# Observability beyond the three pillars

## Continuous Profiling with Alibaba Dragonwell

Sanhong Li | Java Champion, JVM Architect | Alibaba Cloud

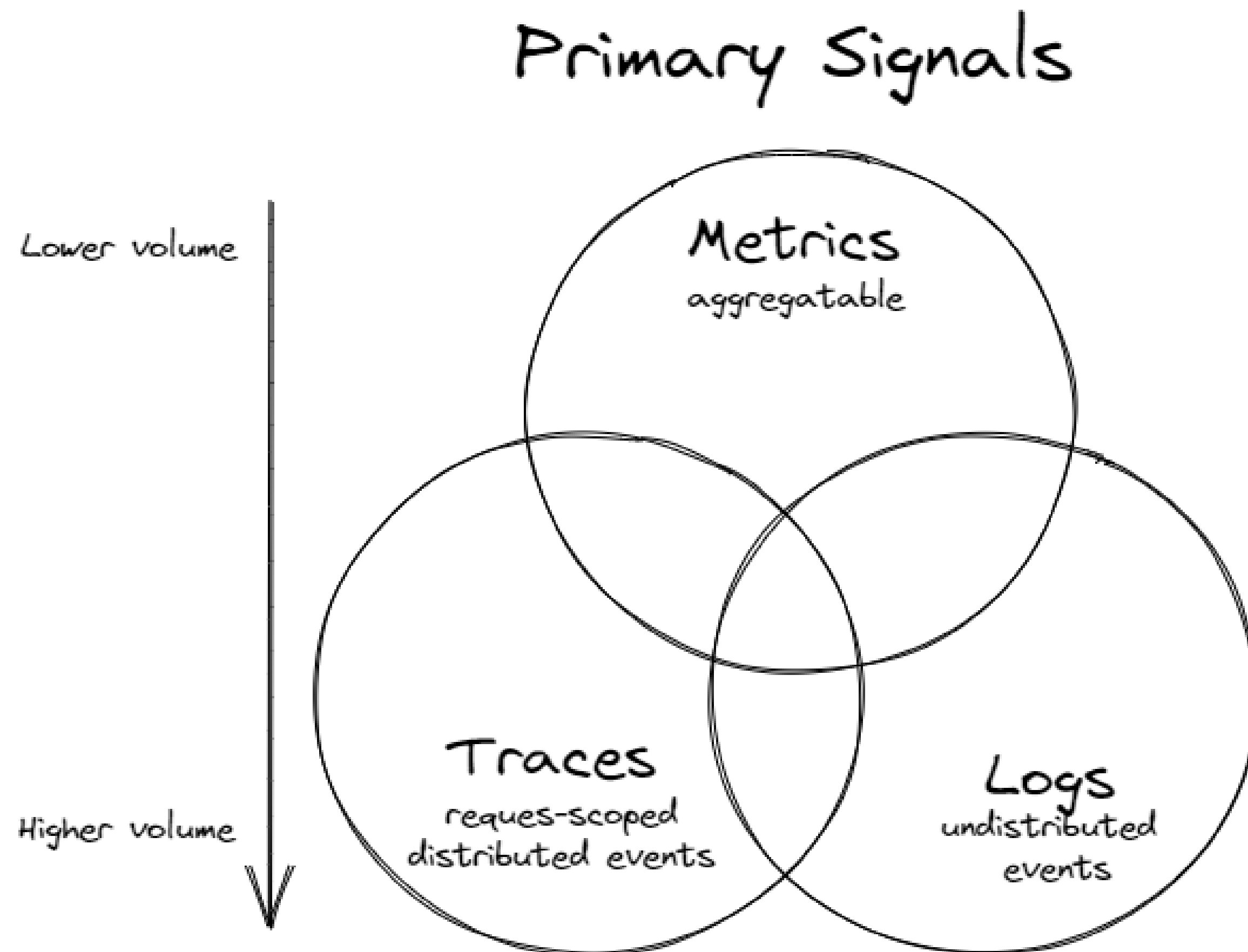# Content

ALIBABA CLOUD
INTELLIGENCE GROUP

# Observability to Software Systems

- "observability" was originally coined in 1960, used to describe mathematical control systems.

- Applied to modern software systems
  - **Understand** the inner workings of your application
  - **Understand** any system state your application may have gotten
  - **Understand** the inner workings and system state **solely** by observing and interrogating with external tools
  - **Understand** the internal state without shipping any new custom code to handle it

Observability Engineering, Charity Majors, Liz Fong-Jones, George Miranda 2022

# (Traditional) Three Observability Pillars

Primary Signals

Lower volume

Higher volume

Metrics
aggregatable

Traces
reques-scoped
distributed events

Logs
undistributed
events

**Metrics**

numeric representation of data measured over intervals of time(e.g via Java MXBean API)

**Traces**

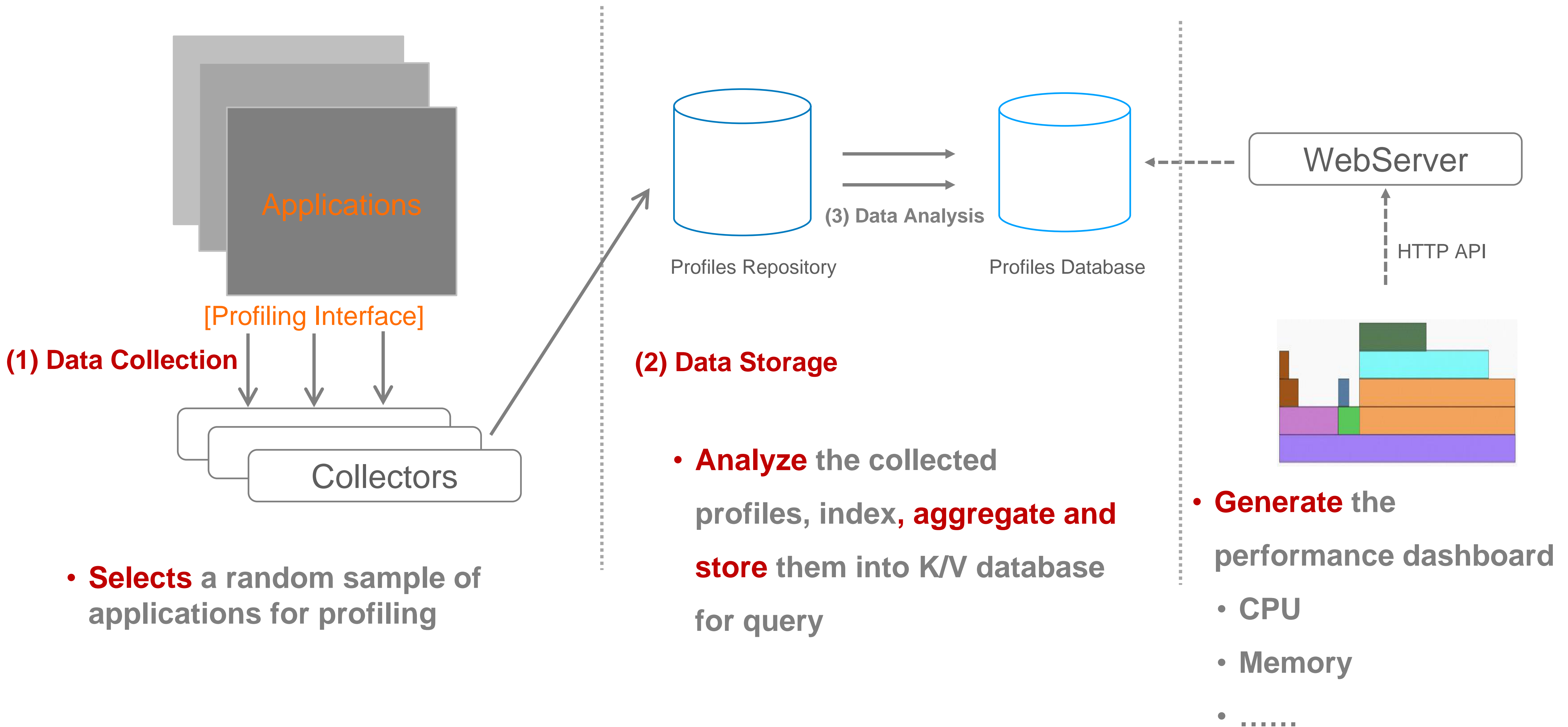representation of a series of causally related distributed events(e.g via JVMTi agent)

**Logs**

immutable, timestamped record of discrete events that happened over time.

(e.g gc log is a typical example)

https://github.com/cncf/tag-observability/blob/main/whitepaper.md#observability-signals

# Continuous Profiling: 4th aspect of Observability

- Google pioneered the <span style="color:red">continuous profiling</span> concept in its own data centers

  - "Google-Wide Profiling: A Continuous Profiling Infrastructure for Data Centers", research paper published by Google in 2010

- Continuous profiling (<span style="color:red">constantly</span> monitors an application's performance <span style="color:red">in real time</span>)

  - <span style="color:red">Executing in a production environment</span>(no need to develop accurate predictive load tests or benchmarks for the production )

  - <span style="color:red">Sampling</span>(low overhead)

https://storage.googleapis.com/pub-tools-public-publication-data/pdf/36575.pdf

# Continuous Profiling: Conceptual Architecture

Applications

[Profiling Interface]

**(1) Data Collection**

Collectors

• **Selects** a random sample of applications for profiling

Profiles Repository

(3) Data Analysis

Profiles Database

**(2) Data Storage**

• **Analyze** the collected profiles, index, **aggregate and store** them into K/V database for query

WebServer

HTTP API

• **Generate** the performance dashboard
  • CPU
  • Memory
  • ......

# Content

# Java Flight Recorder(JFR) History

**2008**

Oracle acquired BEA

**2012**

JRockit Flight Recorder was rebranded

Java Flight Recorder, released in JDK7

(as proprietary commercial offering)

**2020**

Alibaba is working with the community

(Red Hat, Azul and Datadog) to

contribute back to OpenJDK8u

JRockit Flight Recorder
BEA Systems

Oracle acquired Sun

**2010**

Open sourceed in

OpenJDK 11 by Oracle

**2018**

# JFR Workflow Overview

ALIBABA CLOUD
INTELLIGENCE GROUP

profile.jfc

JIT

GC

Runtime

JFR

JFR
Events

JFR dump
(*.jfr)

1. Java Mission Control

2. JFR tool shipped by
JDK by default
- jfr summary
- jfr print

3. Java API

jdk.jfr.consumer

1. Command line via VM
arguments
- -XX:+FlightRecorder
- -XX:StartFlightRecording

2. JCMD comand
- JFR.start
- JFR.stop

3. Java API

jdk.jfr

(since JDK9)

# JFR Events

**Instant Event**
- occur immediately

**Examples from $JAVA_HOME/lib/jfr/default.jfc**

```xml
<event name="jdk.ThreadStart">
  <setting name="enabled">true</setting>
  <setting name="stackTrace">true</setting>
</event>
```

**Duration Event**
- have a start and end

**Timed Event**
- Like Duration Event, but with threshold set

```xml
<event name="jdk.JavaMonitorEnter">
  <setting name="enabled">true</setting>
  <setting name="stackTrace">true</setting>
  <setting name="threshold" control="synchronization-threshold">20 ms</setting>
</event>
```

**Sample Events**
- logged at a regular interval

```xml
<event name="jdk.ExecutionSample">
  <setting name="enabled" control="method-sampling-enabled">true</setting>
  <setting name="period" control="method-sampling-java-interval">20 ms</setting>
</event>
```

# Content

# ARMS Continuous Profiler - Introduction

❑ **Alibaba Dragonwell** – downstream of OpenJDK,
free LTS by Alibaba Cloud

❑ **ARMS** - Application Real-Time Monitoring Service
(ARMS) is an application performance
management (APM) service, written in Java,
running on top of Alibaba Dragonwell

❑ **Continuous Profiler**, based on JFR and async-
profiler technology, part of ARMS

   ✓ Collecting profiling data via Java agent

   ✓ Major features

      ✓ CPU/Allocation/Wall clock profiling

      ✓ Integration with Tracing



CPU    Allocation    Wall clock    Tracing Integration

**Java Agent**

*Continuous Profiler(JFR/Async Profiler based)*

*Application Realtime Monitoring Service(ARMS)*

*Alibaba Dragonwell*

https://www.alibabacloud.com/product/arms

# ARMS Continuous Profiler – Overview

**ALIBABA CLOUD**
**INTELLIGENCE** GROUP
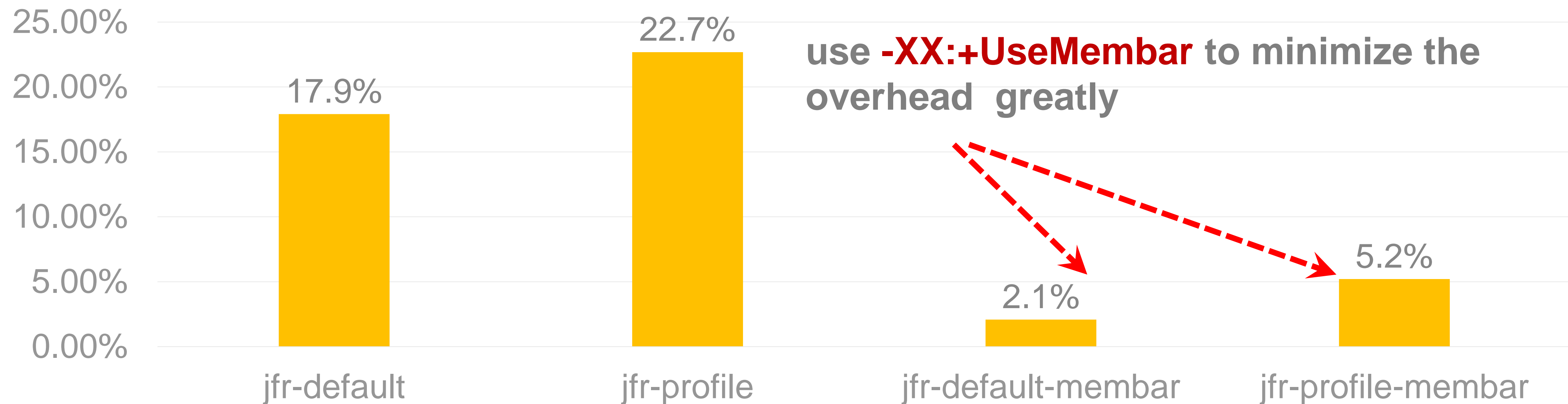
# JFR Overhead Assessment

## SPECjbb 2015 - max-jOPS



- **OS:**  Linux version 3.10.0-1160.80.1.el7.x86_64

- **CPU:**  vCore 24, x86_64, Intel(R) Xeon(R) Platinum 8369B CPU @ 2.70GHz

- **Memory:**  24G

- **JDK:**  OpenJDK **1.8.0_362 64-Bit** Server VM (build 25.362-b08, mixed mode)

- **Flags:**  -XX:+UseConcMarkSweepGC -Xmx10g -Xmn5g -XX:MaxMetaspaceSize=512m -Xloggc:gc.log -XX:+PrintGCDetails -XX:+PrintGCDateStamps
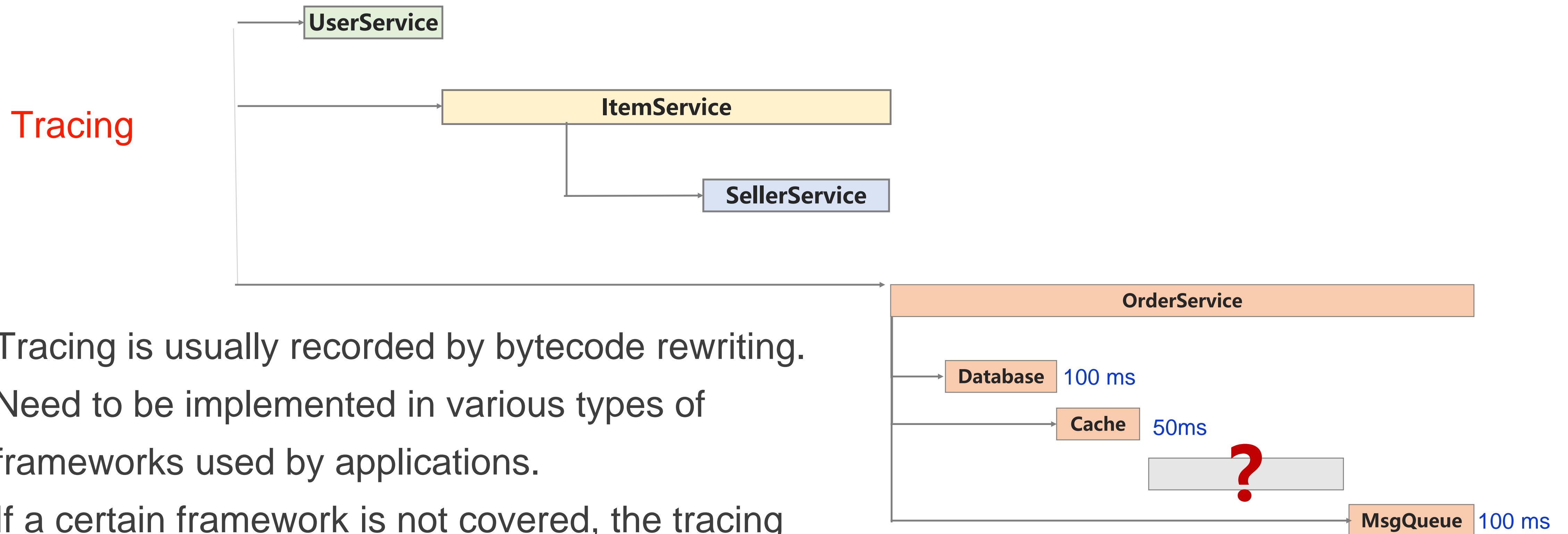
# JFR Overhead Assessment

## SPECjbb 2015 - max-jOPS



use **-XX:+UseMembar** to minimize the overhead greatly

- The cost was introduced by the call to os::serialize_thread_states (which uses mprotect to force a pseudo-memory barrier) by JfrThreadSampler thread. Notes: this sync mechanism has the performance issue in scalability!
- **-XX:+UseMembar** uses a direct memory fence operation, which is more cheap to get the state of java thread
- More detail info:
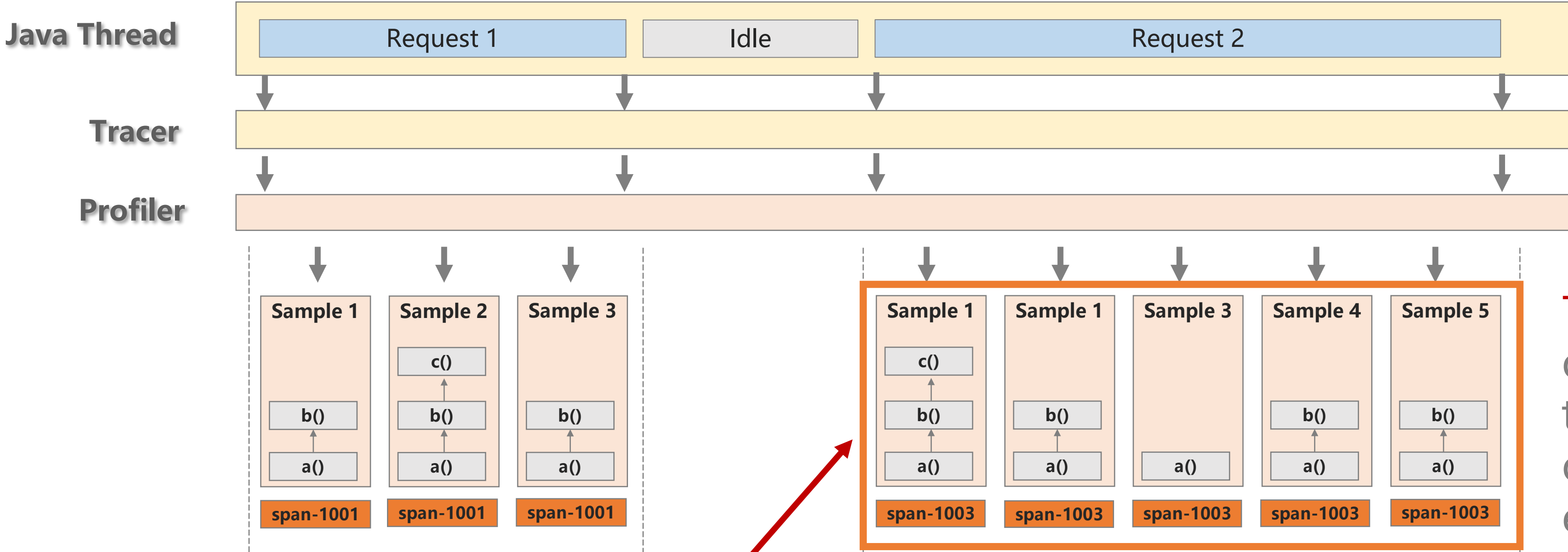  - https://bugs.openjdk.org/browse/JDK-8187812
  - https://bugs.openjdk.org/browse/JDK-8276309

# Contextualized JFR – Corelating with Tracing

Tracing

UserService

ItemService

SellerService

OrderService

Database  100 ms

Cache  50ms

**?**

MsgQueue  100 ms

- Tracing is usually recorded by bytecode rewriting.

- Need to be implemented in various types of frameworks used by applications.

- If a certain framework is not covered, the tracing information will be missed, resulting in monitoring blind spots.
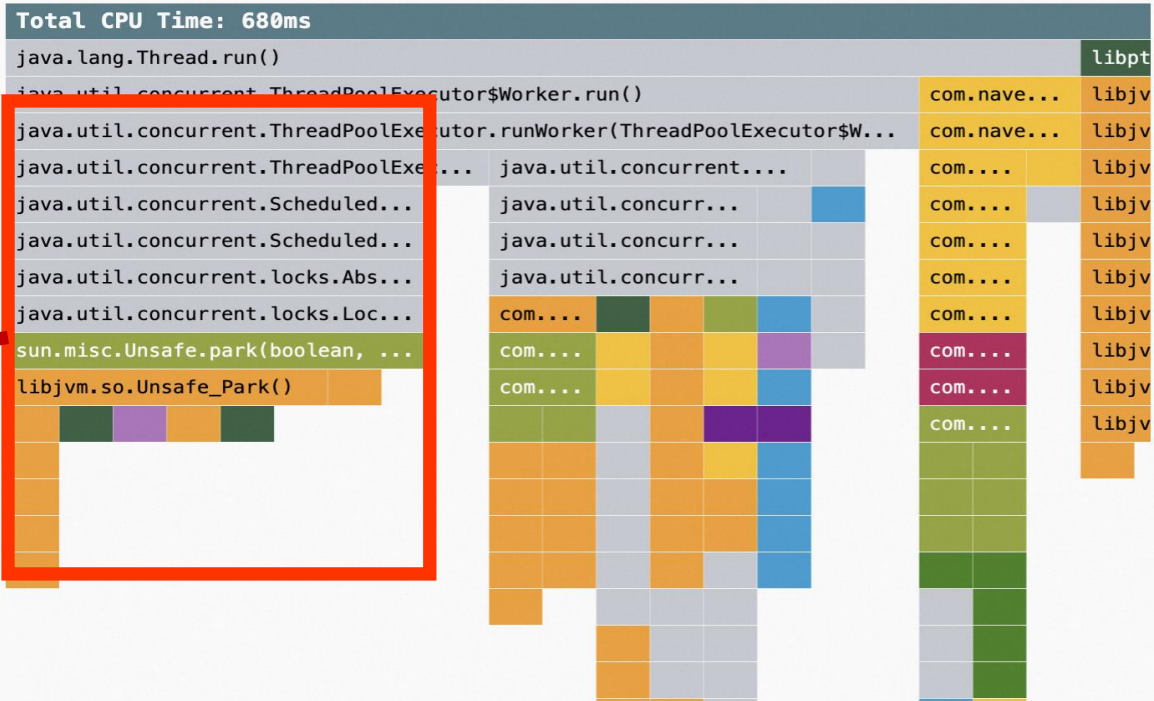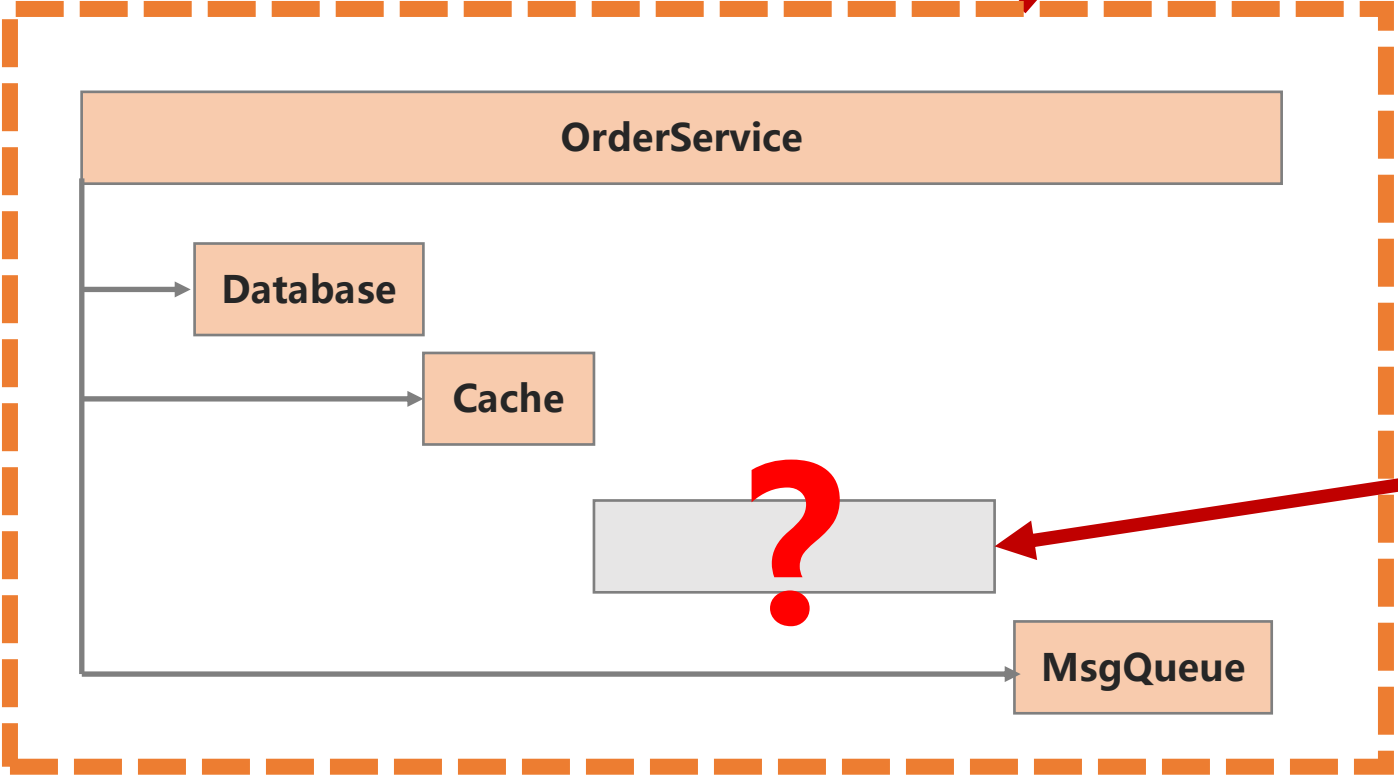
- The challenge is: even if we have collected JFR events at those blind spots, we don't know the their relationship - JFR currently lacks the ability to associate context to JFR events.
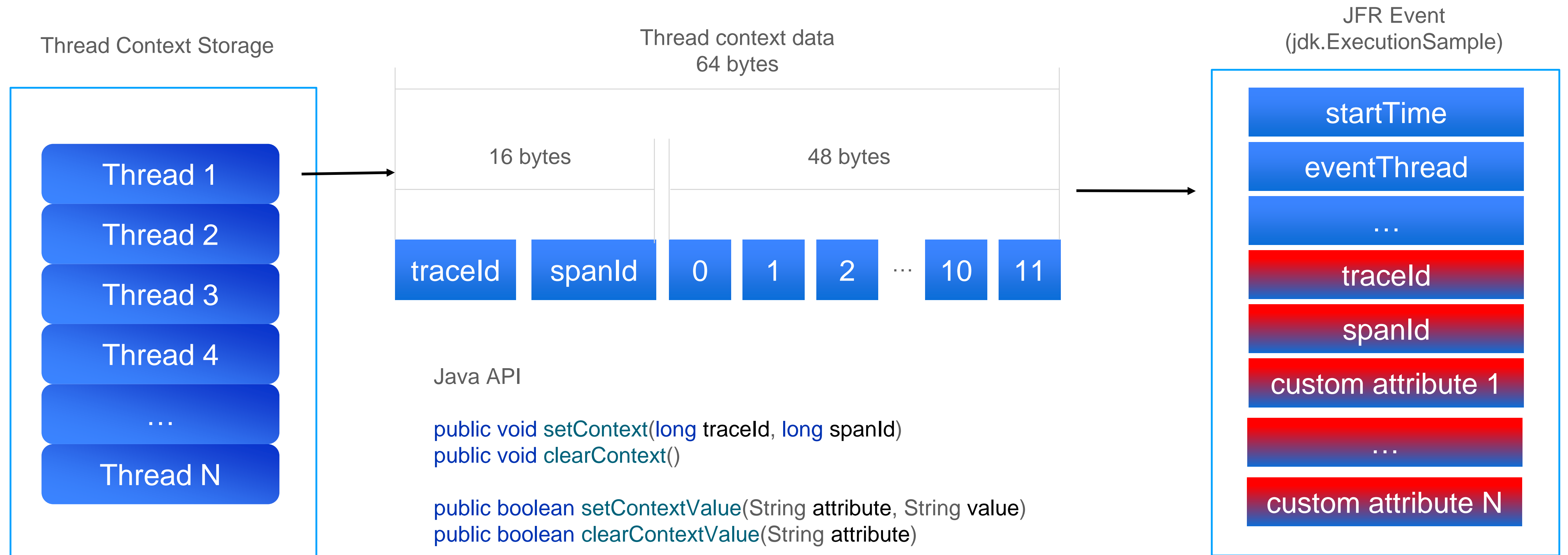
# Contextualized JFR – Corelating with Tracing(2)

**Trace** is composed of one or more **spans**. Span in the trace represents one microservice in the execution path.

# Contextualized JFR – Corelating with Tracing(3)

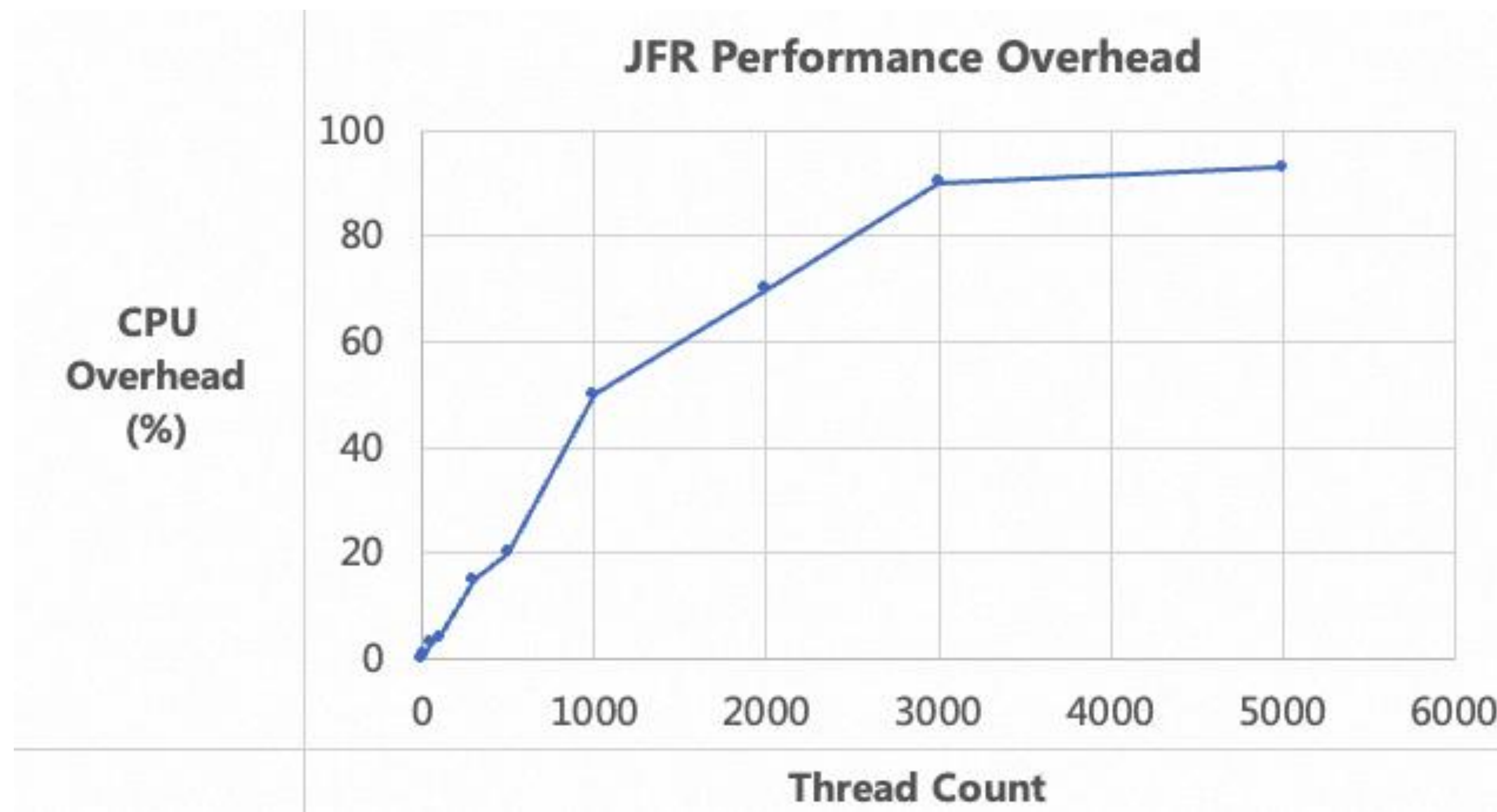- An context implementation example(based on ideas of aync-profiler & Datadog implementation)

Thread Context Storage

Thread context data
64 bytes

JFR Event
(jdk.ExecutionSample)

16 bytes | 48 bytes

Thread 1
Thread 2
Thread 3
Thread 4
…
Thread N

traceId | spanId | 0 | 1 | 2 | … | 10 | 11

startTime
eventThread
…
traceId
spanId
custom attribute 1
…
custom attribute N

Java API

```java
public void setContext(long traceId, long spanId)
public void clearContext()

public boolean setContextValue(String attribute, String value)
public boolean clearContextValue(String attribute)
```

More reference implementation info:
- https://github.com/async-profiler/async-profiler
- https://github.com/DataDog/java-profiler

# Lessons Learned

**#1** For JDK 8, if the number of Java threads is large (for example, more than 500),  CPU overhead may be expensive and is positively related to the number of  threads.



It is highly recommended to add **the -XX:+UseMembar** flag to avoid this problem.
It has been turned on by default since JDK 10: https://bugs.openjdk.org/browse/JDK-8187812

# Lessons Learned(2)

**#2** For JDK 8 and JDK 11, the amount of events for memory allocation may be large. It is not recommended to enable them for applications with fast memory allocation.
- https://bugs.openjdk.org/browse/JDK-8257602

**#3** Before JDK 11.0.7, the OldObjectSample event may create unexpected amount of checkpoint data, cause the JFR file to be very large, and it is not recommended to enable it.
- https://bugs.openjdk.org/browse/JDK-8225797

**#4** Before JDK 11.0.12, the OldObjectSample event is expensive, not suitable for production.
- https://bugs.openjdk.org/browse/JDK-8225797

# Case Study - Object Allocation

Metrics captured from real workloads



cpu%
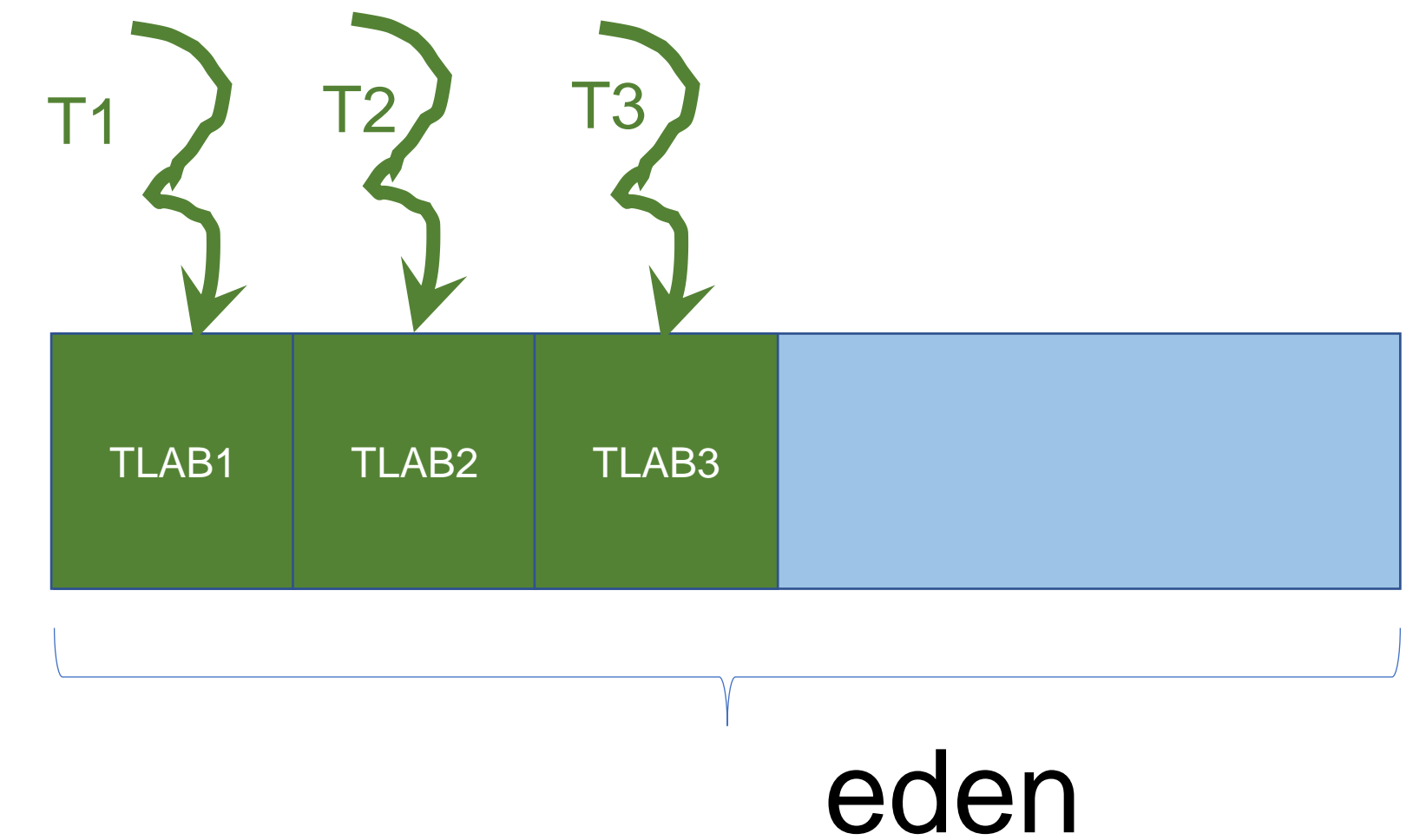
ygc count /10 min

time

GC Spikes

# Case Study - Object Allocation

```
[GC (Allocation Failure) 2018-05-17T21:09:04.953+0800: 16.569: [ParNew: 921899K->53900K(961216K), 0.0412584 secs]
[GC (Allocation Failure) 2018-05-17T21:09:09.686+0800: 21.302: [ParNew: 927756K->61952K(961216K), 0.0493610 secs]
[GC (Allocation Failure) 2018-05-17T21:09:11.642+0800: 23.258: [ParNew: 935808K->61153K(961216K), 0.1264167 secs]
[GC (Allocation Failure) 2018-05-17T21:09:16.322+0800: 27.938: [ParNew: 935009K->74003K(961216K), 0.0779854 secs]
[GC (Allocation Failure) 2018-05-17T21:09:28.447+0800: 40.063: [ParNew: 947859K->66919K(961216K), 0.0559919 secs]
[GC (Allocation Failure) 2018-05-17T21:09:34.607+0800: 46.223: [ParNew: 926011K->87230K(961216K), 0.0436882 secs]
[GC (Allocation Failure) 2018-05-17T21:09:39.122+0800: 50.738: [ParNew: 961086K->87360K(961216K), 0.3830953 secs]
[GC (Allocation Failure) 2018-05-17T21:09:41.372+0800: 52.988: [ParNew: 961216K->87360K(961216K), 0.3958484 secs]
[GC (Allocation Failure) 2018-05-17T21:09:52.437+0800: 64.053: [ParNew: 961216K->87360K(961216K), 0.0797925 secs]
[GC (Allocation Failure) 2018-05-17T21:10:27.194+0800: 98.810: [ParNew: 961216K->87360K(961216K), 0.2047217 secs]
```

GC log cannot tell you what allocated the most objects

# Case Study - Object Allocation

- Support TLAB allocation statistics by

  - EventObjectAllocation**Outside**TLAB

  - EventObjectAllocation**InNew**TLAB

- Used to find out where the allocation pressure is



eden

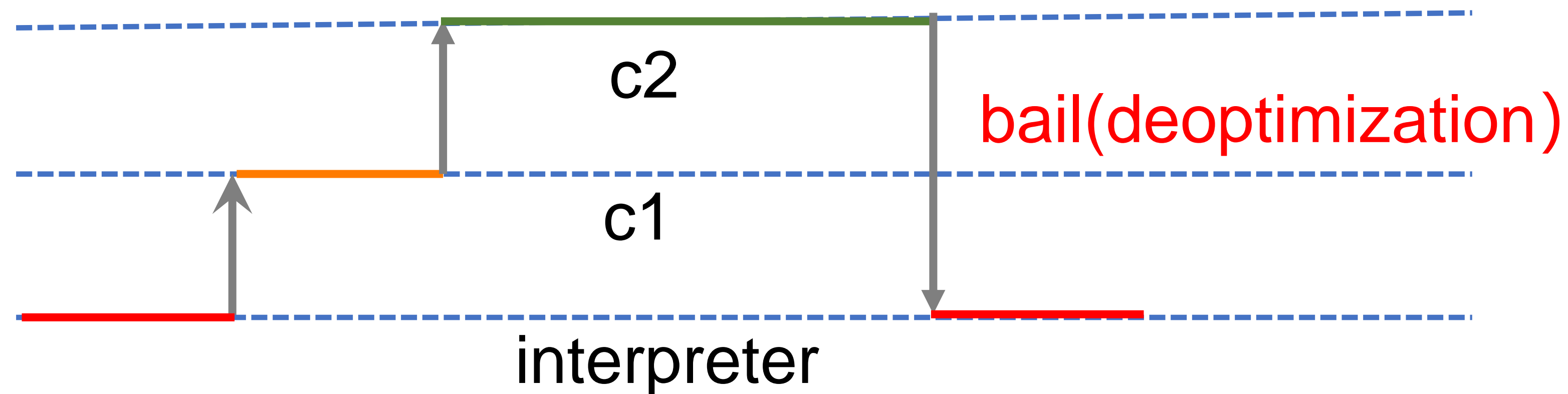| Thread | | Count | Average TLAB Allocation | Average Allocation Outside TLA... | Est. TLAB Allocation | Total Allocation Outside TLABs |
|---|---|---|---|---|---|---|
| EagleEye-StatLogController-writer-thread-1 | | 24,621 | 185 B | 522 B | 188 MiB | 3.42 MiB |
| AsyncAppender-Worker-createParamsAppender-async | | 2,933 | 17.2 KiB | 28 KiB | 191 MiB | 65 MiB |
| pool-55-thread-1 | | 1,783 | 672 B | 909 B | 1.37 MiB | 1.4 MiB |

# Object Allocation JFR Demo

1.  Run server: make object allocations on request arrival.

2.  Use jcmd to enable object allocation event tracking.

3.  Launch client: send requests to sever.

4.  Use jmcx to generate folded stacks from JFR dump file.

5.  Use flamegraph.pl to generate flame graph for object allocation.

# Case Study II - Deoptimization

## Basics concept of Just-in-Time compiler

- Mix mode execution

- Profile Guided Optimization

  - Optimization decision are made dynamically

  - Bail to interpreter if the assumption is wrong



c2

bail(deoptimization)

c1

interpreter

Deoptimization is very expensive if speculation is wrong:
fall back to interpreter and wait for re-compilation

# Case Study II - Deoptimization

- Unstable if case

```
if (condition) {
    // hot path                    A
    // always hit in profile
    ...
} else {
    // cold path                   B
    // profile data assume it's never taken
}
```

trap to runtime system: uncommon branch is taken

# Deoptimization JFR Demo

1.  Run JMH benchmark: an example for Unstable-if deopt

    ✓ Enable JFR setting via –jvmArgs(JMH parameter)

2.  Use jmcx to generate folded stacks from JFR dump file.

3.  Use flamegraph.pl to generate flame graph for object allocation.

ALIBABA CLOUD
INTELLIGENCE GROUP

# Content

# Round-up

1. The Definition of Observability and basics of Continuous Profiling(4th Pillar of Observability)

2. Basics of JFR(Observability tools for JVM applications)

   ✓ JFR Workflow(How to use it) and JFR Events(Understand what they are used for)

3. Alibaba Practice: ARMS Continuous Profiler

4. Two Case Studies : Object Allocations/Deoptimization