

Ultra-deep technical exploration of JVM Garbage Collection (GC) Analysis in Containers

1. JVM vs Container Memory Model – A Deep Conflict

Default JVM Behavior:

The JVM calculates InitialHeapSize, MaxHeapSize, and GC region sizes based on:

- `/proc/meminfo`
- `sysconf(_SC_NPROCESSORS_CONF)`
- `Runtime.getRuntime().availableProcessors()`

Problem:

In a Docker container with limited memory (via cgroup v1/v2), the **container sees limited memory**, but the JVM (especially < Java 10 or Java 8 <u>191</u>) **ignores it** and assumes full host RAM is available.

Side-Effect:

JVM Behavior	Container Reality	Result
JVM allocates 4 GB heap	Container limit is 1 GB	OOMKill (SIGKILL)
JVM spawns 16 GC threads	CPU limit = 1 core	CPU throttling, GC pause inflation
Eden space too large	Low memory capacity	Minor GC frequency skyrockets

2. How JVM GC Ergonomics Are Skewed in Containers

JVM GC heuristics involve:

- CPU core count
- MaxHeapSize
- Allocation rate estimation
- GC pause latency estimation

But in containers:

- JVM overestimates CPU cores (/proc/cpuinfo still shows all vCPUs unless patched)
- Heuristics like -XX:MaxGCPauseMillis become invalid

Example: G1GC misbehavior

- G1 tries to **optimize pause time** using concurrent regions.
 - If Eden is 400 MB (based on 4 GB Xmx assumption), but available RAM is 1 GB, you'll get:
 - Eden GC every few seconds
 - Humongous object problems (objects > 50% of region size)
 - Huge Young GC times (copy pause)
-

3. GC Algorithms and Their Behavior in Containers

G1GC in Containers

Ideal For: Apps with medium memory footprint (256 MB – 2 GB), needing balanced pause/throughput

Tuning Essentials:

-XX:+UseG1GC

-XX:MaxGCPauseMillis=200

-XX:InitiatingHeapOccupancyPercent=30

-XX:+ParallelRefProcEnabled

-XX:+UseStringDeduplication

Problems Observed:

- Incorrect region sizing → High GC pause
 - G1 Humongous object GC every few seconds
 - Eden resizing jitter
-

✅ ZGC in Containers

Ideal For: Ultra-low latency services (e.g., financial APIs), where <10ms pause SLA

Needs: Java 11+, more CPU cores

Flags:

-XX:+UseZGC

-XX:ZUncommitDelay=5 # Aggressive memory release to OS

-XX:SoftMaxHeapSize=512m

Notes:

- Does not compact
 - Entire GC is concurrent
 - High CPU overhead — not ideal for constrained containers
-

✅ Shenandoah GC in Containers

Ideal For: RedHat-based JVMs (OpenJ9), low-latency needs

Flags:

-XX:+UseShenandoahGC

-XX:+ShenandoahUncommit

-XX:ShenandoahUncommitDelay=1000

Strength: Compacting GC with low STW pause

Weakness: JDK-specific tuning required; lower maturity in open source compared to G1/ZGC

❌ Parallel GC in Containers

Strength: High throughput GC for batch jobs

Weakness in Containers:

- STW GC only (no concurrency)
 - Long pause times if Eden is large
-

- Worst fit for REST or real-time apps
-

4. Full GC Log Interpretation in Containerized JVMs

Enable Logging:

-Xlog:gc*:file=/var/log/gc.log:time,uptime,level,tags

Log Sample (Java 11):

[5.015s][info][gc,start] GC(0) Pause Young (Normal) (G1 Evacuation Pause)

[5.015s][info][gc,heap] GC(0) Eden regions: 8->0(12)

[5.015s][info][gc,heap] GC(0) Survivor regions: 0->3(3)

[5.015s][info][gc,heap] GC(0) Old regions: 0->0

[5.015s][info][gc,phases] GC(0) Pause Young (Normal) 21.457ms

Interpretation:

- GC(0): First GC event
- Eden: 8 → 0 regions (Eden cleared)
- Survivor: Promoted 3 regions
- Old regions unchanged
- GC duration: 21.457ms

In container, observe:

- Humongous allocation → GC(0) Pause Full (G1 Humongous Allocation)
 - If multiple young GCs per second → over-allocation
-

5. Heap Sizing Math in Cgroup-Aware JVM

Assume container memory = 1 GiB (1024 Mi)

-XX:MaxRAMPercentage=75.0 → 768 Mi heap max

-XX:InitialRAMPercentage=50.0 → 512 Mi initial heap

Heap layout with G1GC (rough):

- Eden: ~256 Mi
- Survivor: ~32 Mi
- Old gen: ~480 Mi
- Metaspace: separate from heap (~64 Mi default limit)

Edge Case:

- Netty/Redis client allocating direct memory → needs:

-XX:MaxDirectMemorySize=256m

- Else, native memory exhaustion despite healthy heap

💣 6. Case Study – Memory Leak + GC Stall in Microservice

Context:

- Microservice in K8s (Java 11), 1Gi limit
- GC log: frequent minor GCs, occasional full GCs
- P99 latency > 2s
- Heap used: 200 Mi
- Container killed by OOM killer

Investigation:

1. Enabled -XX:NativeMemoryTracking=summary
2. Ran:

jcmd <pid> VM.native_memory summary

3. Found:

Total: reserved=1024MB, committed=980MB

- Java Heap: 400MB

- Class: 120MB

- Thread Stack: 64MB

- Arena + Code Cache + Direct Buffer = 400MB+

4. Heap dump: Not leaking on-heap.
5. DirectBufferLeak: Netty client had unclosed channels

Fix:

- Added shutdown hooks to release Netty buffer
- Added:

-XX:MaxDirectMemorySize=128m

Result:

- GC normalized
- No more OOMKills
- STW pauses reduced by 90%

7. Overlay GC with Load Test Timeline

Tools:

- JMeter → InfluxDB → Grafana
- GC logs parsed → Logstash → Elasticsearch

Overlay Example:

Time	GC Pause (ms)	Avg Resp Time	Errors
13:02:35	155 ms (Full GC)	982 ms	5xx
13:03:10	82 ms (Young GC)	441 ms	0
13:05:20	213 ms (Full GC)	1.2 sec	Spike

You can now **pinpoint latency spikes to GC events.**

8. Monitoring GC in Containers (Prometheus + Grafana)

Metrics Exposed (Micrometer):

- jvm_gc_pause_seconds_sum
- jvm_gc_memory_allocated_bytes

- `jvm_memory_used_bytes{area="heap"}`

JMX Exporter alternative:

`lowercaseOutputName: true`

rules:

- pattern: "java.lang:type=GarbageCollector,name=(.*)"
name: jvm_gc_collection_seconds
type: summary

Grafana Dashboards:

- GC pause per collector
- GC count (Young, Full)
- Heap used vs committed
- Container RSS vs JVM heap → detect native leaks

✅ Always Check:

- JVM sees cgroup values → check `jcmd <pid> VM.info`
- GC log enabled and collected to external mount
- $Xmx + \text{direct memory} + \text{metaspace} \leq \text{cgroup limit}$
- GC type matches workload
- Native memory isn't leaking (NMT, jcmd)
- Logs + metrics overlay confirm GC impact

TIPS: JVM GC in Containers

TIP 1: Heap Pressure Monitoring ≠ Container Memory Pressure

Even if JVM's heap is **well within Xmx**, your container might get OOMKilled.

Why?

- Metaspace, thread stacks, direct memory, NIO buffers, and JIT code cache live **outside the Java heap**
- Linux cgroups count **total RSS**, not just heap
- JVM GC won't help if **native memory** is leaking

Action:

- Use `-XX:NativeMemoryTracking=summary`
- Cap native memory:
- `-XX:MaxDirectMemorySize=128m`
- `-XX:MaxMetaspaceSize=128m`

TIP 2: PreTouch and Memory Fragmentation

Containers under memory pressure can fragment physical memory. JVMs use `mmap()` and `madvise()` to request pages.

Symptom:

- Slow GC starts
- Major page faults
- High STW delay before actual GC begins

Fix:

`-XX:+AlwaysPreTouch`

 This **forces physical page allocation** upfront, reducing fragmentation-related GC hiccups.

TIP 3: GC Thread Pools Must Respect CPU Quotas

When CPU quota is 1, JVM might still try to use default GC threads (based on `os::active_processor_count()` from `/proc/cpuinfo`)

Action:

`-XX:ParallelGCThreads=1`

`-XX:ConcGCThreads=1`

`-XX:ActiveProcessorCount=1`

Use `jcmd <pid> VM.flags` to verify. Over-provisioned GC threads will **steal CPU from the app** in low-vCPU containers.

TIP 4: GC Log Timestamp Normalization for Correlation

GC logs print time since JVM startup:

6.123s: [GC pause (G1 Evacuation Pause) ...]

But logs like application logs, Splunk, or InfluxDB dashboards use **epoch timestamps**.

Action:

Use startup timestamp + GC log to convert relative times to absolute:

`JVM_START_EPOCH + gc_time_offset = Real timestamp`

Useful when overlaying with:

- JMeter run logs
 - Kubernetes liveness probe failures
 - Service mesh (Istio/Envoy) telemetry
-

TIP 5: Avoid System.gc() STW in containers

When System.gc() is called, it triggers a **Full GC**, which is especially disastrous under limited CPU/memory.

Action:

-XX:+DisableExplicitGC

Alternative:

-XX:+ExplicitGCInvokesConcurrent

This forces a **concurrent Full GC** in G1 instead of STW.

TIP 6: Use jstat to Track GC Behavior in Real Time (Even in Prod Containers)

docker exec -it <container> jstat -gc <pid> 1000

Fields of interest:

- S0U, S1U (Survivor space usage)
- EU (Eden space used)
- OC (Old gen capacity)
- YGC, FGC (young/full GC counts)
- YGCT, FGCT (time spent in GC)

Look For:

- EU frequently reaching EC → Eden overflow
 - FGC climbing → too frequent Full GCs
 - GC time >10% of app time → throughput loss
-

TIP 7: Configure and Export GC Logs to External Volumes in Docker

GC logs written inside container are lost unless **persisted or streamed**.

Dockerfile + CMD:

ENV JAVA_OPTS="-Xlog:gc*:file=/logs/gc.log:time,tags"

VOLUME /logs

CMD ["java", "\$JAVA_OPTS", "-jar", "app.jar"]

Mount with:

docker run -v /host/logs:/logs ...

Then tail it or ship via FluentBit/Logstash to Splunk/Elasticsearch.

TIP 8: Humongous Allocations in G1GC = GC Nightmare

In G1GC:

- Any object > 50% of region size → Humongous Allocation
- Stored in contiguous old-gen regions
- Not cleaned by minor GC
- Triggers frequent Full GCs

Action:

-XX:G1HeapRegionSize=2m

OR:

- Break down large object graphs (esp. strings, JSON, lists)
- Watch for:

[Full GC (G1 Humongous Allocation)]

In GC logs

TIP 9: Restart Avoidance by Tracking RSS at Runtime

GC analysis often misses the fact that **RSS grows over time** even if heap usage is flat.

Track:

ps -o pid,rss,cmd | grep java

Correlate JVM heap (from JMX or Prometheus) with RSS.

If:

- JVM heap = 300 MB
- RSS = 800 MB → Memory leak outside heap

Likely Culprits:

- Netty ByteBuf leaks
- Unsafe memory (sun.misc.Unsafe)
- Native libraries

TIP 10: G1GC Tuning Heuristics in Containers

In G1GC, most GC issues stem from:

- **Incorrect region sizing**
- **Improper pause prediction**
- **Eden sizing instability**

Tuning Combo:

-XX:+UseG1GC

-XX:+ParallelRefProcEnabled

-XX:MaxGCPauseMillis=100

-XX:InitiatingHeapOccupancyPercent=30

-XX:G1ReservePercent=10

-XX:G1HeapRegionSize=2m

-XX:+AlwaysPreTouch

This balances:

- Fast startup (pre-touch)
- Pause control (<100ms)
- Heap occupancy
- Safety buffer

TIP 11: Capture Heap + Native Delta Together

Many tools show either heap or native memory.

Ideal Script for Heap Leak + Native Leak:

```
jmap -dump:format=b,file=/tmp/heap.hprof <pid>
```

```
jcmd <pid> VM.native_memory summary > /tmp/nmt.log
```

Then:

- Use Eclipse MAT for heap dump
 - Parse nmt.log to find arena/class/native growth
 - Combine to trace total memory leak
-

TIP 12: Simulate GC Pressure for Canary Testing

Before releasing:

- Deploy app in K8s with low memory/cpu
- Run synthetic traffic (JMeter, Locust)
- Observe:
 - GC frequency
 - GC pause time
 - Heap promotion trend
 - NMT growth

Canary Alerting:

- GC pause > 200ms → warn
 - Full GC count > 5/hr → warn
 - Native memory growth > 20% over baseline → critical
-

TIP 13: Instrument GC Pause Breakdown with Phases

Use:

`-Xlog:gc*,gc+phases=debug`

Will print:

GC(11) Pause Young (Normal) (G1 Evacuation Pause) 50.123ms

Pre Evacuate Collection Set: 0.1ms

Evacuate Collection Set: 25.2ms

Post Evacuate Collection Set: 3.1ms

Other: 21.7ms

Use this to:

- Pinpoint evacuation bottlenecks
 - Correlate GC phase with CPU stalls (e.g., poor GC thread CPU share)
-

TIP 14: Split Critical and Non-Critical Pods for GC Experimentation

In K8s:

- Use **different node pools or namespaces**
- Run “experimental GC tuned” versions side-by-side
- Compare:
 - GC pause time
 - 5xx error rates
 - JMeter throughput metrics

Use Istio traffic mirroring or blue-green deployment to split traffic.

TIP 15: Full GC Cost Estimation

To estimate GC stall budget:

Formula:

Pause time \approx Object Count \times Object Size \times Copy Time per MB

E.g.:

- 100k live objects of 1 KB
- Copy time = 2 ms/MB

Then:

- GC pause = $(100k \times 1 \text{ KB}) / 1024 \times 2 \text{ ms} = \sim 200 \text{ ms}$

Predict STW pause impact before prod rollout.

FINAL TECHNICAL TAKEAWAYS

Default JVM ≠ Container-Aware

→ Always use:

-XX:+UseContainerSupport

-XX:MaxRAMPercentage=75.0

-XX:ActiveProcessorCount=<CPU limit>

GC Thrashing = Eden too large / Humongous allocations

→ Tune G1GC:

-XX:+UseG1GC

-XX:MaxGCPauseMillis=100

-XX:G1HeapRegionSize=2m

Heap ≠ Total Memory Usage

→ Native leaks/OOMKill despite free heap?

Use:

-XX:NativeMemoryTracking=summary

-XX:MaxDirectMemorySize=<limit>

jcmd <pid> VM.native_memory summary

Overlay GC logs with JMeter/test timeline

→ Align latency spikes with GC pause

→ Use -Xlog:gc*:file=... + dashboards (Grafana/Splunk)

STW Pause Cause? Break it Down!

→ Use:

-Xlog:gc*,gc+phases=debug

→ Analyze evacuation time vs GC phase CPU stalls

Pre-touch to prevent page faults:

-XX:+AlwaysPreTouch

Mount GC logs in Docker

→ Use volume mounts:

-v /host/logs:/logs +

-Xlog:gc*:file=/logs/gc.log

Monitor These GC Metrics in Grafana:

- jvm_gc_pause_seconds
 - jvm_memory_used_bytes
 - jvm_gc_memory_allocated_bytes
 - jvm_gc_collection_seconds_count
-

Test GC under stress BEFORE prod

→ Canary deploy → Load → Watch GC logs + heap + RSS + p99

Checklist for GC in Containers

- Container-aware heap + CPU sizing
 - GC algorithm tuned for load
 - GC pause below SLA (<200ms)
 - Heap + native memory budgeted
 - GC logs + metrics integrated into observability stack
-