演讲主题

陆家靖

收钱吧框架工具团队负责人、 SkyWalking PMC Member

"基于SkyWalking Agent的
持续性能剖析与交互式诊断"

# 目录
## CONTENTS

**01**

# 持续性能剖析
# Continuous Profiling

# GOOGLE-WIDE PROFILING:
# A CONTINUOUS PROFILING
# INFRASTRUCTURE FOR DATA CENTERS

GOOGLE-WIDE PROFILING (GWP), A CONTINUOUS PROFILING INFRASTRUCTURE FOR

DATA CENTERS, PROVIDES PERFORMANCE INSIGHTS FOR CLOUD APPLICATIONS. WITH

NEGLIGIBLE OVERHEAD, GWP PROVIDES STABLE, ACCURATE PROFILES AND A

DATACENTER-SCALE TOOL FOR TRADITIONAL PERFORMANCE ANALYSES. FURTHERMORE,

GWP INTRODUCES NOVEL APPLICATIONS OF ITS PROFILES, SUCH AS APPLICATION-

PLATFORM AFFINITY MEASUREMENTS AND IDENTIFICATION OF PLATFORM-SPECIFIC,

MICROARCHITECTURAL PECULIARITIES.

FIG. GWP published by Google in 2010：low overhead, stable, accurate, scalable
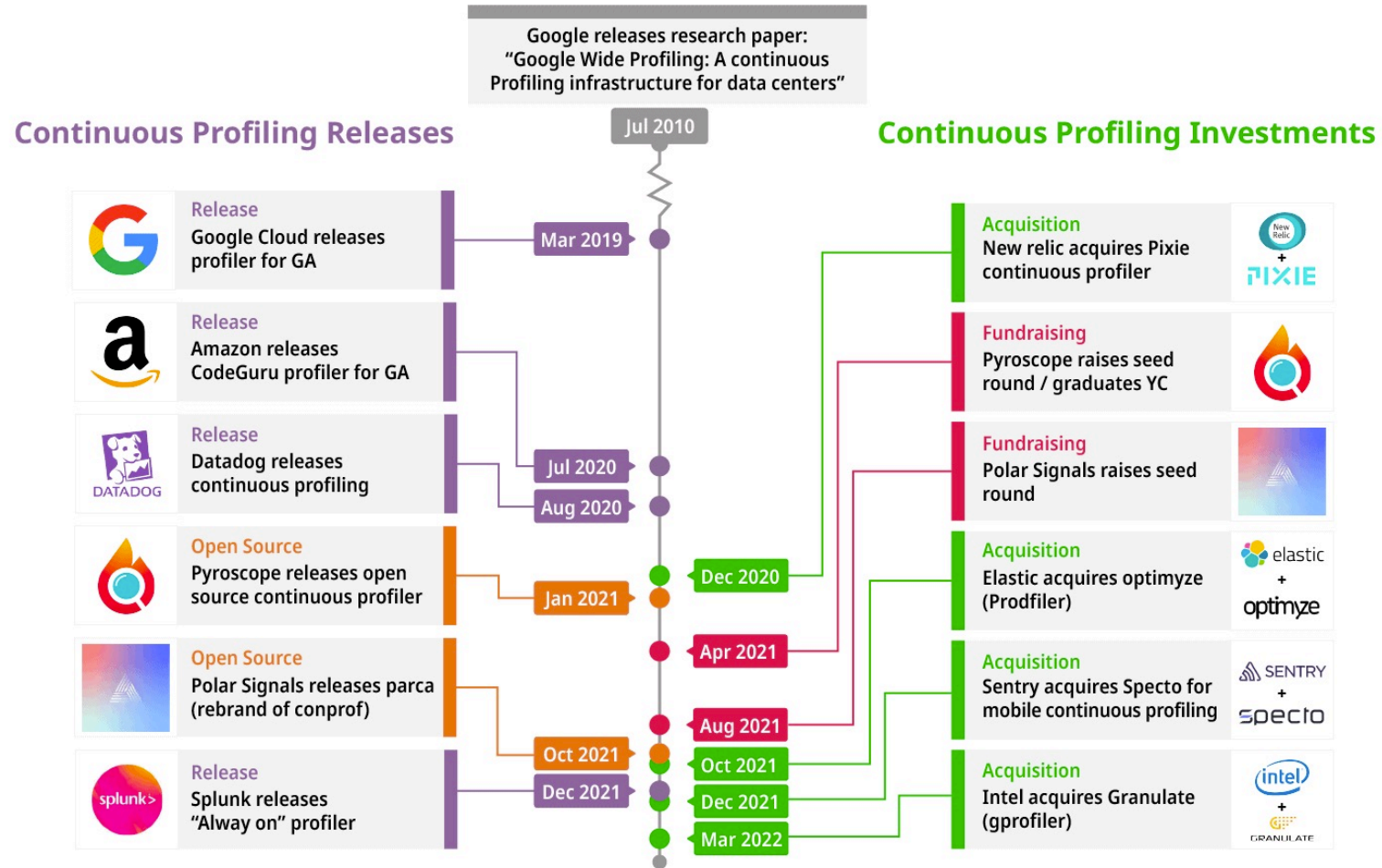
# Continuous Profiling的发展史

FIG. Since GWP，many major vendors have joined "Continuous Profiling"：
Pyroscope is an open-source solution，acquired by Grafana Lab on 2023-03-15
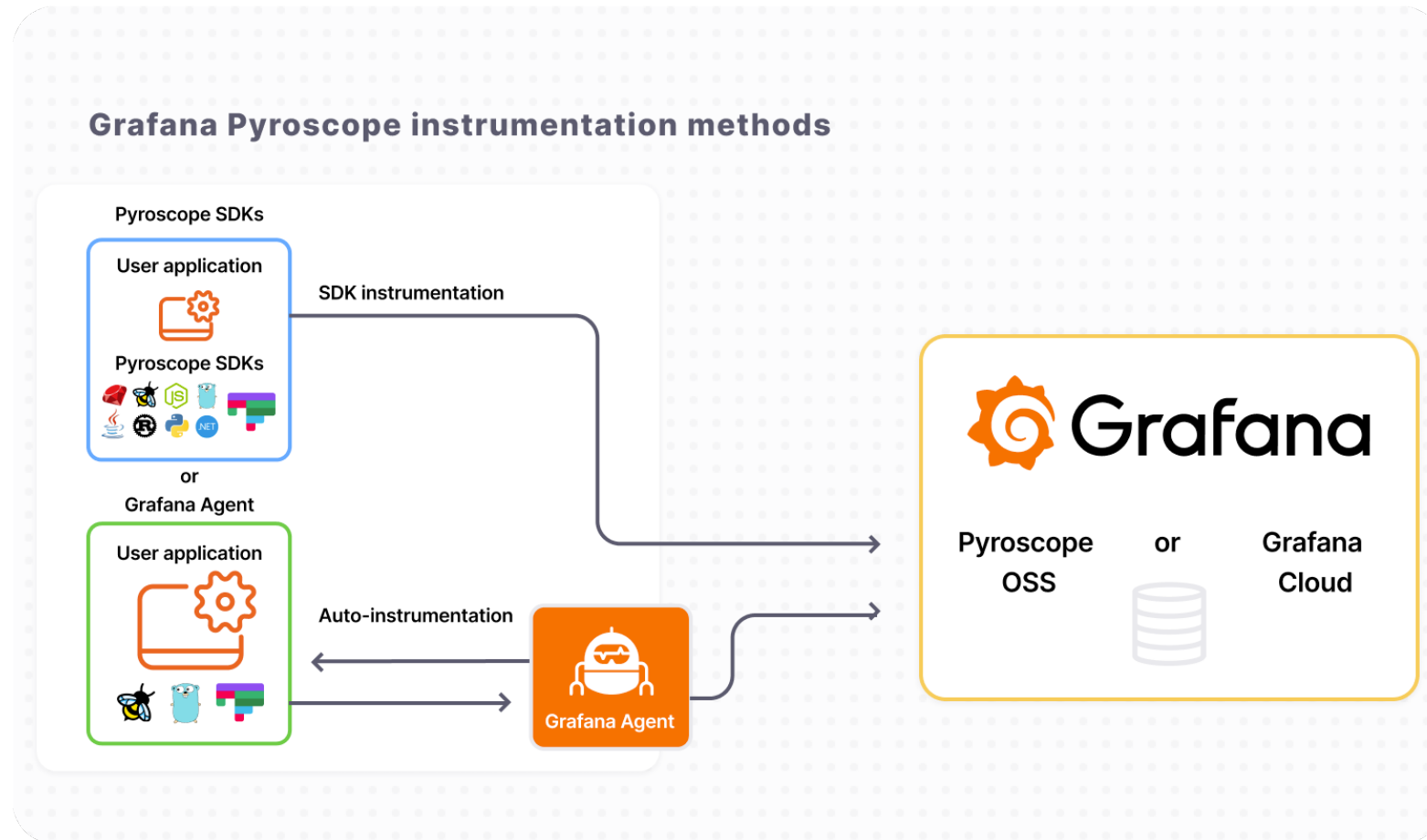
# Grafana Pyroscope

FIG. Architecture of the Grafana Pyroscope

# Java: How to collect? Java Flight Recorder

- Capture both JVM and application data
  - GC
  - Synchronization
  - Compiler
  - CPU usage
  - Exceptions
  - I/O
- Sampling-based profiler
  - Very low overhead: 2-3%
- Buffers
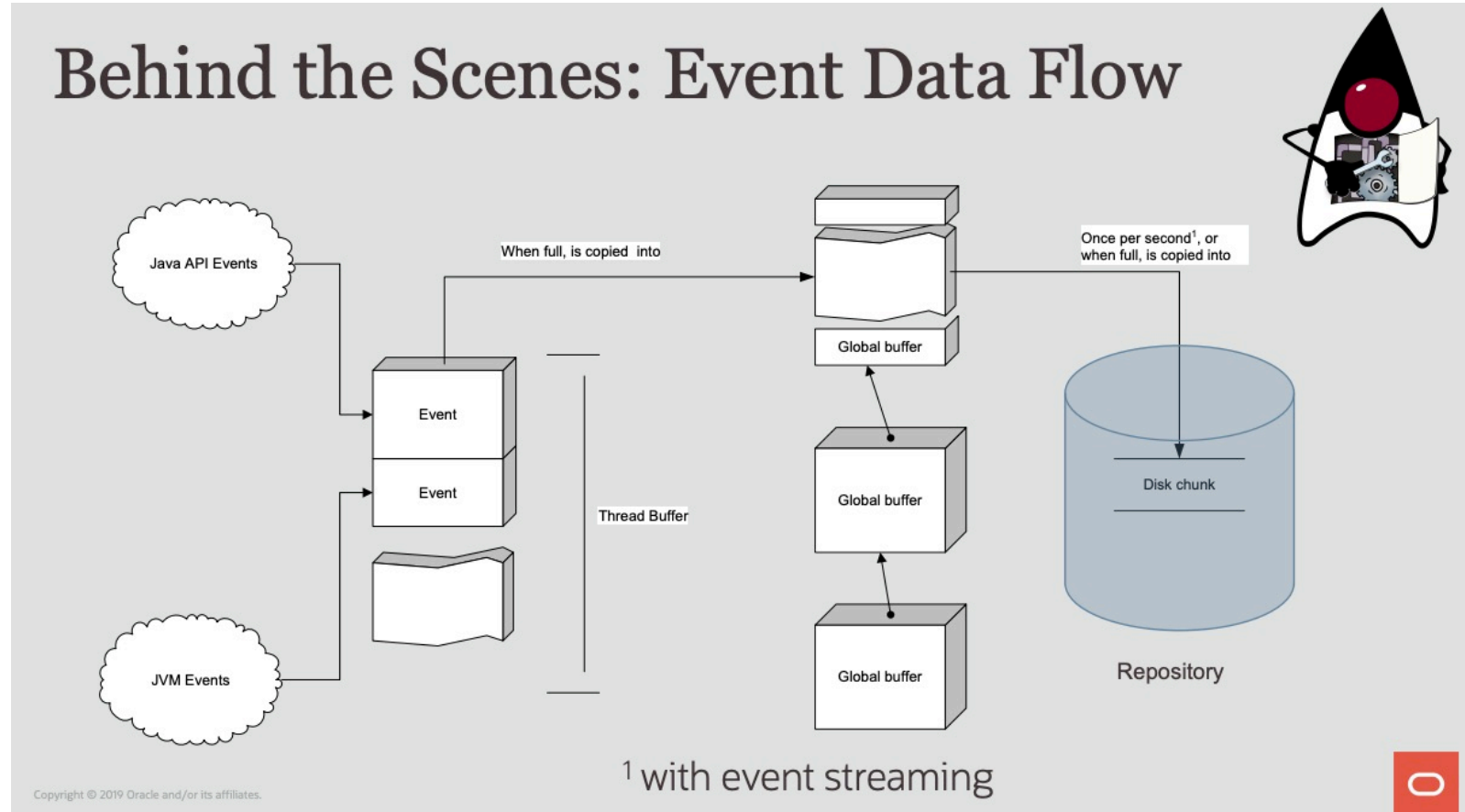  - Thread Buffer
  - Global Buffer
  - Repository (Disk chunk)



FIG. How JFR works in the background: API events and JVM events as sources. https://www.infoq.com/presentations/monitoring-jdk-jfr

# Java: How to collect? Java Flight Recorder

```
jdk.ExecutionSample {
  startTime = 2023-02-13T05:53:01.646060063Z
  sampledThread = "http-nio-8080-exec-482" (javaThreadId = 12559)
  state = "STATE_RUNNABLE"
  contextId = 0
  stackTrace = [
    java.util.LinkedHashMap.entrySet() line: 635
    java.util.HashMap.putMapEntries(Map, boolean) line: 513
    java.util.HashMap.<init>(Map) line: 491
    io.netty.bootstrap.AbstractBootstrap.copiedMap(Map) line: 429
    io.netty.bootstrap.AbstractBootstrap.options() line: 417
    ...
  ]
}
```

Event ID
Timestamp (CPU ticks)
Duration (CPU ticks)
Thread ID
StackTrace ID
Event Specific Payload

FIG. The anatomy of a JFR event and a typical example

# Java: How to collect? Async Profiler

## async-profiler 🔗

This project is a low overhead sampling profiler for Java that does not suffer from Safepoint bias problem. It features HotSpot-specific APIs to collect stack traces and to track memory allocations. The profiler works with OpenJDK, Oracle JDK and other Java runtimes based on the HotSpot JVM.

async-profiler can trace the following kinds of events:

- CPU cycles
- Hardware and Software performance counters like cache misses, branch misses, page faults, context switches etc.
- Allocations in Java Heap
- Contented lock attempts, including both Java object monitors and ReentrantLocks

# Java: How to analyze? FlameGraph

https://github.com/brendangregg/FlameGraph



FIG. A typical flamegraph

# Java: How to analyze? JDK Mission Control

https://www.azul.com/products/components/azul-mission-control/

# Arch Overview

FIG. Overview of the system design

# JFR Reader: read events w/ jfr mod

```java
package org.example;

import jdk.jfr.consumer.RecordedEvent;
import jdk.jfr.consumer.RecordingFile;


import java.nio.file.Paths;
import java.util.List;

public class App {
    public static void main(String[] args) throws Exception {
        List<RecordedEvent> events = RecordingFile.readAllEvents(Paths.get( first: "/path/to/jfr"));
        for (final RecordedEvent event : events) {
            // process...
        }
    }
}
```

```
jdk.ExecutionSample {
  startTime = 2023-02-13T05:53:01.646060063Z
  sampledThread = "http-nio-8080-exec-482" (javaThreadId = 12559)
  state = "STATE_RUNNABLE"
  contextId = 0
  stackTrace = [
    java.util.LinkedHashMap.entrySet() line: 635
    java.util.HashMap.putMapEntries(Map, boolean) line: 513
    java.util.HashMap.<init>(Map) line: 491
    io.netty.bootstrap.AbstractBootstrap.copiedMap(Map) line: 429
    io.netty.bootstrap.AbstractBootstrap.options() line: 417
    ...
  ]
}
```

FIG. Read all events and then decode (JDK 8u262+)

# JFR Reader: build call stack



```java
5 usages
12  public class Tree {

    2 usages
13      private final TreeNode root = new TreeNode( name: "");

    1 usage
15      public void insertStackFrames(List<RecordedFrame> frames, long v) {
16          TreeNode n = this.root;
17          for (final RecordedFrame frame : Lists.reverse(frames)) {
18              n.total += v;
19              final RecordedMethod m = frame.getMethod();
20              final String frameStr = m.getType().getName() + "." + m.getName();
21              n = n.insertString(frameStr);
22          }
23          // Leaf.
24          n.total += v;
25          n.self += v;
26      }
```

```java
29      public static class TreeNode {
30          @Getter
31          private final String name;

        2 usages
32          private long total;

        1 usage
33          private long self;

        5 usages
34          private final List<TreeNode> childrenNodes;
```

FIG. Build the call stack (Tree with treeNode as children)

# JFR Reader: build call stack (~80M)

```
8  ▷     public class App {
9  ▷         public static void main(String[] args) throws Exception {
10               Tree callStack = new Tree();
11  [63.35 MB]
12               for (final RecordedEvent event : RecordingFile.readAllEvents(Paths.get( first: "/Users/megrez/Downloads/fcb95fca8e9
13  [647.24 MB] v       if (event ≠ null) {
14                           decodeEvent(event, callStack);
15                       }
16               }
17
         1 usage
18  @        public static void decodeEvent(RecordedEvent event, Tree callStack) {
19  v            switch (event.getEventType().getName()) {
20                   case "jdk.ObjectAllocationInNewTLAB":
21  [🔥 42.76 GB]            callStack.insertStackString(event.getStackTraceFrames(), event.getLong( name: "allocationSize"));
22                       break;
23                   }
24               }
25       }
```

~120 million

FIG. Memory issue: large heap size while building the call stack
with millions of (allocation) events

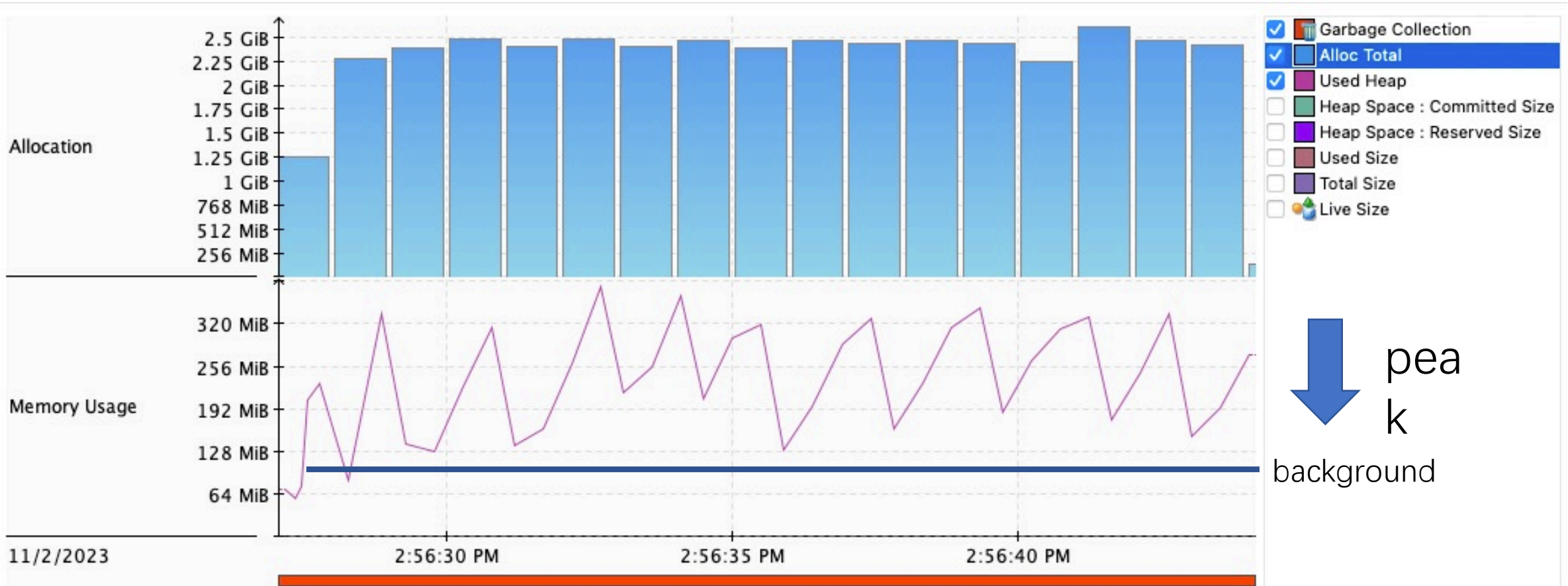# JFR Reader: build call stack (~80M)



FIG. Memory issue: large heap size while building the call stack
with millions of (allocation) events

# JFR Reader: Iterator pattern



```java
public class App {
    public static void main(String[] args) throws Exception {
        Tree callStack = new Tree();
        try (RecordingFile recordingFile = new RecordingFile(Paths.get( first: "/Users/megrez/Downloads/fcb95fca8e93dde7d921
            while (recordingFile.hasMoreEvents()) {
                final RecordedEvent event = recordingFile.readEvent();
                if (event != null) {
                    decodeEvent(event, callStack);
                }
            }
        }
    }

    1 usage
    public static void decodeEvent(RecordedEvent event, Tree callStack) {
        switch (event.getEventType().getName()) {
            case "jdk.ObjectAllocationInNewTLAB":
                callStack.insertStackString(event.getStackTrace().getFrames(), event.getLong( name: "allocationSize"));
                break;
        }
    }
}
```

FIG. Process RecordEvent one by one

# JFR Reader: Iterator pattern

FIG. Memory issue: large heap size while building the call stack with millions of (allocation) events

# JFR Reader: Slow!

```
5 usages
12  public class Tree {
        2 usages
13      private final TreeNode root = new TreeNode( name: "");
14

        1 usage
15      public void insertStackFrames(List<RecordedFrame> frames, long v) {
16          TreeNode n = this.root;
17          for (final RecordedFrame frame : Lists.reverse(frames)) {
18              n.total += v;
19              final RecordedMethod m = frame.getMethod();
20              final String frameStr = m.getType().getName() + "." + m.getName();
21              n = n.insertString(frameStr);
22          }
23          // Leaf.
24          n.total += v;
25          n.self += v;
26      }
```

| | |
|---|---|
| 17 | 300 ms |
| 19 | 2,079 ms |
| 20 | ⚠ 10,236 ms |
| 21 | 1,229 ms |
| 22 | 30 ms |

FIG. Performance issue: most time spent on building frame names

# JFR Reader: use raw references

```java
public class StackTrace {
    // 方法ID
    public final long[] methods;
    // 每个byte表示对应的方法类型，有INTERPRETED，JIT_COMPILED等
    public final byte[] types;
    // 每个int表示方法所在的行号和bci
    public final int[] locations;
    // ...
}
```

FIG. use references instead of materialized stack trace

# JFR Reader: use raw references

通过JDK原生的方式读取      通过async-profiler读取

| Event1 | stackTrace1 |
| Event2 | stackTrace1 |
| Event3 | stackTrace1 |
| Event4 | stackTrace2 |
| Event5 | stackTrace2 |
| Event6 | stackTrace3 |

size=240万

| Event1 | stackTraceId1 |
| Event2 | stackTraceId1 |
| Event3 | stackTraceId1 |
| Event4 | stackTraceId2 |
| Event5 | stackTraceId2 |
| Event6 | stackTraceId3 |

size=240万

Set
- stackTrace1
- stackTrace2
- stackTrace3

size=3万

FIG. 2,000,000 alloc events share 30,000 stacktraces

# JFR Reader: binary search O(logN)?

```
2 usages
private TreeNode insertString(String stackStr) {
    final int i = Collections.binarySearch(
            // lazy transform
            Lists.transform(this.childrenNodes, TreeNode::getName),
            stackStr);
    if (i < 0) { // not found
        final int insertionPoint = -(i + 1);
        this.childrenNodes.add(insertionPoint, new TreeNode(stackStr));
        return this.childrenNodes.get(insertionPoint);
    }
    return this.childrenNodes.get(i);
}
```

FIG. Another performance issue: too many binary searches
during insertion even if binary search has O(logN) complexity

# JFR Reader: insert first

pyroscope

# JFR Reader: aggregate first

# JFR Reader: final round

FIG. Final result: use **<100M** heap, and  finish parsing
**<1 second**

# JFR Reader: What about large JFR file?

# JFR Reader: What about large JFR file?

chunksize=10m



FIG. ChunkSize can be controlled by parameter

# JFR Reader: What about large JFR file?

# One more thing: correlation

**02**

# 交互式诊断
# Interactive Diag.

# How to diag. a CPU spike

██有一个节点 cpu 和 young gc 次数 遥遥领先,
看起来很奇怪

# How to diag. a CPU spike: Arthas

# How to integrate SkyWalking with Arthas

FIG. 将 Apache SkyWalking 与 Arthas 集成 By 魏翔
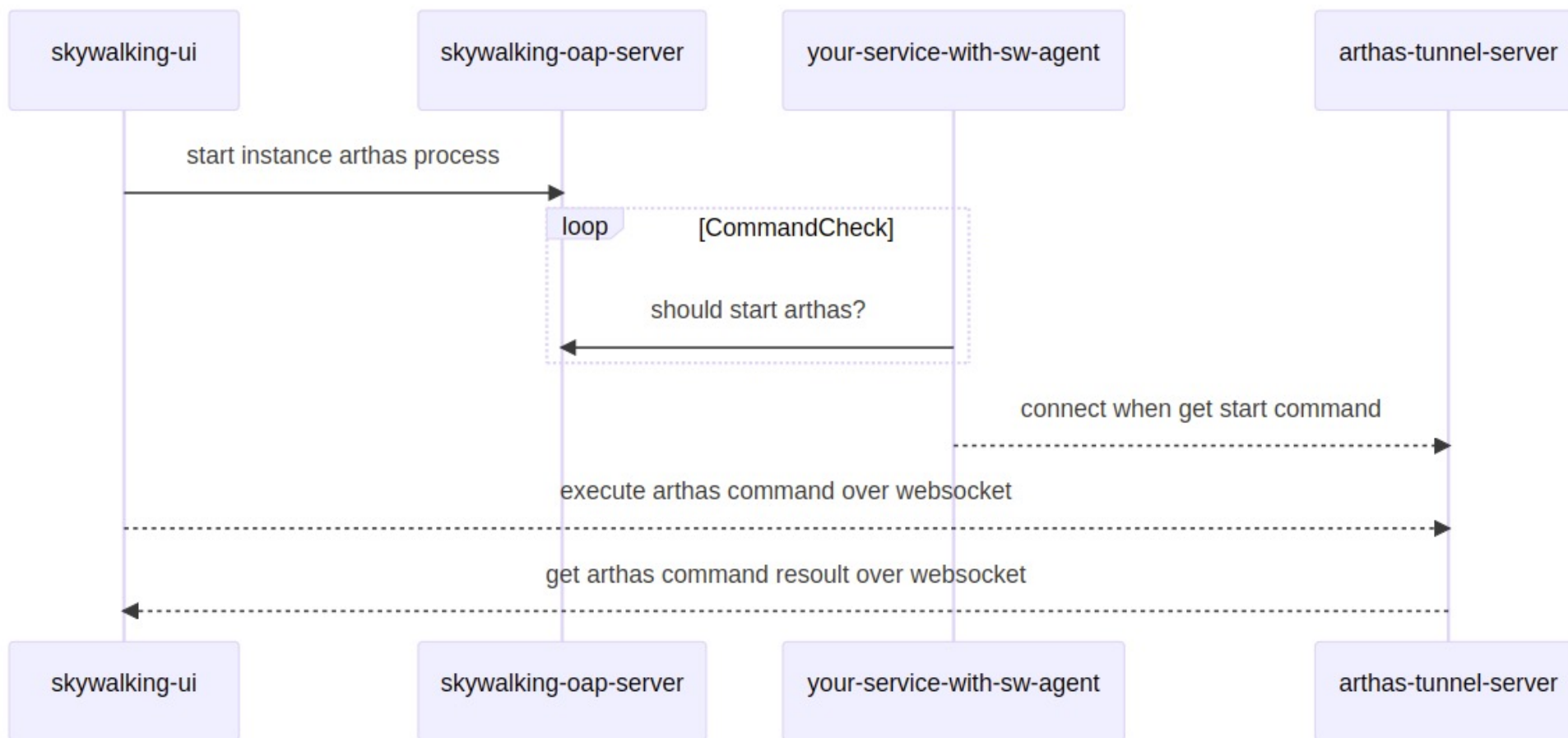https://skywalking.apache.org/zh/2023-09-17-integrating-skywalking-with-arthas/
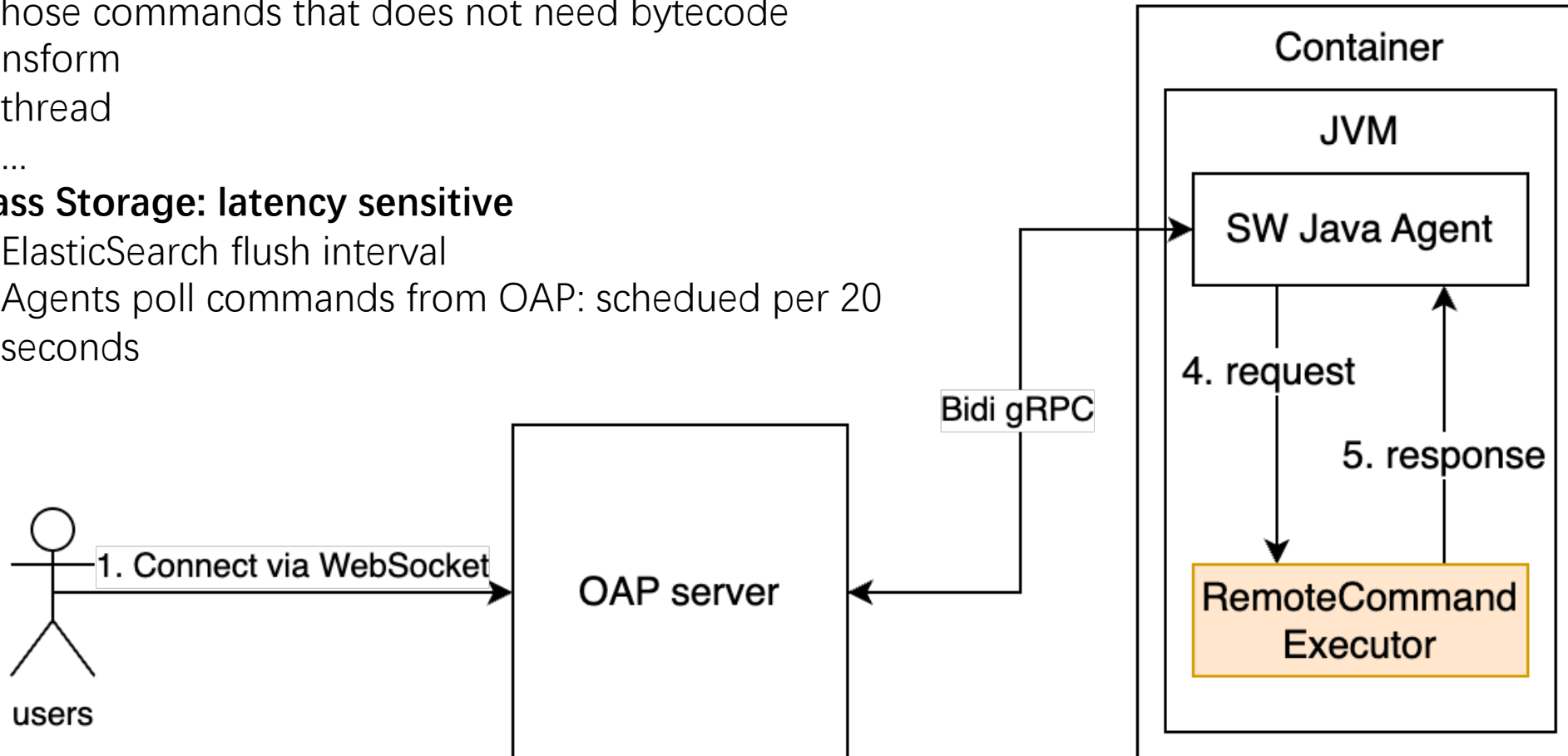
# How to integrate SkyWalking with Arthas

- For those commands that does not need bytecode retransform
  - thread
  - ...
- **Bypass Storage: latency sensitive**
  - ElasticSearch flush interval
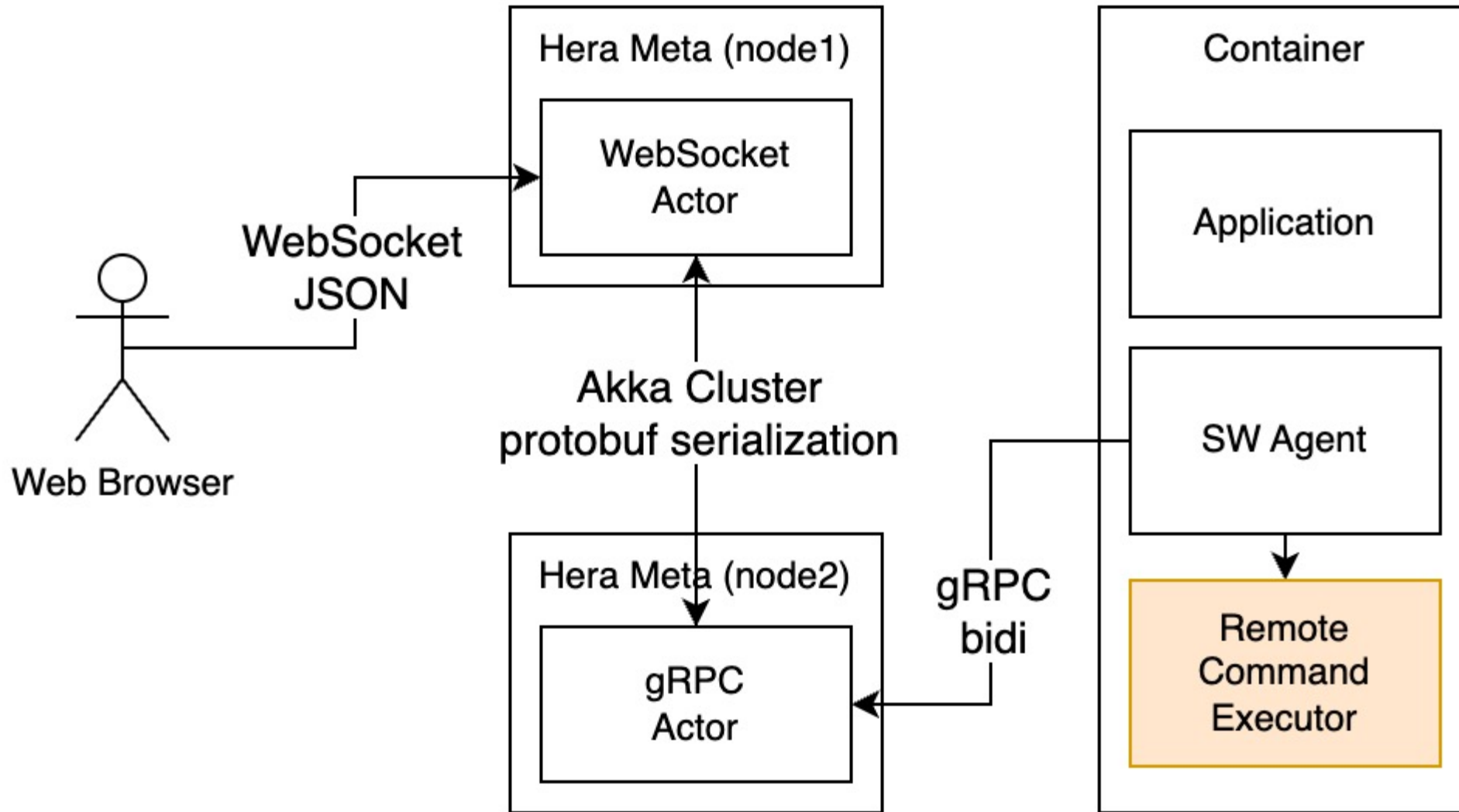  - Agents poll commands from OAP: schedued per 20 seconds

# Protocol Design: bidi over unary

```
29   service ProfileTask {
30
31       // query all sniffer need to execute profile task commands
32       rpc getProfileTaskCommands (ProfileTaskCommandQuery) returns (common.v1.Commands) {
33       }
34
35       // collect dumped thread snapshot
36       rpc collectSnapshot (stream ThreadSnapshot) returns (common.v1.Commands) {
37       }
38
39       // report profiling task finished
40       rpc reportTaskFinish (ProfileTaskFinishReport) returns (common.v1.Commands) {
41       }
42
43   }
```

```
29   service RemoteCommandTask {
30     // collect remote command result
31     rpc executeRemoteCommand (stream RemoteCommandRequest) returns (stream RemoteCommandResponse) {
32     }
33   }
```

# What about distributed OAP?

# What about retransform?

https://github.com/apache/skywalking/blob/master/docs/en/FAQ/Compatible-with-other-javaagent-bytecode-processing.md#compatibility-with-other-java-agent-bytecode-processes

## Problem 🔗

1. When using the SkyWalking agent, some other agents, such as Arthas, can't work properly. #4858

2. The retransform classes in the Java agent conflict with the SkyWalking agent, as illustrated in this demo

## Cause 🔗

The SkyWalking agent uses ByteBuddy to transform classes when the Java application starts. ByteBuddy generates auxiliary classes with different random names every time.

When another Java agent retransforms the same class, it triggers the SkyWalking agent to enhance the class again. Since the bytecode has been regenerated by ByteBuddy, the fields and imported class names have been modified, and the JVM verifications on class bytecode have failed, the retransform classes would therefore be unsuccessful.

## Resolution 🔗
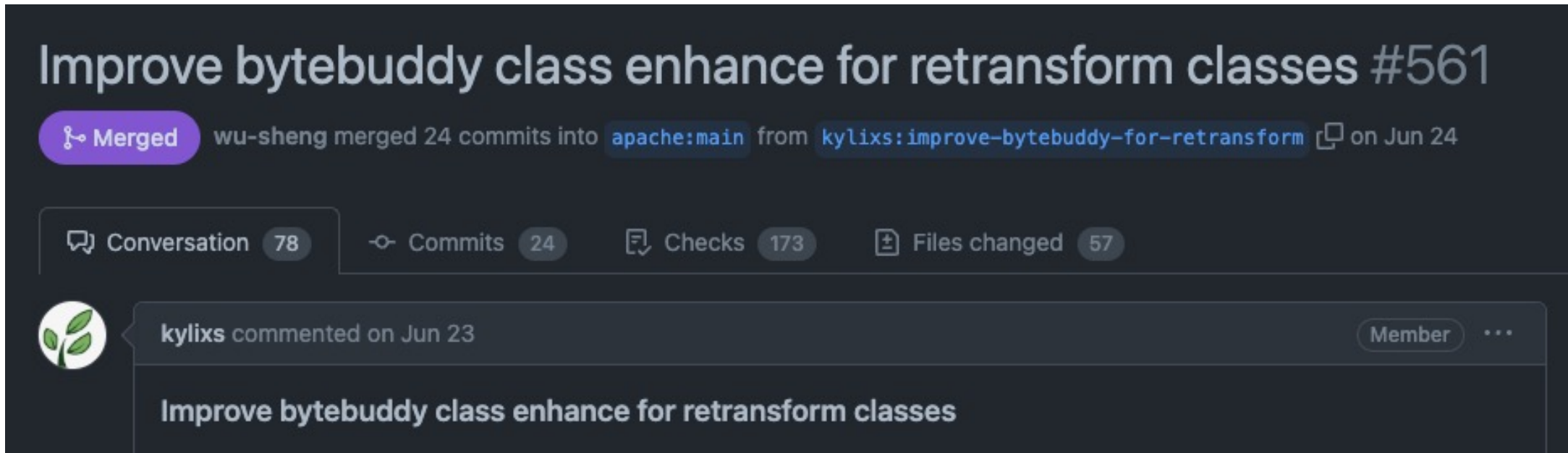
**1. Enable the class cache feature**

Add JVM parameters:

```
-Dskywalking.agent.is_cache_enhanced_class=true -Dskywalking.agent.class_cache_mode=MEMORY
```

- For those commands that does need bytecode retransform,
  - watch: observe method exec (parameter, result, exception...)
  - trace: trace method exec path
  - monitor: stat method exec (not real time)
- Main idea
  - For TypeDescription: always perfer bytecode from TypePool to reflection API
  - For aux. fields/methods: use stable prefix/suffix instead of random ones

# Changes in 9.0: perf issue (resolved)

FIG. Using POOL_FIRST TypeDescription strategy in SW Java 9.0 caused almost double application launch time and larger heap size. Resolved in PR #637.