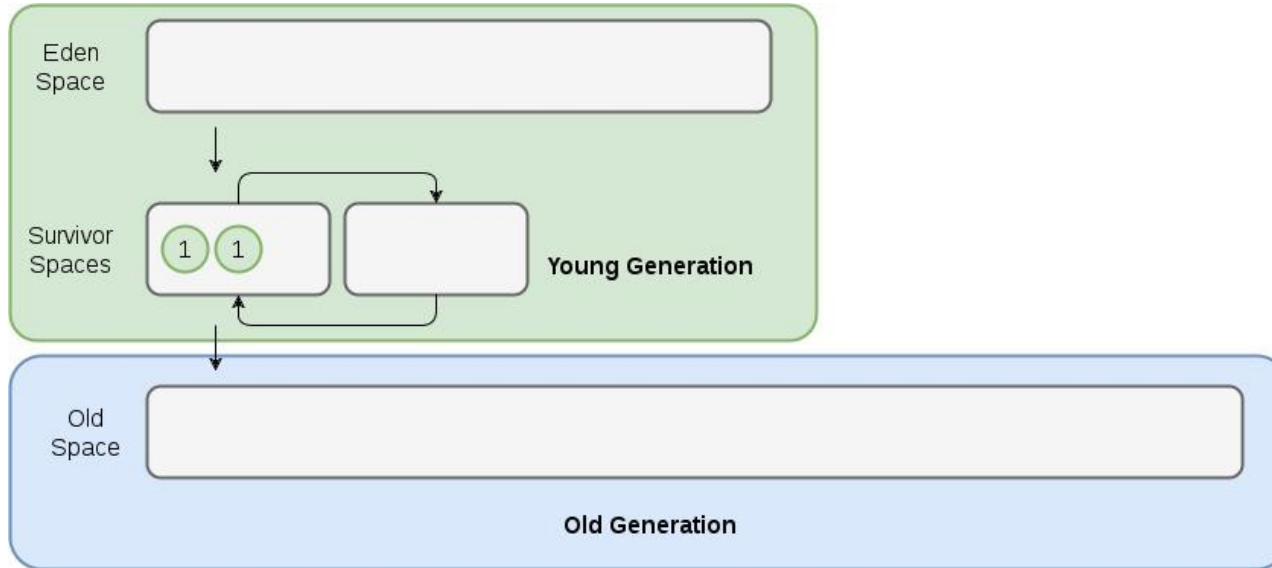
- Two generations:
 - **Young** and **Old**
- Surviving objects are copied to
 - **Survivor** spaces and then to
 - the **Old** generation.

OpenJDK HotSpot Generational GCs (PS, CMS, G1)



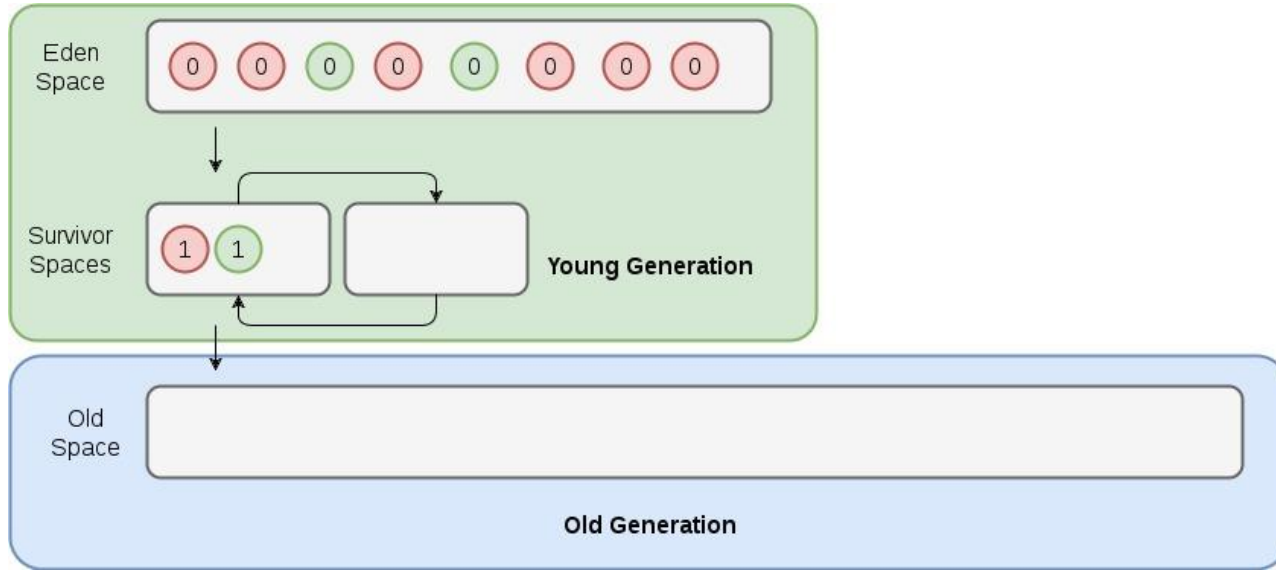
- Two generations:
 - **Young** and **Old**
- Surviving objects are copied to
 - **Survivor** spaces and then to
 - the **Old** generation.

OpenJDK HotSpot Generational GCs (PS, CMS, G1)



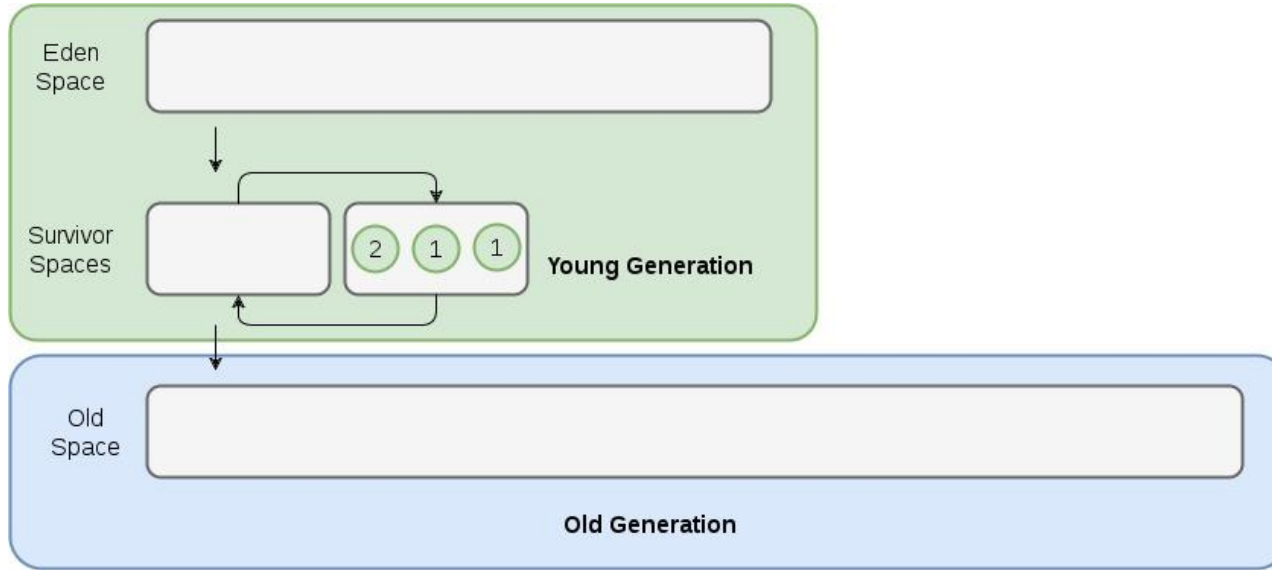
- Two generations:
 - **Young** and **Old**
- Surviving objects are copied to
 - **Survivor** spaces and then to
 - the **Old** generation.

OpenJDK HotSpot Generational GCs (PS, CMS, G1)



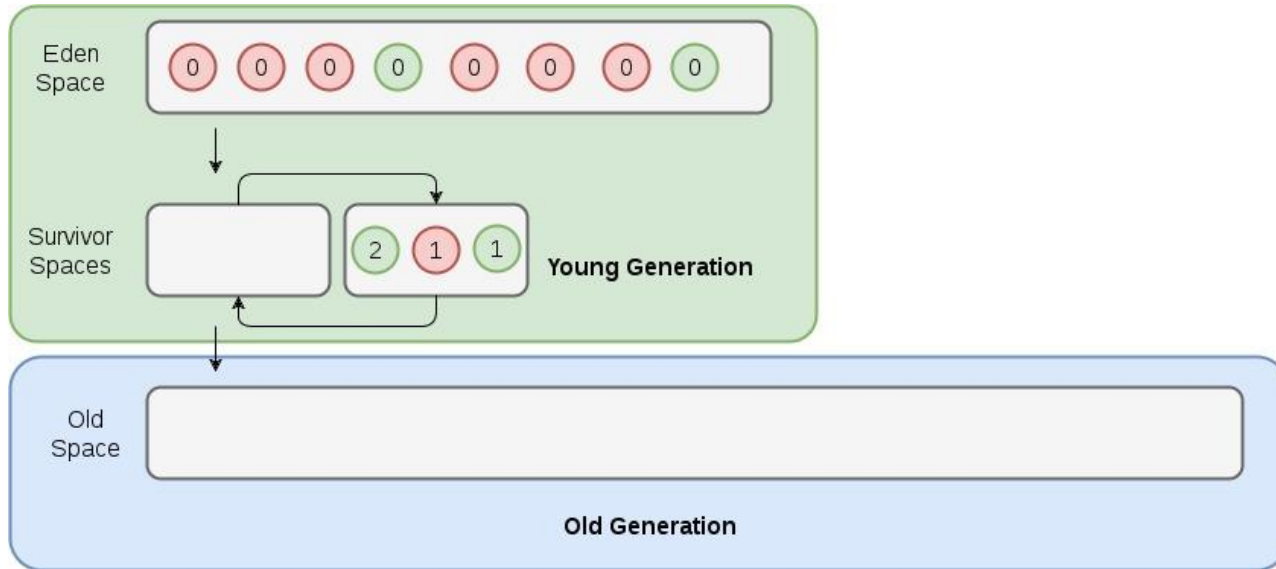
- Two generations:
 - **Young** and **Old**
- Surviving objects are copied to
 - **Survivor** spaces and then to
 - the **Old** generation.

OpenJDK HotSpot Generational GCs (PS, CMS, G1)



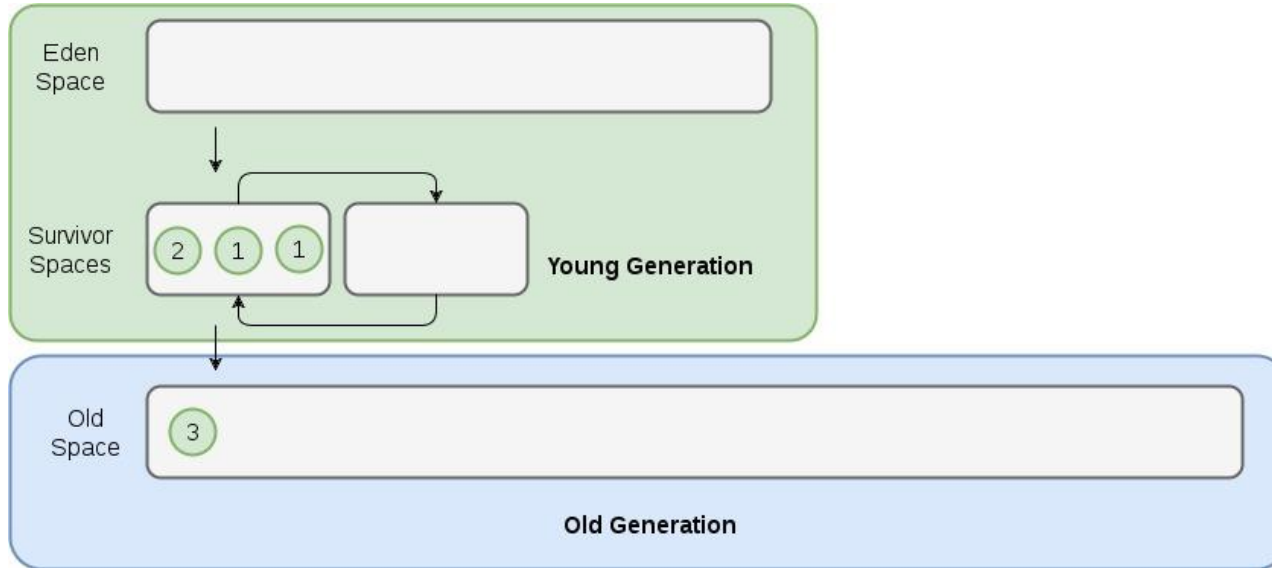
- Two generations:
 - **Young** and **Old**
- Surviving objects are copied to
 - **Survivor** spaces and then to
 - the **Old** generation.

OpenJDK HotSpot Generational GCs (PS, CMS, G1)



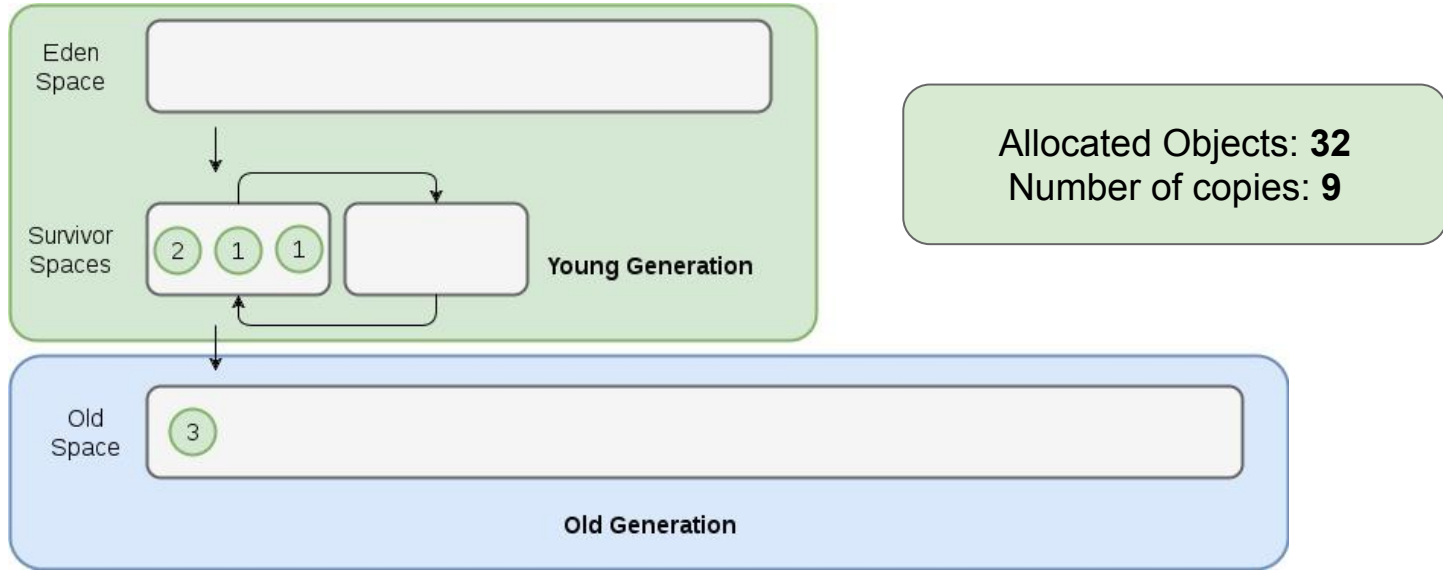
- Two generations:
 - **Young** and **Old**
- Surviving objects are copied to
 - **Survivor** spaces and then to
 - the **Old** generation.

OpenJDK HotSpot Generational GCs (PS, CMS, G1)



- Two generations:
 - **Young** and **Old**
- Surviving objects are copied to
 - **Survivor** spaces and then to
 - the **Old** generation.

OpenJDK HotSpot Generational GCs (PS, CMS, G1)



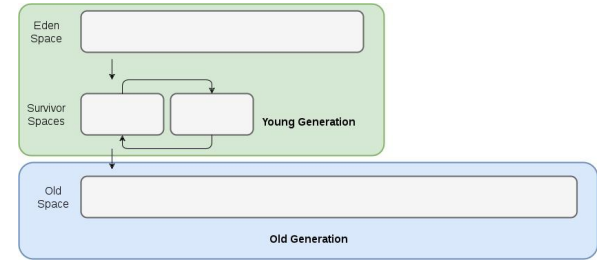
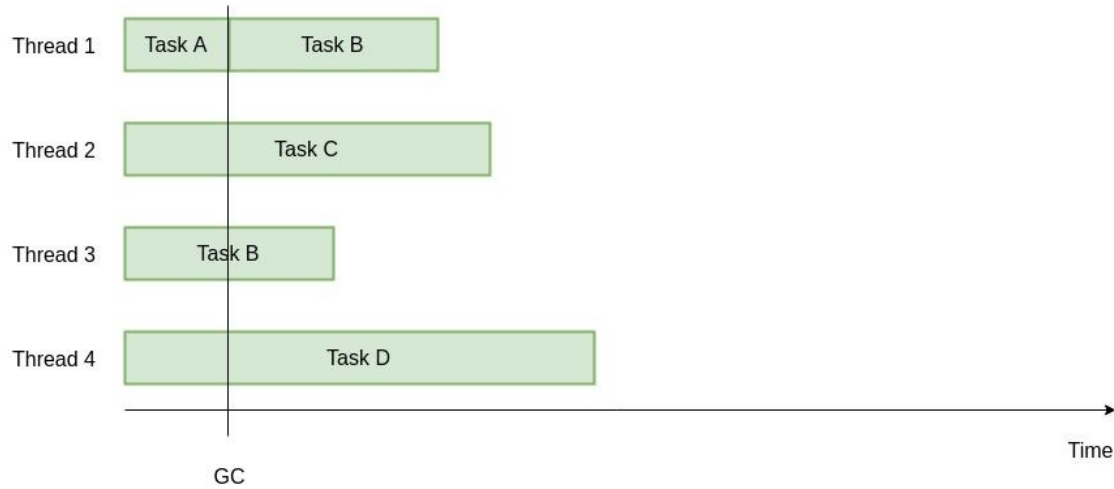
- Two generations:
 - **Young** and **Old**
- Surviving objects are copied to
 - **Survivor** spaces and then to
 - the **Old** generation.

Big Data Application (simplification)

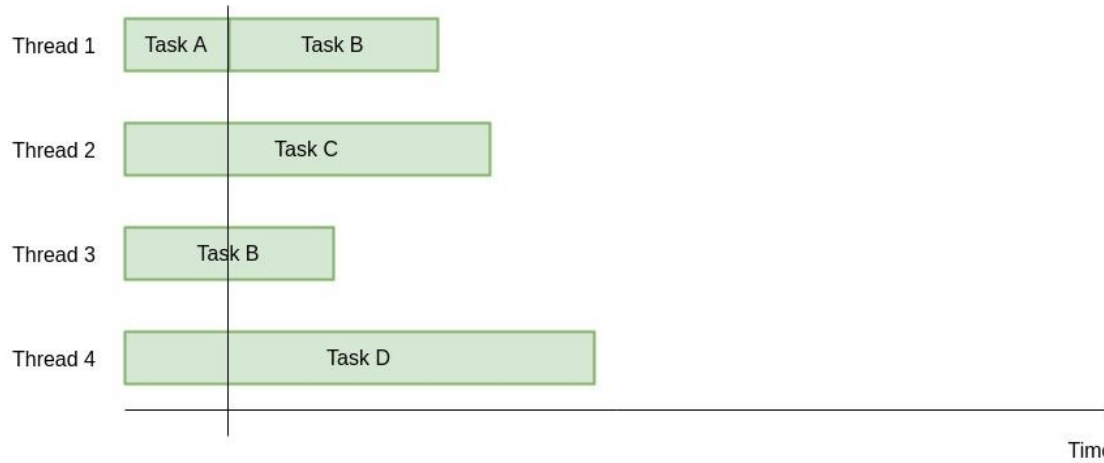
```
1 public void runTask(enum TaskType tt) {  
2  
3     // Allocates memory to hold Working Set  
4     WorkingItem[] buffer = new WorkingItem[WS_SIZE];  
5  
6     // Loads Working Set  
7     DataProvider.load(tt, buffer);  
8  
9     // Process Working Set  
10    Result r = DataProcessor.process(tt, buffer);  
11  
12    // Pushes results from computation  
13    Output.push(r);  
14 }
```

- 4 threads (one per core), running 'runTask' method in loop
- Each task consumes 500 MB of memory (Working Set size)
- Eden is 2GB in size
- Tasks can take different amounts of time to finish

Big Data Application in HotSpot GCs

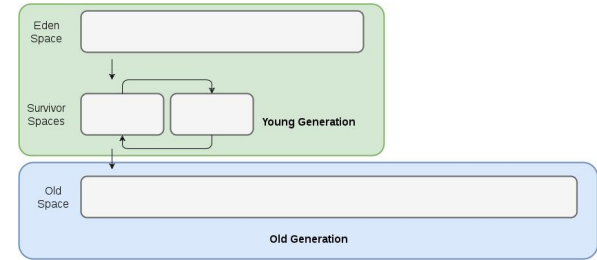


Big Data Application in HotSpot GCs

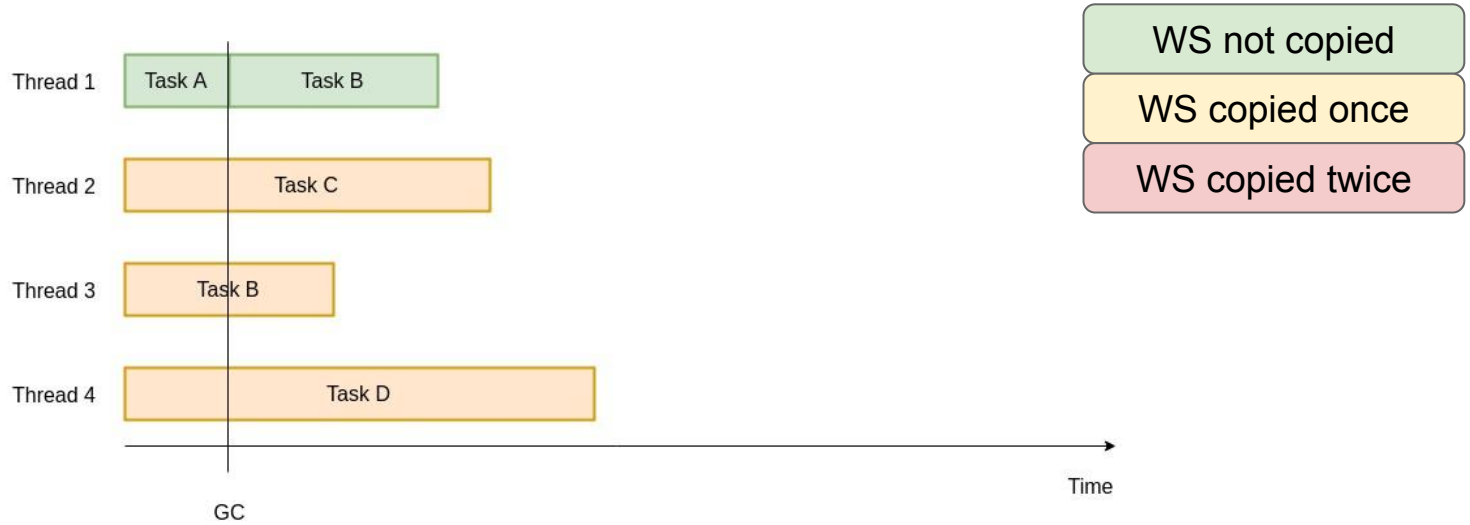


GC

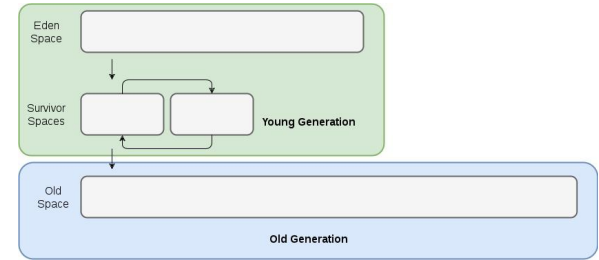
Young generation is full and Thread 1 needs more memory to allocate WS for Task B.



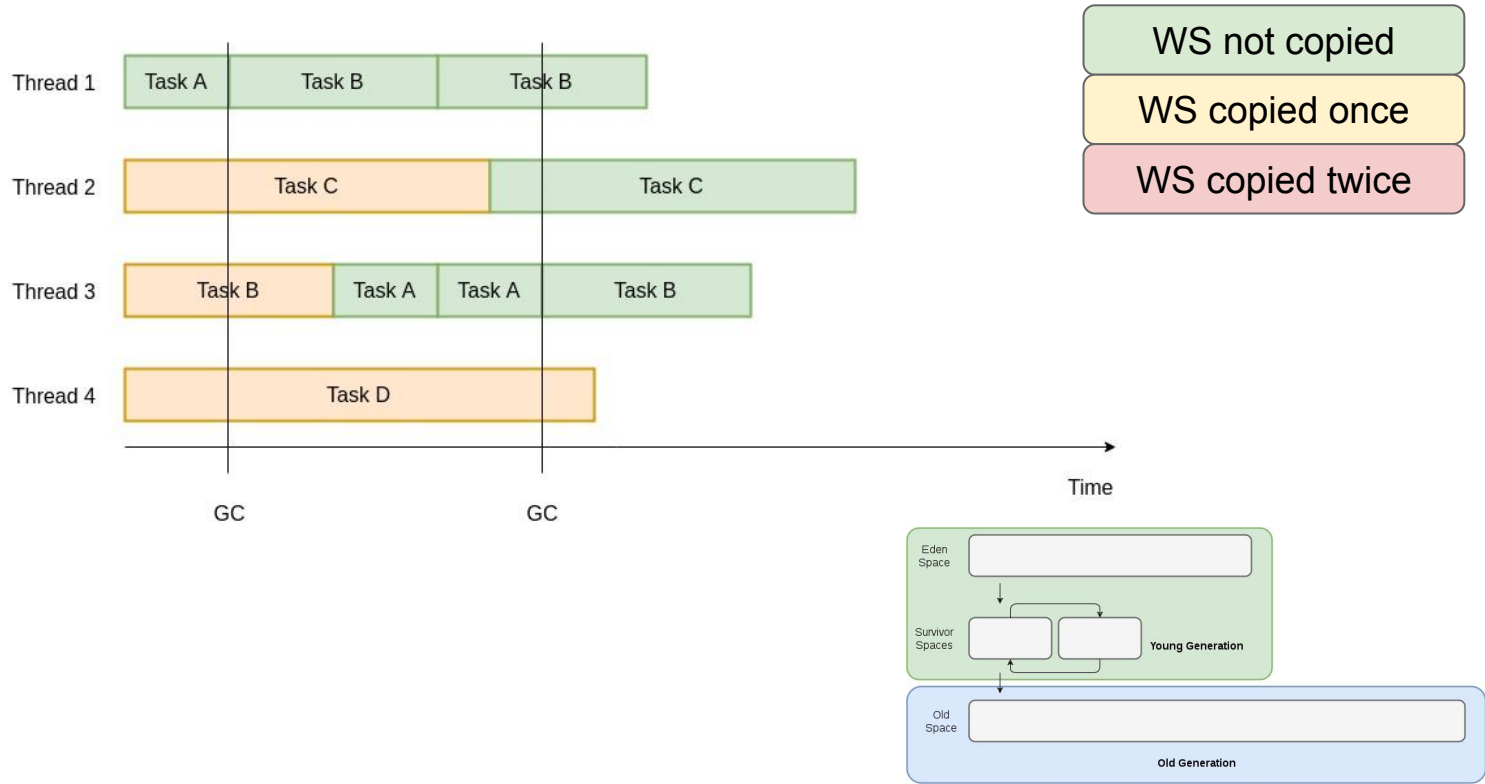
Big Data Application in HotSpot GCs



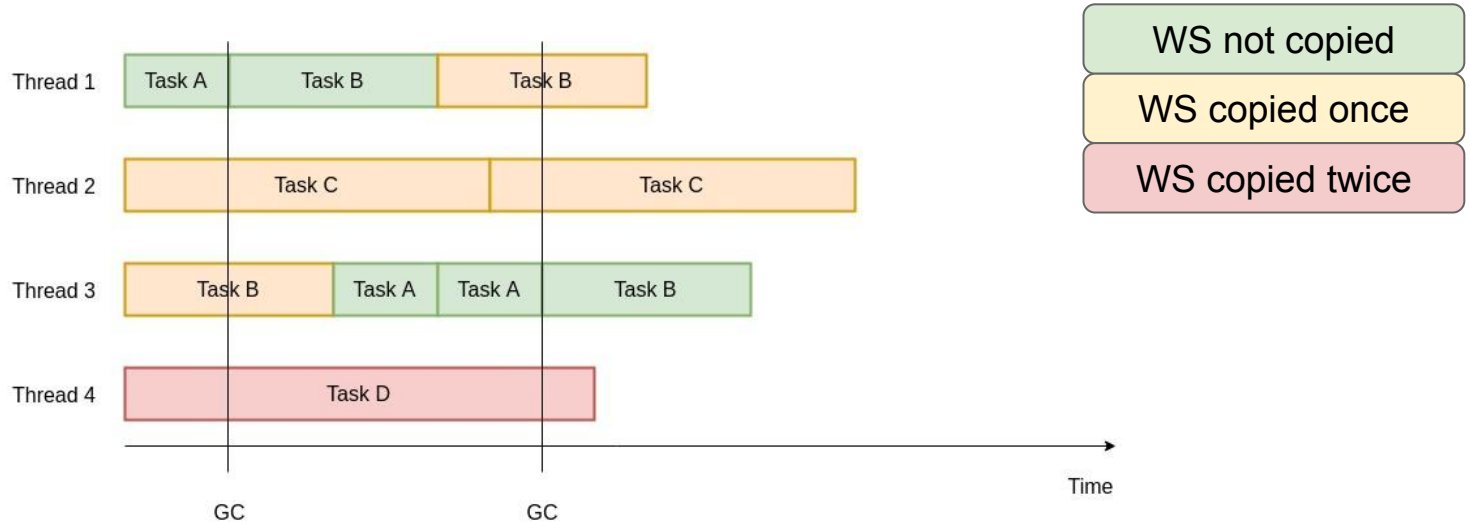
Copies 3 WS = 1500 MB!



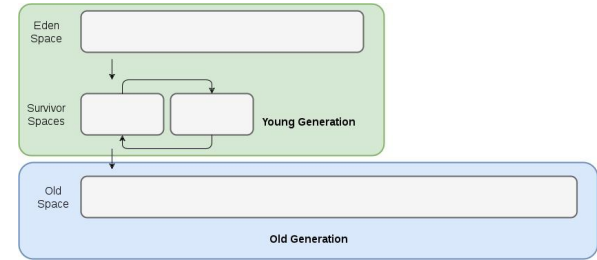
Big Data Application in HotSpot GCs



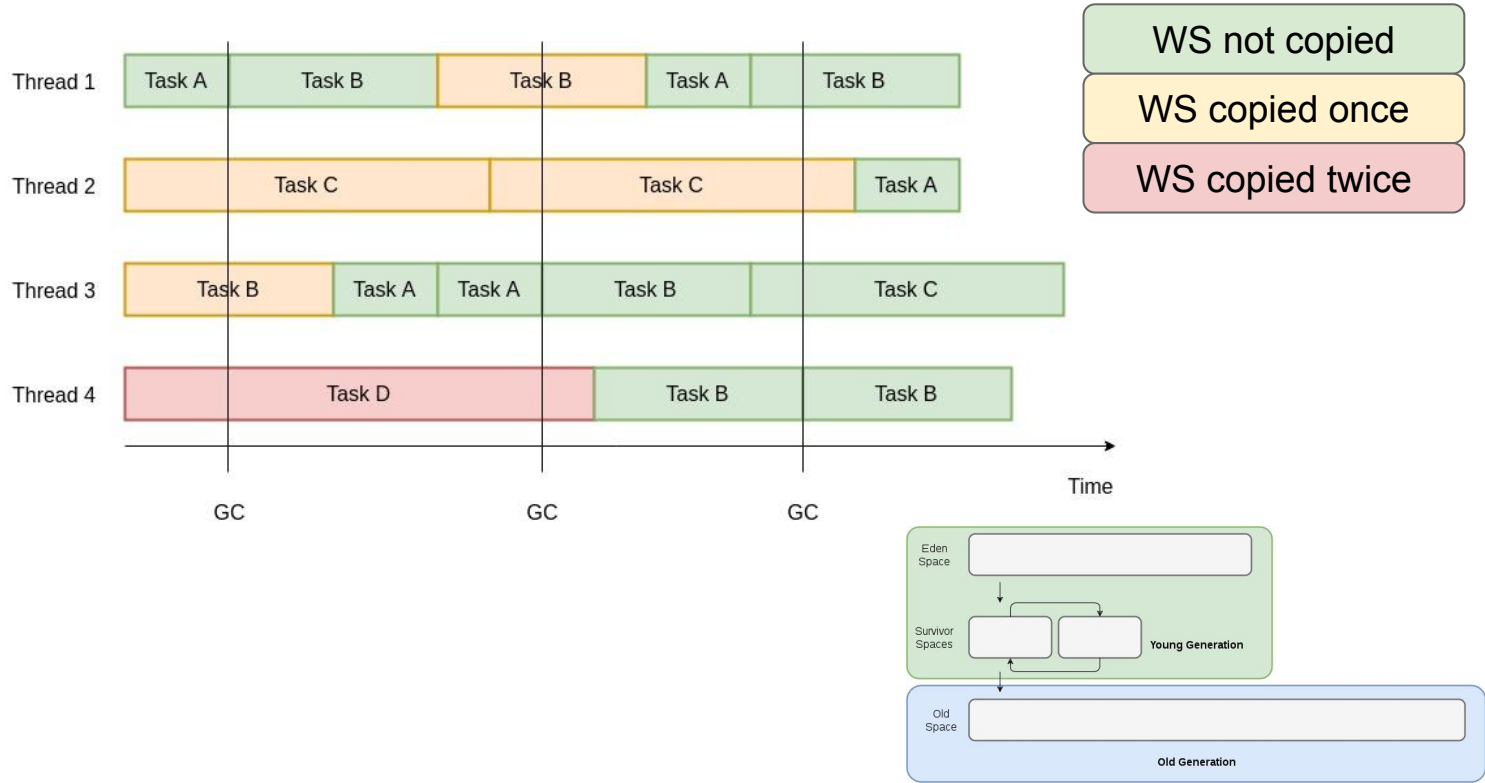
Big Data Application in HotSpot GCs



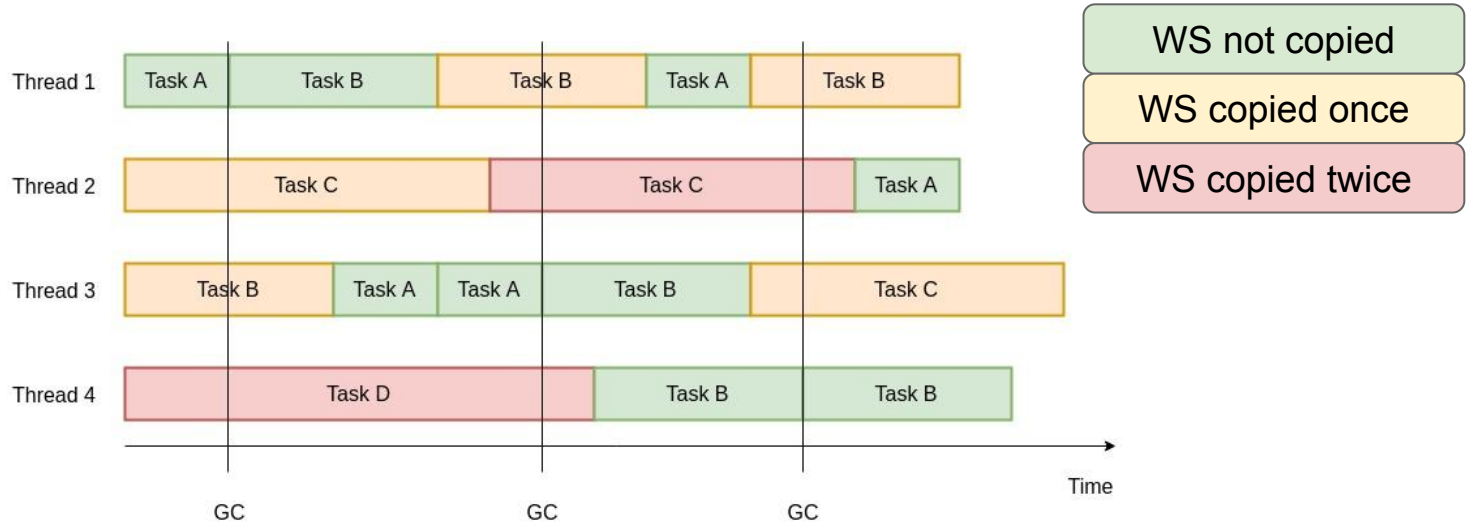
Copies 3 WS = 1500 MB!



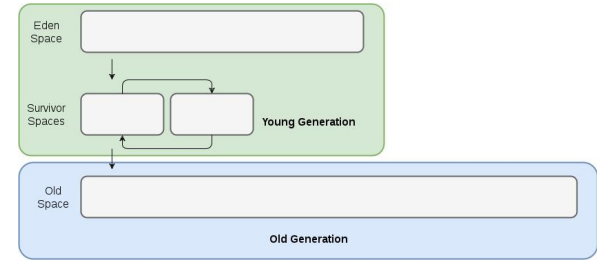
Big Data Application in HotSpot GCs



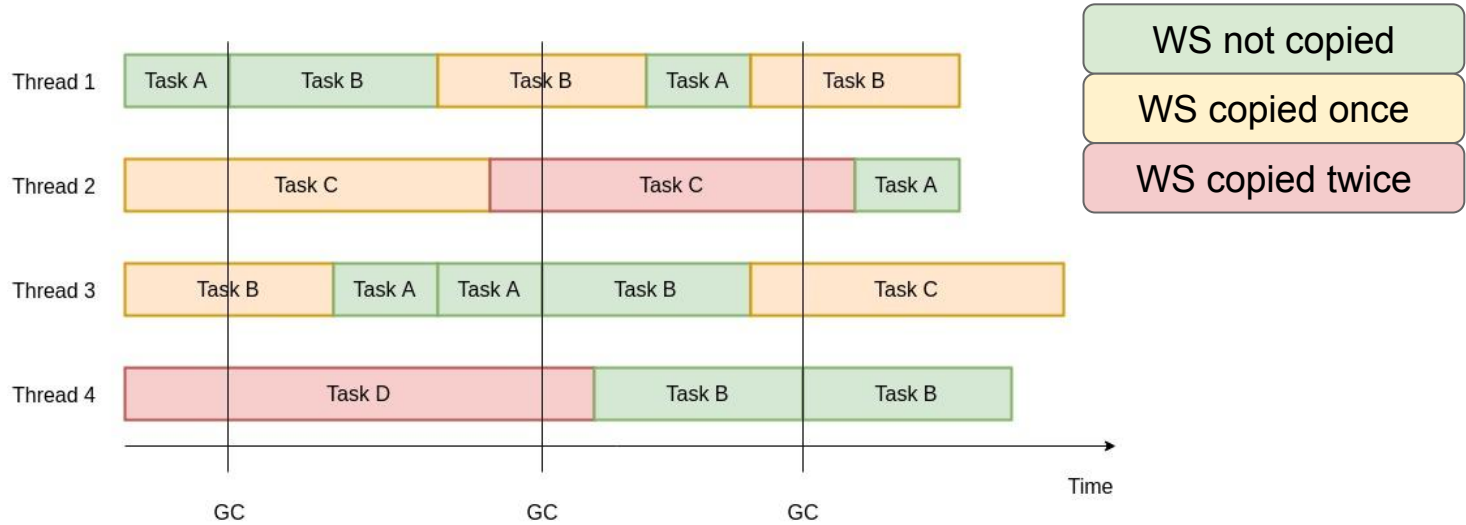
Big Data Application in HotSpot GCs



Copies 3 WS = 1500 MB!

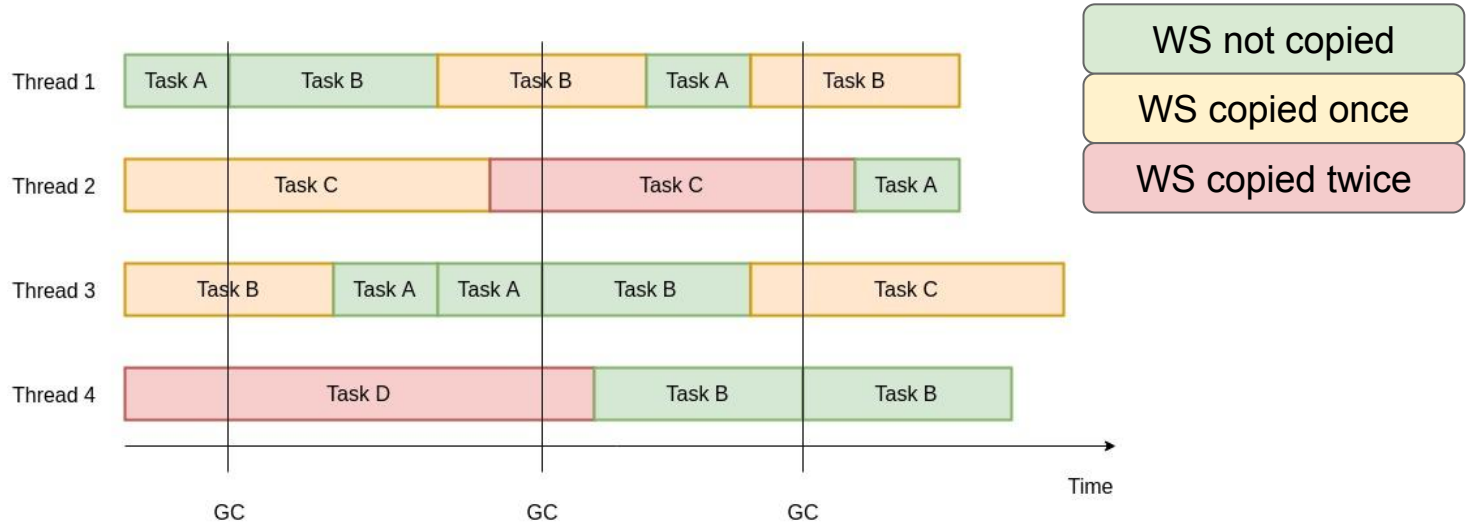


Big Data Application in HotSpot GCs



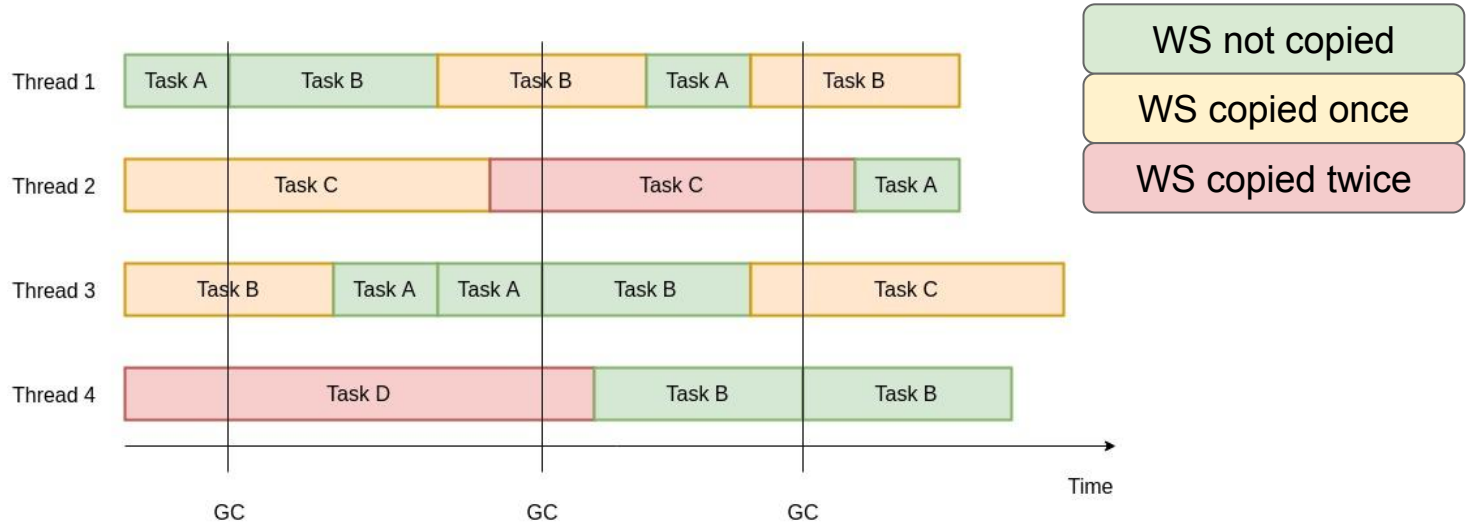
Object copy per GC cycle: 1500 MB
Total amount of object copy: 4500 MB

Big Data Application in HotSpot GCs



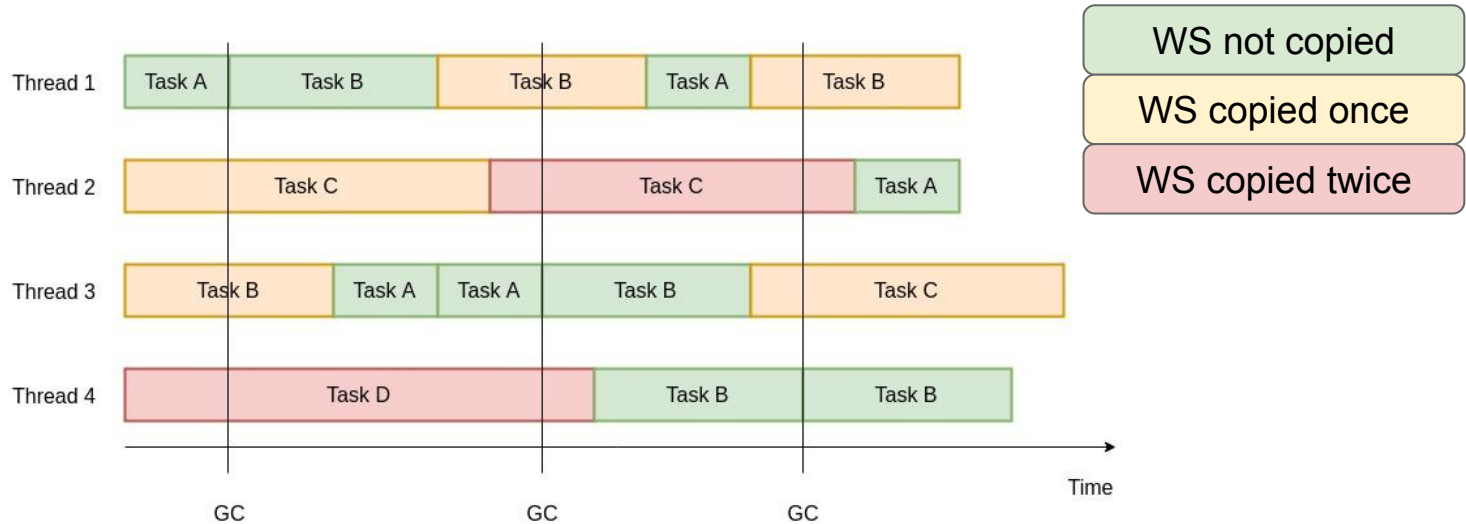
Object copy per GC cycle: 1500 MB
Total amount of object copy: 4500 MB
Assuming average RAM bandwidth of 20GB/s (DDR4)

Big Data Application in HotSpot GCs



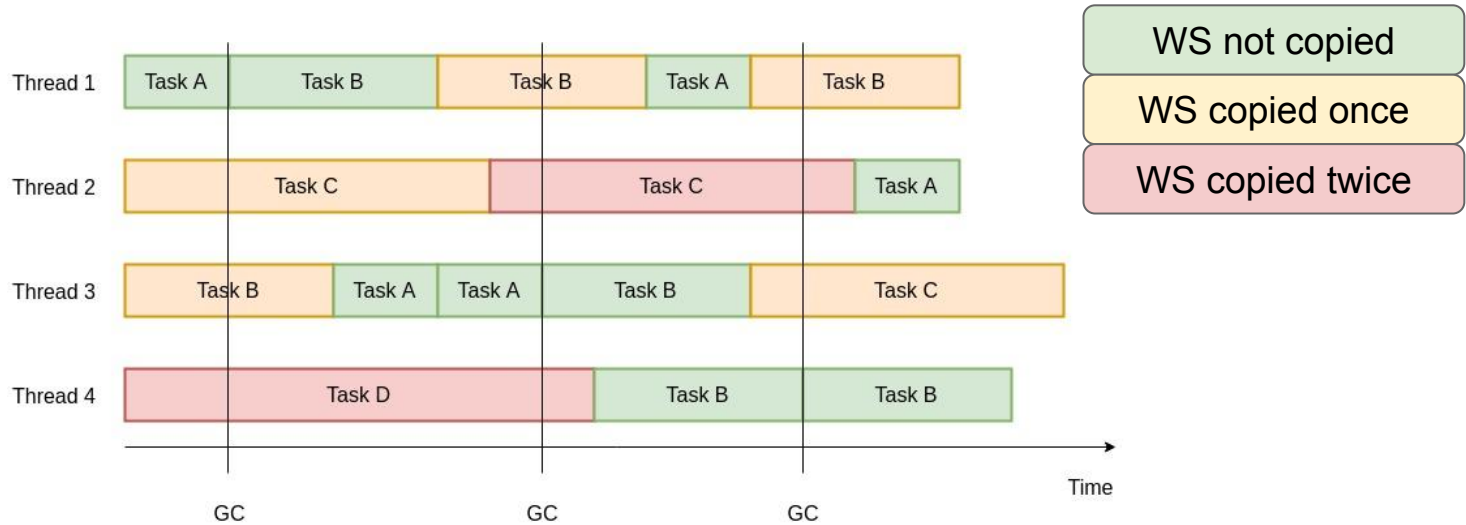
Object copy per GC cycle: 1500 MB
Total amount of object copy: 4500 MB
Assuming average RAM bandwidth of 20GB/s (DDR4)
4 Threads, Eden 2GB = copy 3 tasks (1500 MB) \approx **150 ms**

Big Data Application in HotSpot GCs



Object copy per GC cycle: 1500 MB
Total amount of object copy: 4500 MB
Assuming average RAM bandwidth of 20GB/s (DDR4)
4 Threads, Eden 2GB = copy 3 tasks (1500 MB) \approx **150 ms**
8 Threads, Eden 4GB = copy 7 tasks (3500 MB) \approx **350 ms**

Big Data Application in HotSpot GCs



Object copy per GC cycle: 1500 MB

Total amount of object copy: 4500 MB

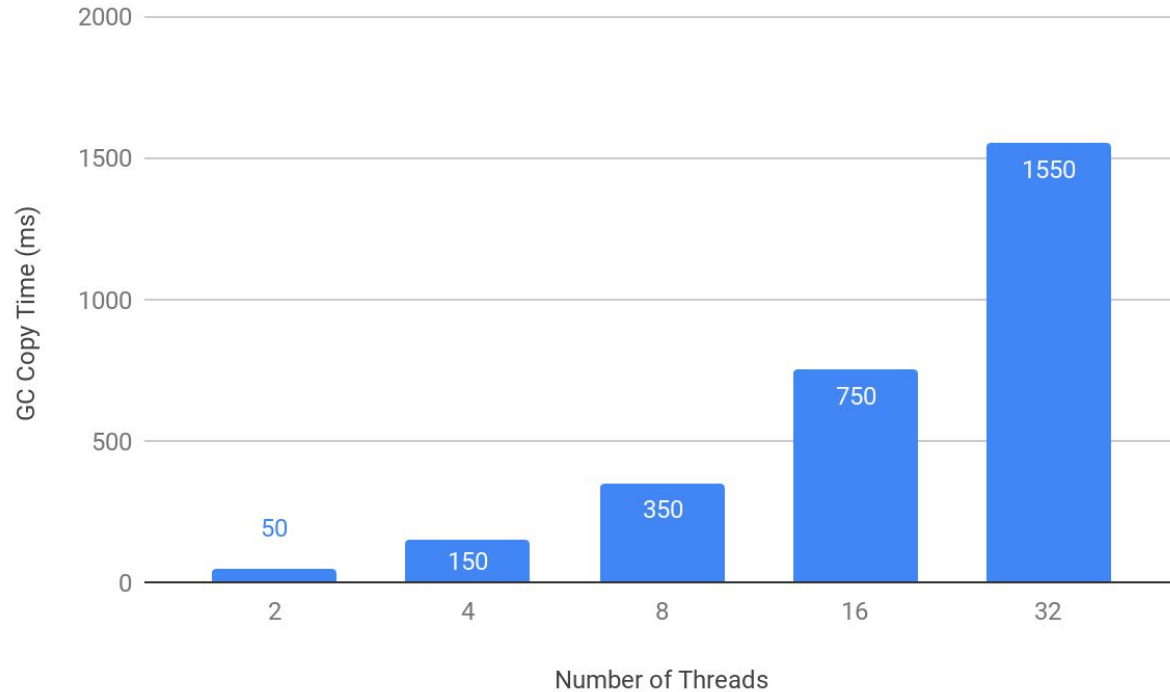
Assuming average RAM bandwidth of 20GB/s (DDR4)

4 Threads, Eden 2GB = copy 3 tasks (1500 MB) \approx **150 ms**

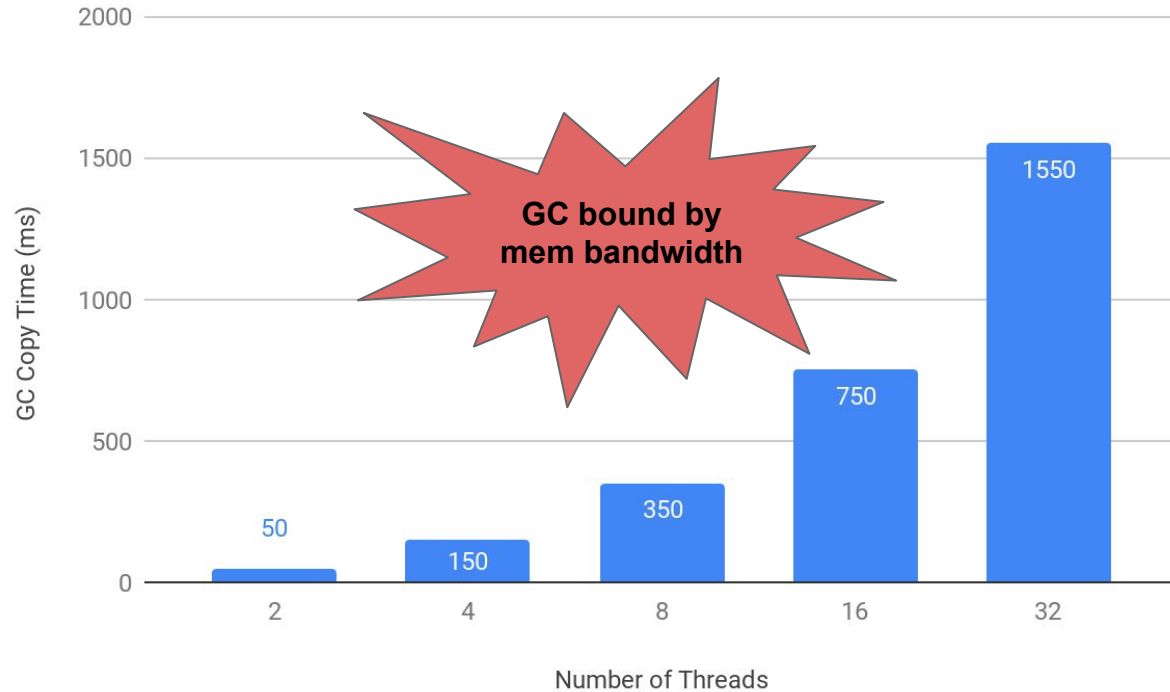
8 Threads, Eden 4GB = copy 7 tasks (3500 MB) \approx **350 ms**

16 Threads, Eden 8GB = copy 15 task (7500 MB) \approx **750 ms**

Big Data Application in HotSpot GCs

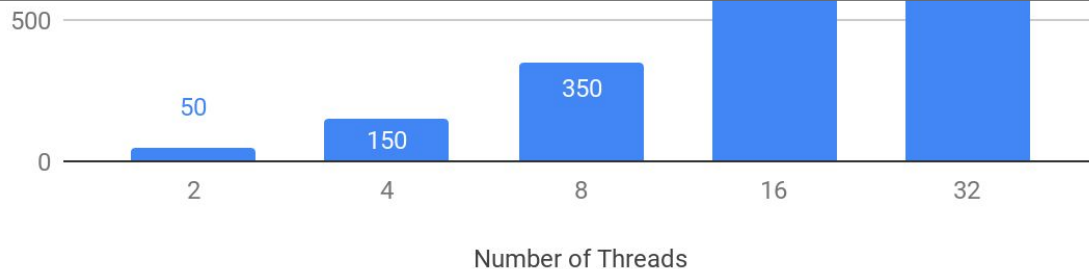


Big Data Application in HotSpot GCs



Big Data Application in HotSpot GCs

Goal: Reduce Application Pauses
by reducing Object Copying
(no negative impact on throughput; no programmer effort)



How to Avoid en-masse Object Copying

- Attempt 1: Heap Resizing

- ✓ Increase Young generation size giving more time for objects to die;
- ! Does not solve the problem, eventually the Young gen will get full and objects will be copied.

How to Avoid en-masse Object Copying

- Attempt 1: Heap Resizing

- ✓ Increase Young generation size giving more time for objects to die;
- ! Does not solve the problem, eventually the Young gen will get full and objects will be copied.

- Attempt 2: Reduce Task/Working Set size

- ✓ Reduces the amount of object copying since the WS is smaller;
- ! Increases overhead as more tasks and coordination is necessary to process smaller tasks.

How to Avoid en-masse Object Copying

- Attempt 1: Heap Resizing

- ✓ Increase Young generation size giving more time for objects to die;
- ! Does not solve the problem, eventually the Young gen will get full and objects will be copied.

- Attempt 2: Reduce Task/Working Set size

- ✓ Reduces the amount of object copying since the WS is smaller;
- ! Increases overhead as more tasks and coordination is necessary to process smaller tasks.

- Attempt 3: Reuse data objects (object pooling)

- ✓ Avoids allocating new memory for future Tasks;
- ! Requires major rewriting of applications combined with very unnatural Java programming style.

How to Avoid en-masse Object Copying

- Attempt 4: Off-heap memory
 - ✓ Reduces GC effort as data objects can reside in off-heap
 - ! Objects describing data objects still reside in the GC-managed heap
 - ! Requires manual memory management (defeats the purpose of running inside a managed heap).

How to Avoid en-masse Object Copying

- **Attempt 4: Off-heap memory**
 - ✓ Reduces GC effort as data objects can reside in off-heap
 - ! Objects describing data objects still reside in the GC-managed heap
 - ! Requires manual memory management (defeats the purpose of running inside a managed heap).

- **Attempt 5: Region-based/Scope-based memory allocation**
 - ✓ Limits object's reachability by scope/region;
 - ! Does not allow objects to freely move between scopes. Bag-of-tasks only, no support for DB!

How to Avoid en-masse Object Copying

- **Attempt 4: Off-heap memory**
 - ✓ Reduces GC effort as data objects can reside in off-heap
 - ! Objects describing data objects still reside in the GC-managed heap
 - ! Requires manual memory management (defeats the purpose of running inside a managed heap).
- **Attempt 5: Region-based/Scope-based memory allocation**
 - ✓ Limits object's reachability by scope/region;
 - ! Does not allow objects to freely move between scopes. Bag-of-tasks only, no support for DB!
- **Attempt 6: Completely Concurrent Collectors (C4, Shenandoah, ZGC)**
 - ✓ Greatly reduced pause times
 - ! High throughput overhead (~30% for Cassandra workloads)

How to Avoid en-masse Object Copying

Takeaway:

- Avoiding massive object copying is non-trivial!
- Existing solutions only alleviate the problem!
- Existing solutions might work in some scenarios but do not provide a general solution.

Proposed Solution

- **Solution:**
 - Allocate objects with similar lifetimes close to each other
 - Reducing memory fragmentation
 - Reducing object promotion
 - As a consequence, object copying is reduced!

Proposed Solution

- **Solution:**

- Allocate objects with similar lifetimes close to each other
 - Reducing memory fragmentation
 - Reducing object promotion
- As a consequence, object copying is reduced!

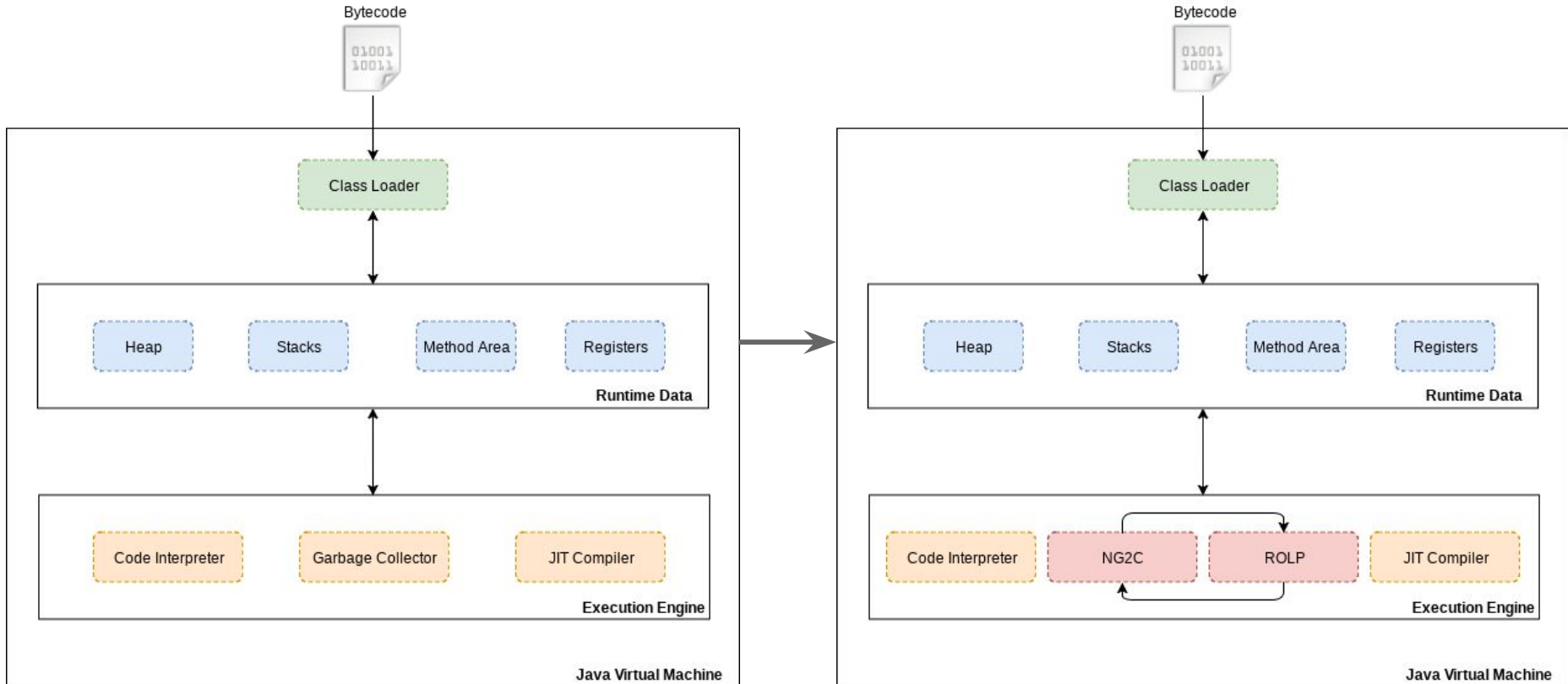
- **Hypothesis:**

- Objects allocated through the same allocation context have similar lifetimes;
- Allocation context is a tuple of:
 - Allocation site (line of bytecode)
 - Call graph state (stack state)

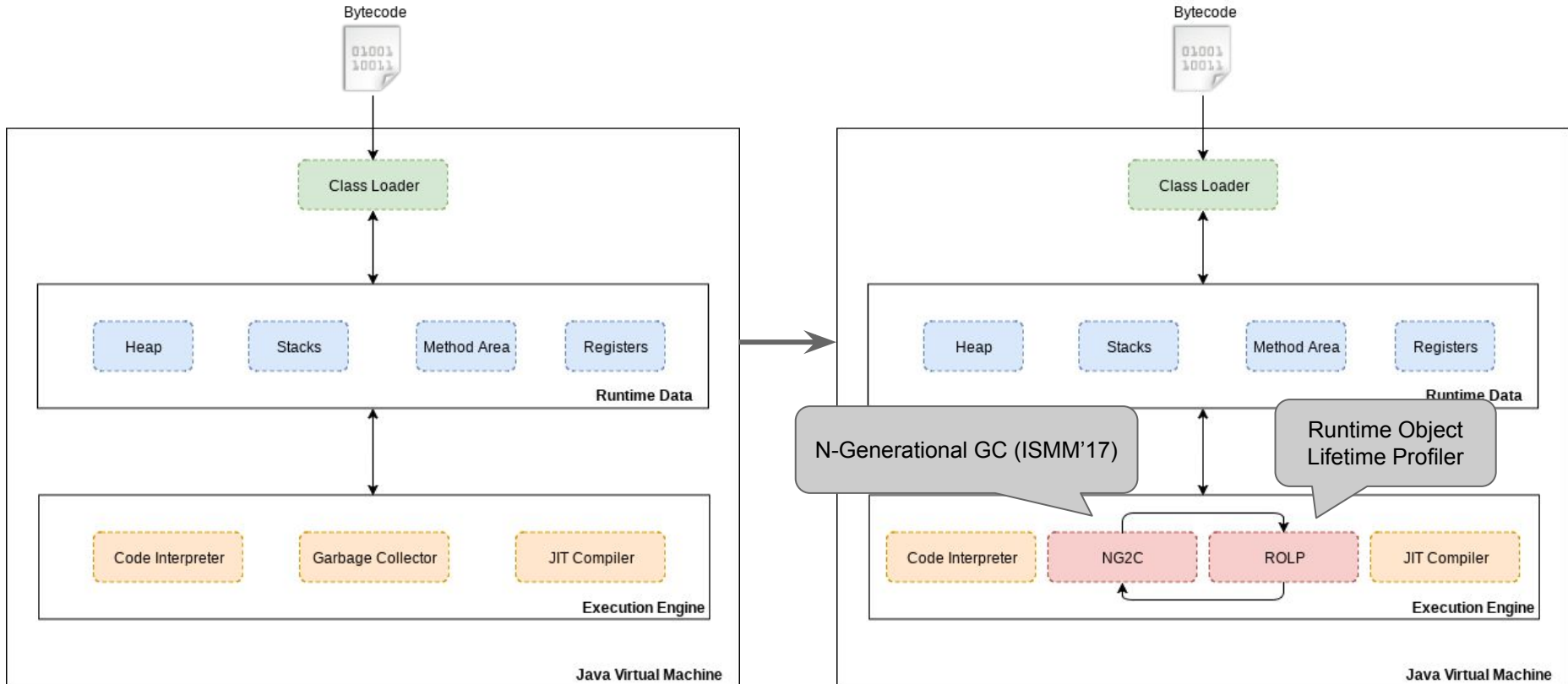
Big Data Application (simplification)

```
1 public void runTask(enum TaskType tt) {
2
3     // Allocates memory to hold Working Set
4     WorkingItem[] buffer = new WorkingItem[WS_SIZE];
5
6     // Loads Working Set
7     DataProvider.load(tt, buffer);
8
9     // Process Working Set
10    Result r = DataProcessor.process(tt, buffer);
11
12    // Pushes results from computation
13    Output.push(r);
14 }
```

Solution - NG2C + ROLP

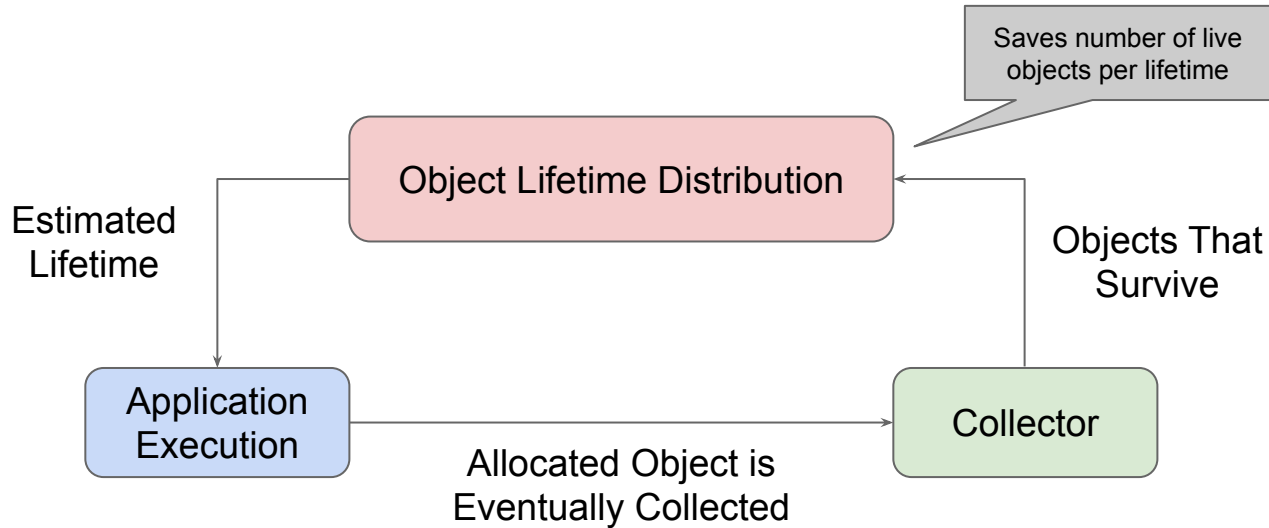


Solution - NG2C + ROLP



Runtime Object Lifetime Profiler (overview)

- Profiler needs to answer a single question
 - How long will objects allocated through a particular allocation context live?

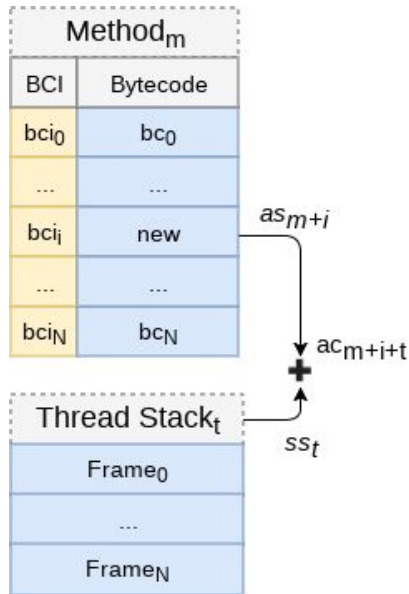


Runtime Object Lifetime Profiler (workflow)

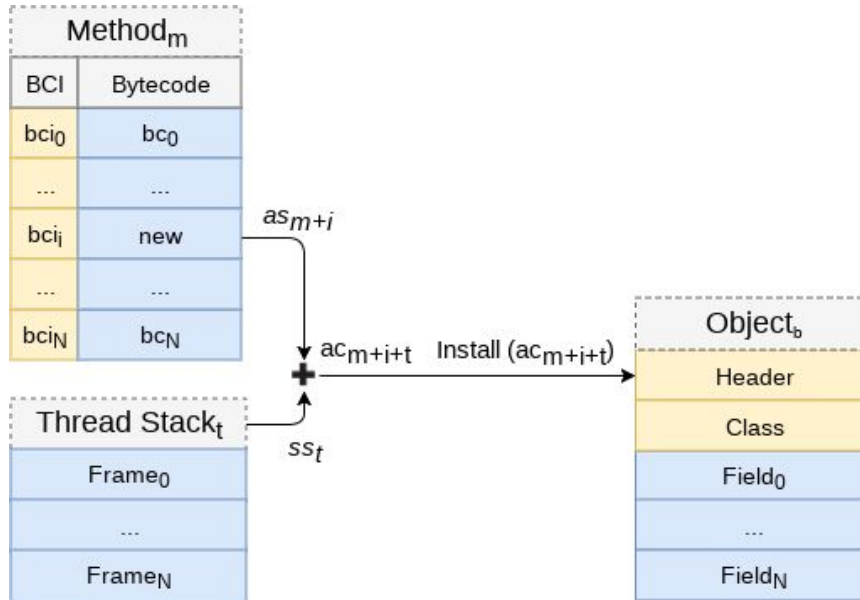
Method _m	
BCI	Bytecode
bci ₀	bc ₀
...	...
bci _j	new
...	...
bci _N	bc _N

Thread Stack _t
Frame ₀
...
Frame _N

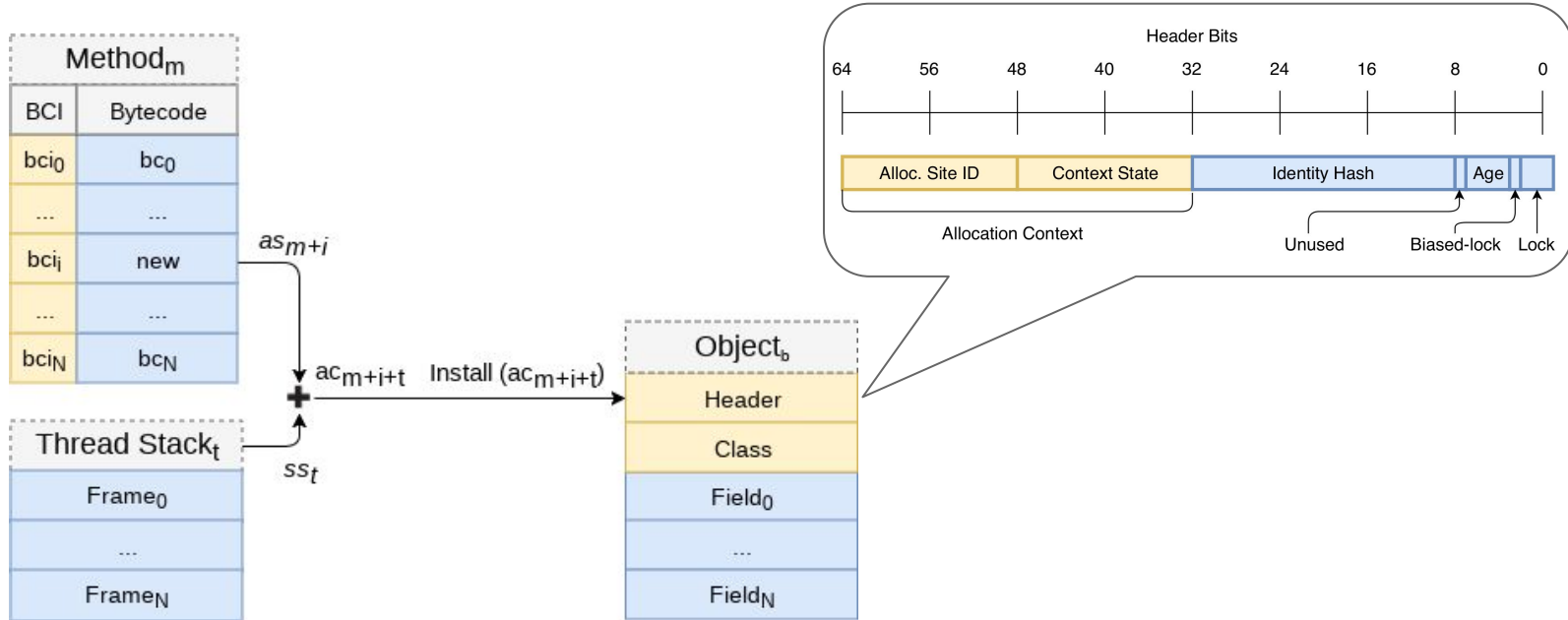
Runtime Object Lifetime Profiler (workflow)



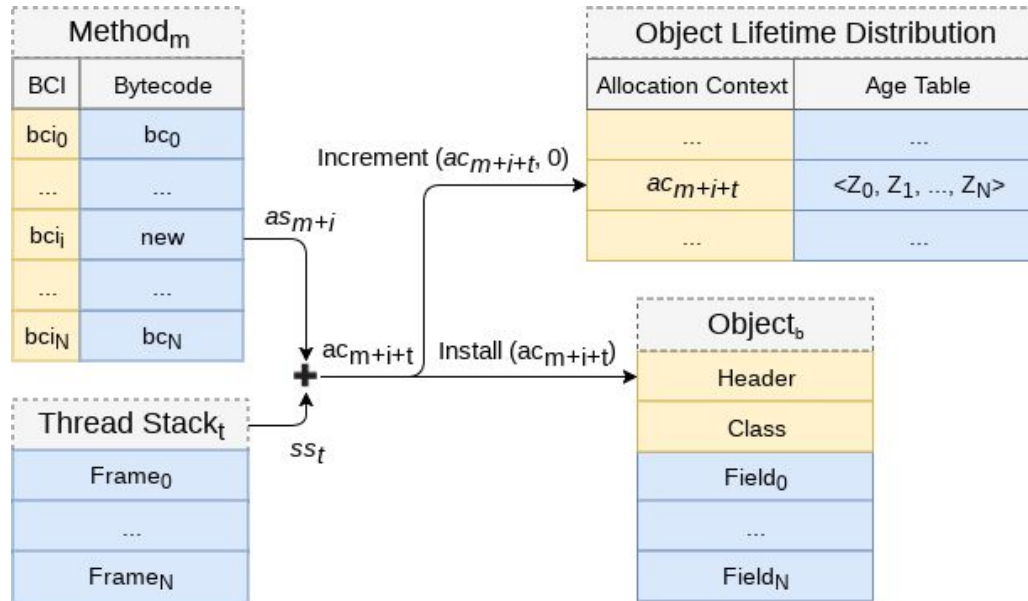
Runtime Object Lifetime Profiler (workflow)



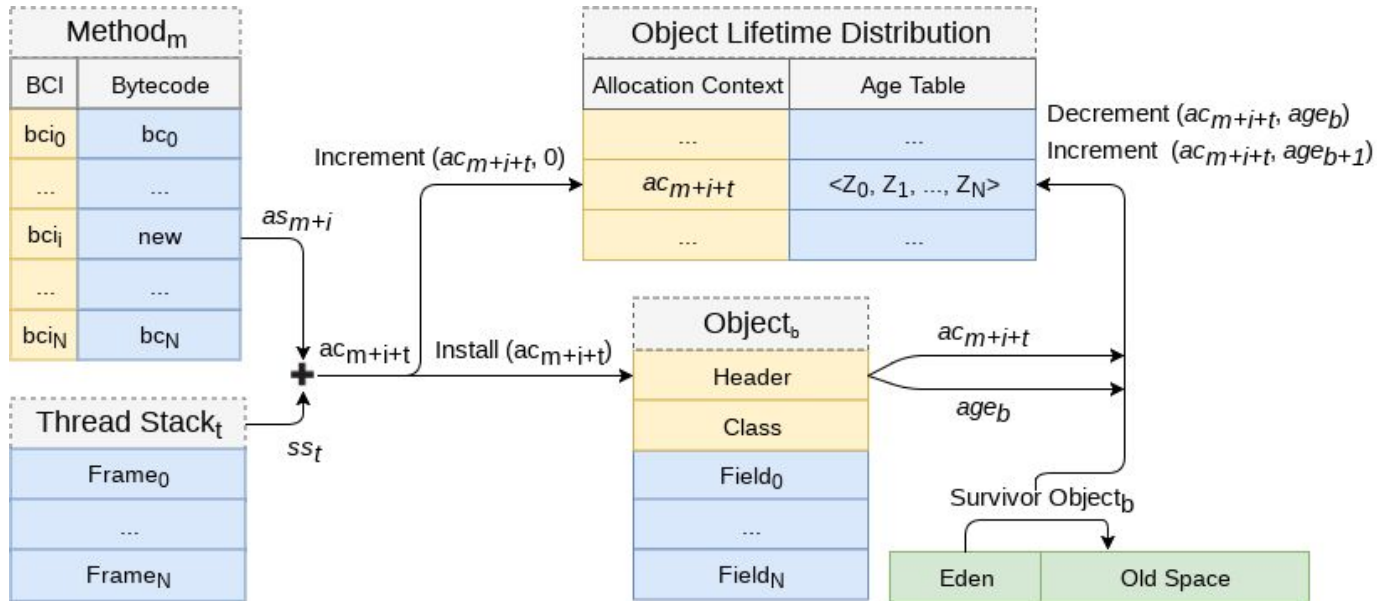
Runtime Object Lifetime Profiler (workflow)



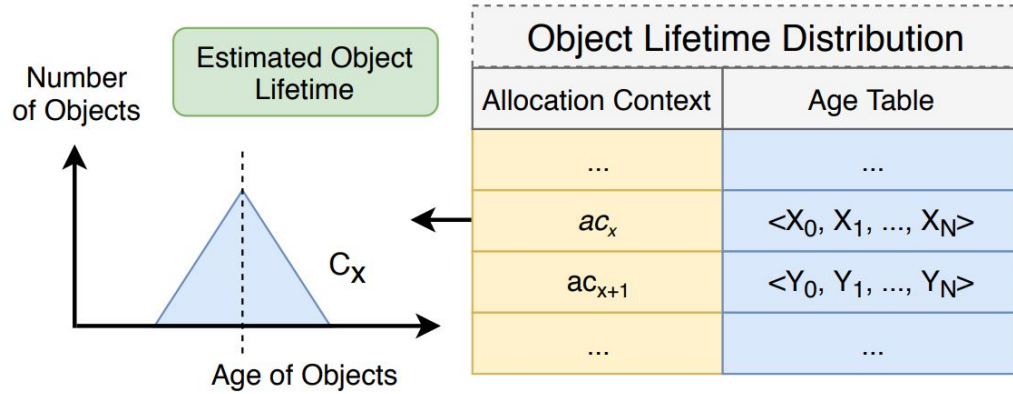
Runtime Object Lifetime Profiler (workflow)



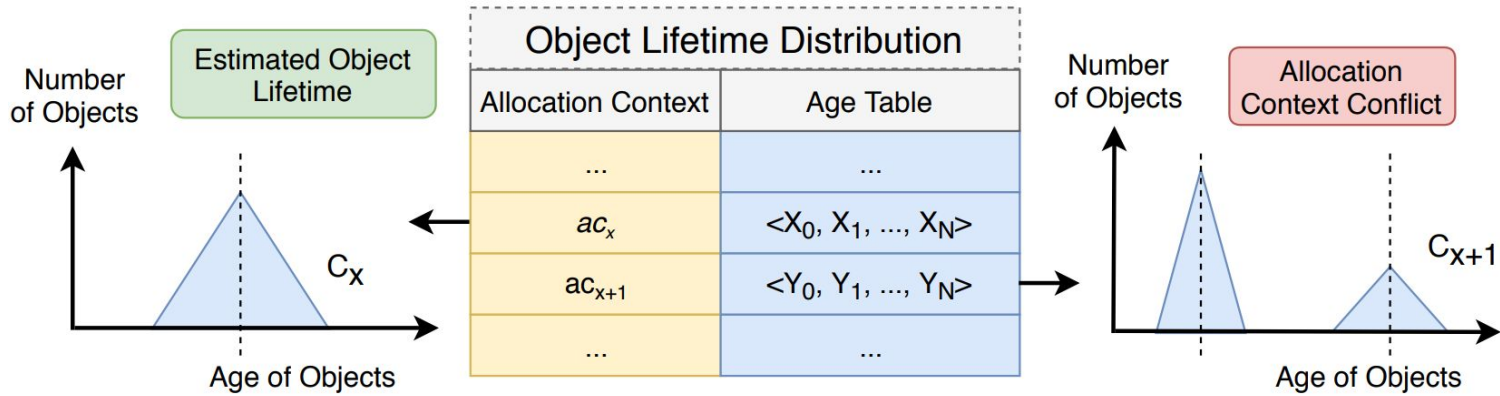
Runtime Object Lifetime Profiler (workflow)



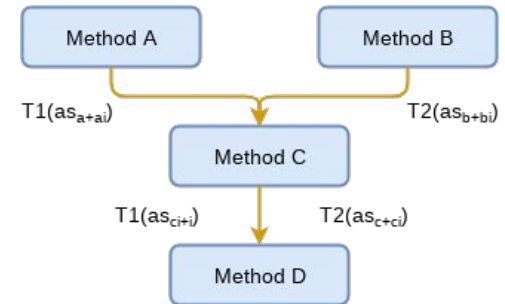
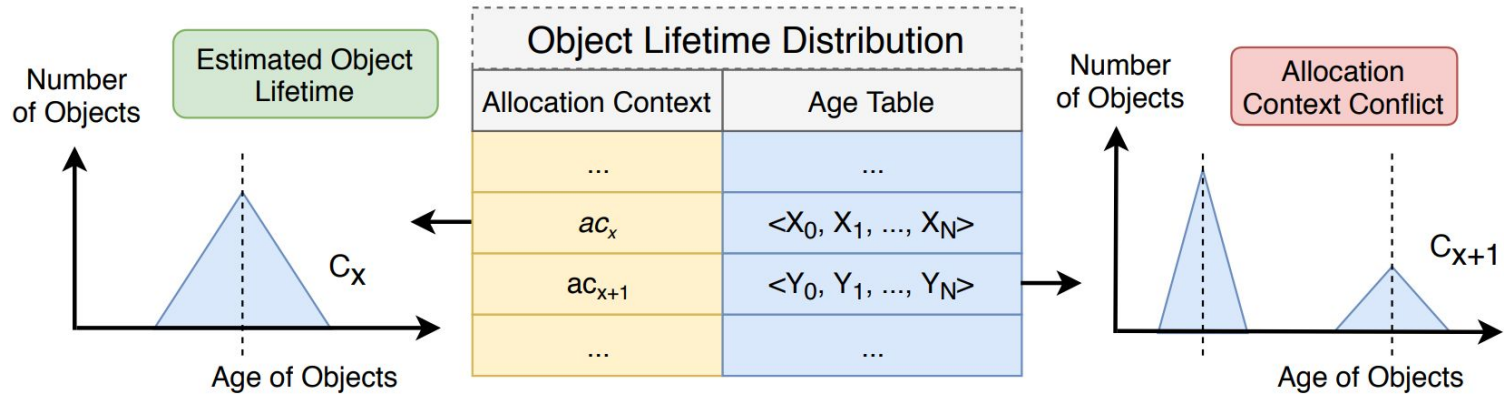
Runtime Object Lifetime Profiler (conflicts)



Runtime Object Lifetime Profiler (conflicts)



Runtime Object Lifetime Profiler (conflicts)

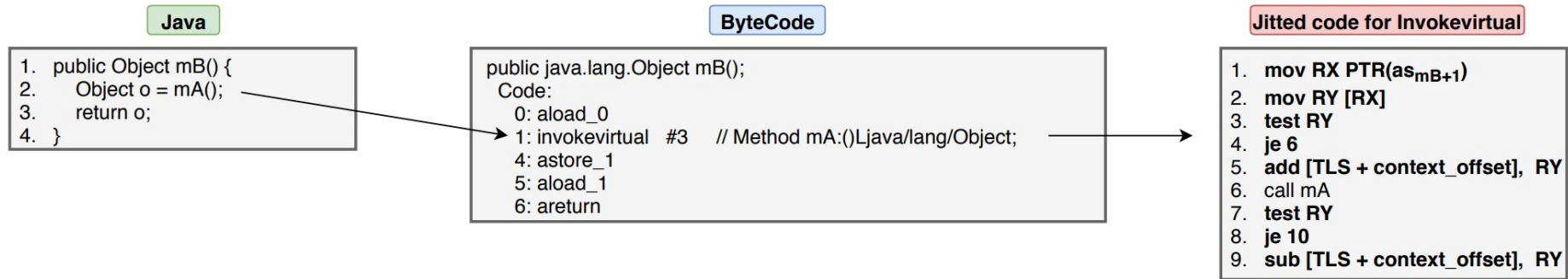


Runtime Object Lifetime Profiler (context tracking)

- Method calls are wrapped with context tracking code (update thread stack state)
- Context tracking is very expensive
 - Only method calls that can resolve conflicts are profiled (next slide)
 - Jitted code can dynamically enable or disable profiling

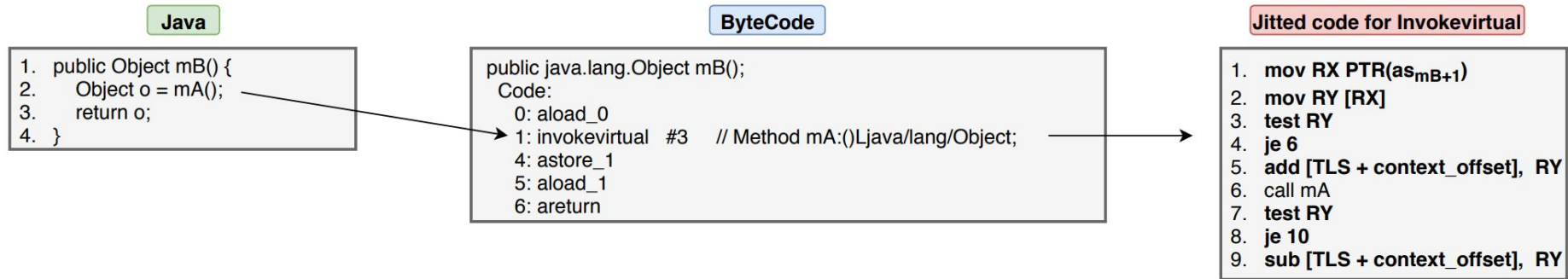
Runtime Object Lifetime Profiler (context tracking)

- Method calls are wrapped with context tracking code (update thread stack state)
- Context tracking is very expensive
 - Only method calls that can resolve conflicts are profiled
 - Jitted code can dynamically enable or disable profiling



Runtime Object Lifetime Profiler (context tracking)

- Method calls are wrapped with context tracking code (update thread stack state)
- Context tracking is very expensive
 - Only method calls that can resolve conflicts are profiled
 - Jitted code can dynamically enable or disable profiling



Method call profiling can be dynamically turned on or off!

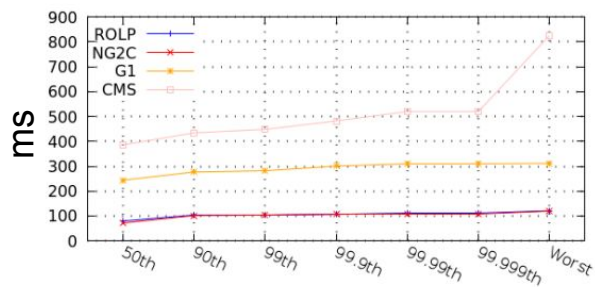
Runtime Object Lifetime Profiler (implementation)

- ROLP is implemented on OpenJDK HotSpot 8
 - Industrial JVM
- ROLP is integrated with NG2C
- ROLP is meant to be running in production workloads
 - Several of implementation/performance optimizations
 - Avoid inlined methods
 - Properly Handling Exceptions
 - Properly Handling On-Stack Replacement
 - Reducing Profiling Overhead for very large applications
 - Shutdown survivor tracking code to reduce overhead
 - Improving the scalability of the Object Lifetime Distribution table
 - ...

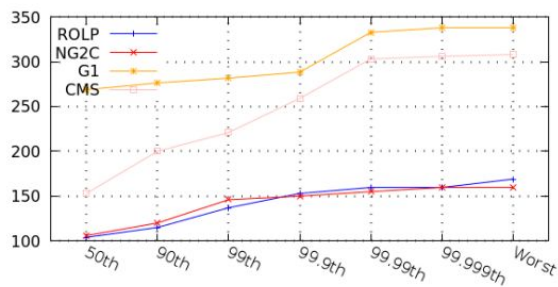
Runtime Object Lifetime Profiler (evaluation)

- Evaluate ROLP's performance compared to:
 - G1 - best solution in OpenJDK, current default GC (ISMM'04)
 - NG2C - multi-generational GC (ISMM'17) - requires programmer effort and knowledge
 - CMS - concurrent mark-sweep - throughput oriented
- Big Data Platforms & Workloads:
 - Cassandra (**Key-Value Store**)
 - YCSB: Write-Intensive (75% writes), Read-Write (50% writes), Read-Intensive (75% reads)
 - Lucene (**In-Memory Indexing Tool**)
 - Read/Write transactions on Wikipedia dump (33M documents): Write-intensive (80% writes)
 - GraphChi (**Graph Processing Engine**)
 - Twitter graph dump (42M vertexes, 1.5B edges): PageRank, Connected Components
 - Environment:
 - Intel Xeon E5505, 16GB RAM
 - Heap/Young Size: 12/2GB

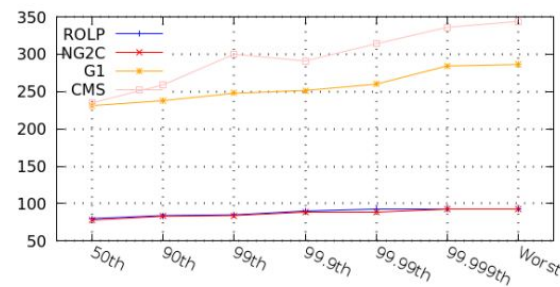
Runtime Object Lifetime Profiler (pausetime percentiles)



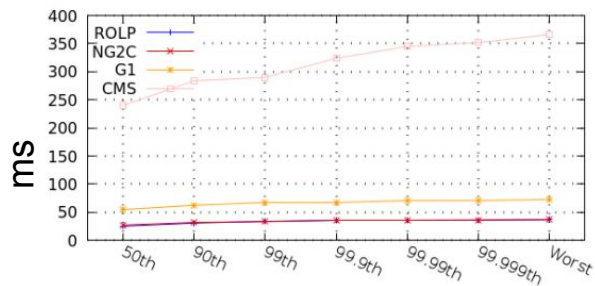
(a) Cassandra WI



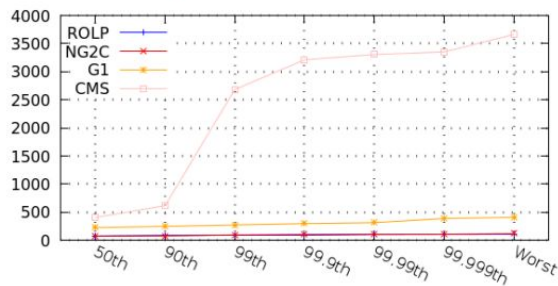
(b) Cassandra WR



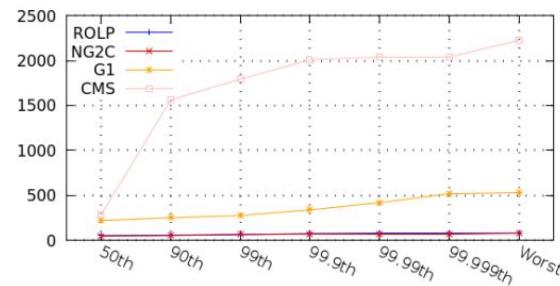
(c) Cassandra RI



(d) Lucene



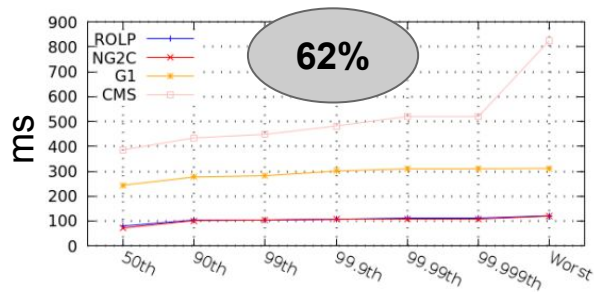
(e) GraphChi CC



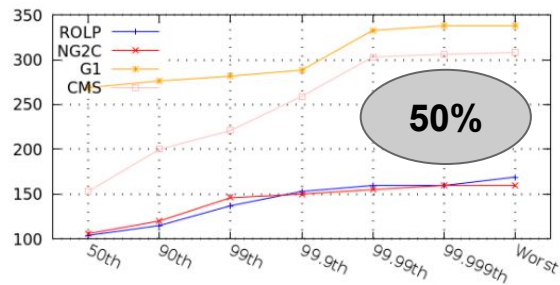
(f) GraphChi PR

Lower is Better

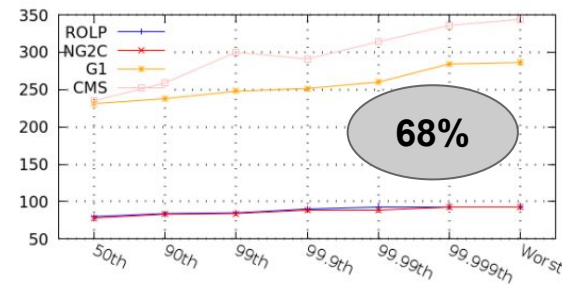
Runtime Object Lifetime Profiler (pausetime percentiles)



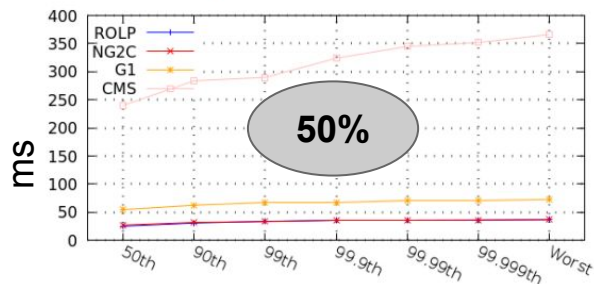
(a) Cassandra WI



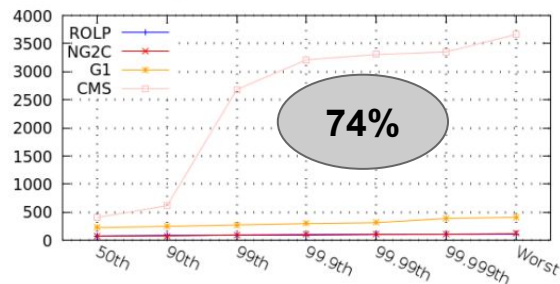
(b) Cassandra WR



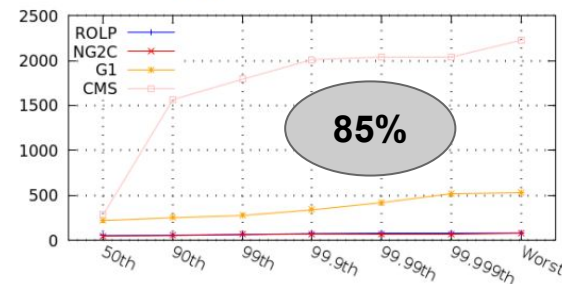
(c) Cassandra RI



(d) Lucene



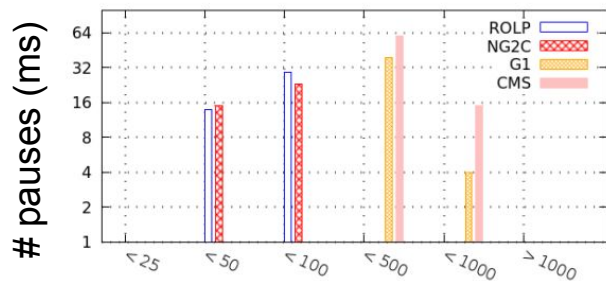
(e) GraphChi CC



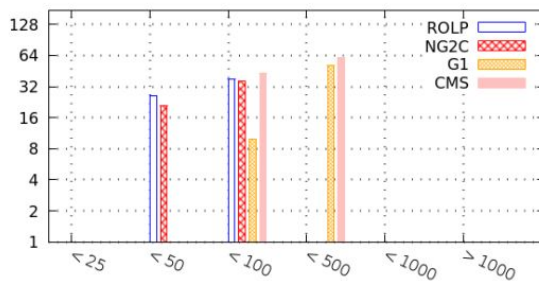
(f) GraphChi PR

Lower is Better

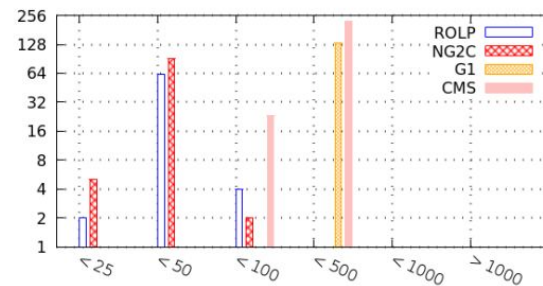
Runtime Object Lifetime Profiler (pausetime distribution)



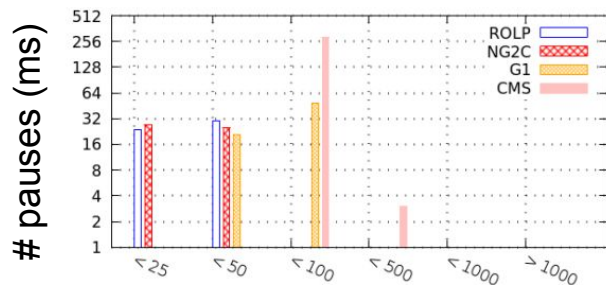
(a) Cassandra WI



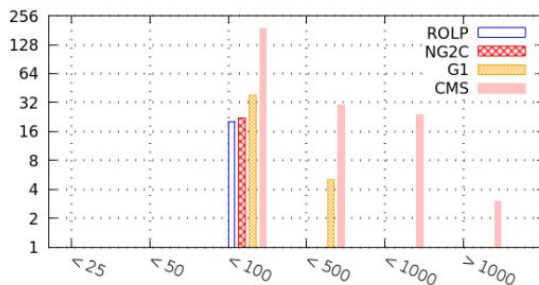
(b) Cassandra WR



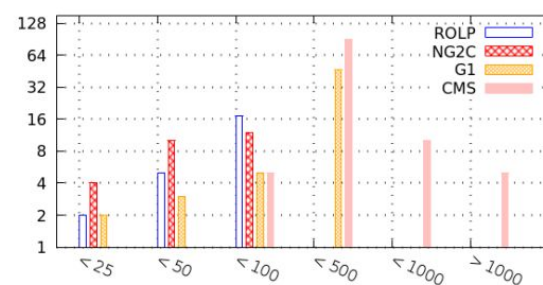
(c) Cassandra RI



(d) Lucene



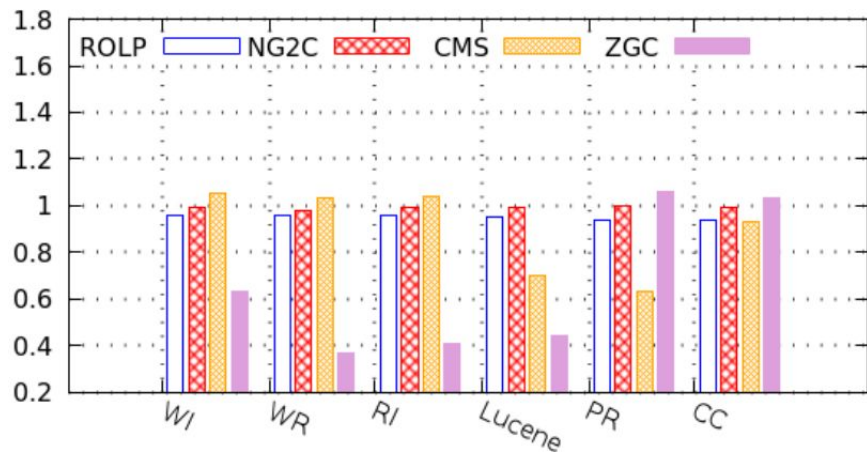
(e) GraphChi CC



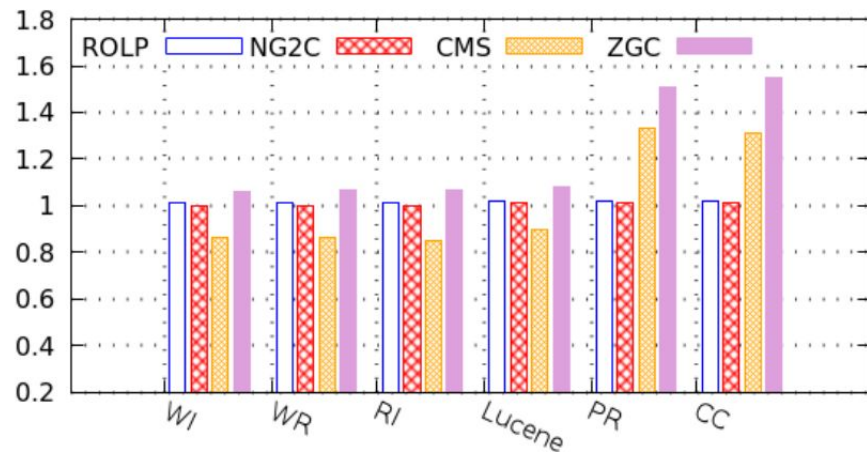
(f) GraphChi PR

Left and Lower is Better

Runtime Object Lifetime Profiler (throughput & memory)



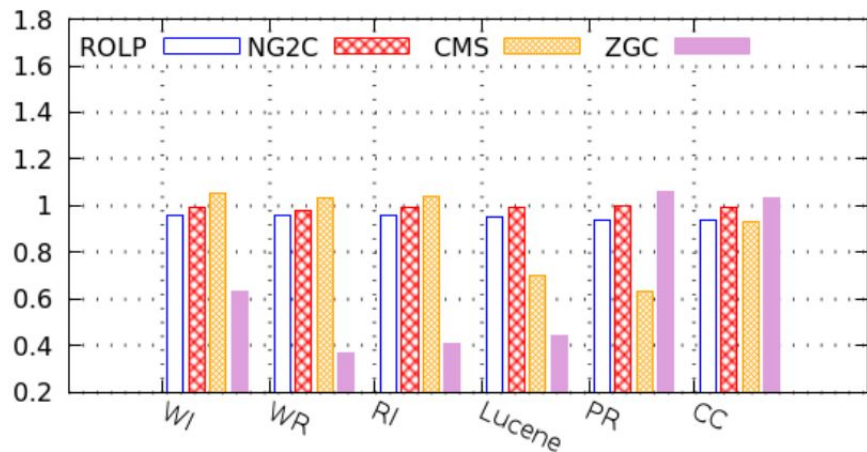
Throughput normalized to G1



Max Memory normalized to G1

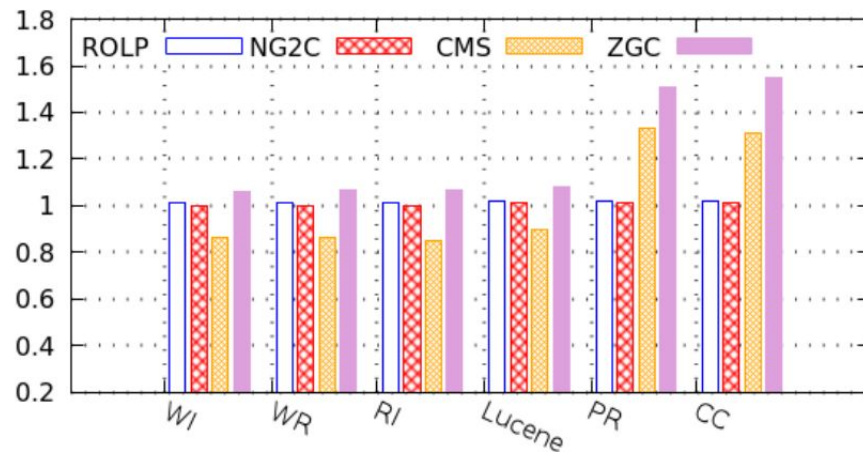
Lower is Better

Runtime Object Lifetime Profiler (throughput & memory)



Throughput normalized to G1

Up to 6%



Max Memory normalized to G1

Up to 2%

Lower is Better

Conclusions



- Big Data applications suffer from high long tail latencies
- Taking advantage of the proposed hypothesis leads to great reductions in pause times
 - More detailed results in the paper
- ROLP can significantly reduce application pauses with
 - Negligible throughput and memory overhead
 - No code access necessary
 - No programmer effort
- ROLP + NG2C is a JVM drop-in replacement

Conclusions



- Big Data applications suffer from high long tail latencies
- Taking advantage of the proposed hypothesis leads to great reductions in pause times
 - More detailed results in the paper
- ROLP can significantly reduce application pauses with
 - Negligible throughput and memory overhead
 - No code access necessary
 - No programmer effort
- ROLP + NG2C is a JVM drop-in replacement

Thank you for your time. Questions?

Rodrigo Bruno

email: rodrigo.bruno@inf.ethz.ch

webpage: rodrigo-bruno.github.io

code: github.com/rodrigo-bruno/rolp