# 后分布式追踪时代的性能问题定位
# -- 方法级性能剖析

Sheng Wu 吴晟, Founding Engineer, Tetrate.io

Apache SkyWalking Founder & V.P. Apache Member
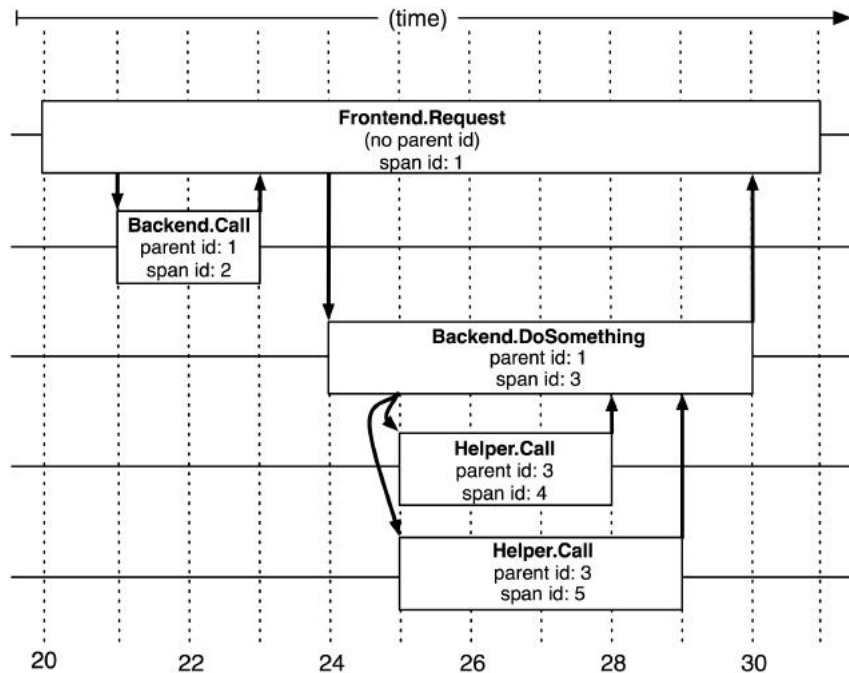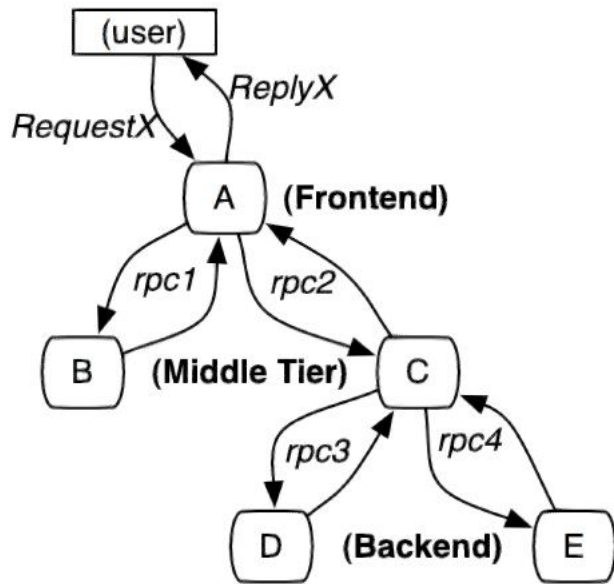
POWERED BY APACHE

Skywalking

tetrate

# （微）服务化、云原生化、容器化

- 扩容和管理需求

- 分布式程序加剧

- 依赖复杂

- 问题诊断困难

- 传统手段失效，特别是传统日志

# 2010 Google Dapper

- Google搜索引擎的分布式追踪探索
- 一篇广为流传的论文 https://research.google/pubs/pub36356/
- 2012 Twitter实现Zipkin
- 2012 Naver开源Pinpoint
- 2015 个人SkyWalking开源
- 2017 SkyWalking加入Apache孵化器
- 2019 SkyWalking成为Apache顶级项目

# 分布式追踪简介 – 带上下文的日志

为什么SkyWalking和传统分布式追踪系统不同？

Distributed Tracing
V.S.
Distributed Tracing based APM

# SkyWalking 原生概念

SkyWalking provides observability capabilities for **service**(s), **service instance**(s), **endpoint**(s). The terms Service, Instance and Endpoint are used everywhere today, so it is worth defining their specific meanings in the context of SkyWalking:

- **Service**. Represents a set/group of workloads which provide the same behaviours for incoming requests. You can define the service name when you are using instrument agents or SDKs. SkyWalking can also use the name you define in platforms such as Istio.
- **Service Instance**. Each individual workload in the Service group is known as an instance. Like `pods` in Kubernetes, it doesn't need to be a single OS process, however, if you are using instrument agents, an instance is actually a real OS process.
- **Endpoint**. A path in a service for incoming requests, such as an HTTP URI path or a gRPC service class + method signature.

# SkyWalking 原生概念

- **Span**
  - Entry/Local/Exit
  - Peer
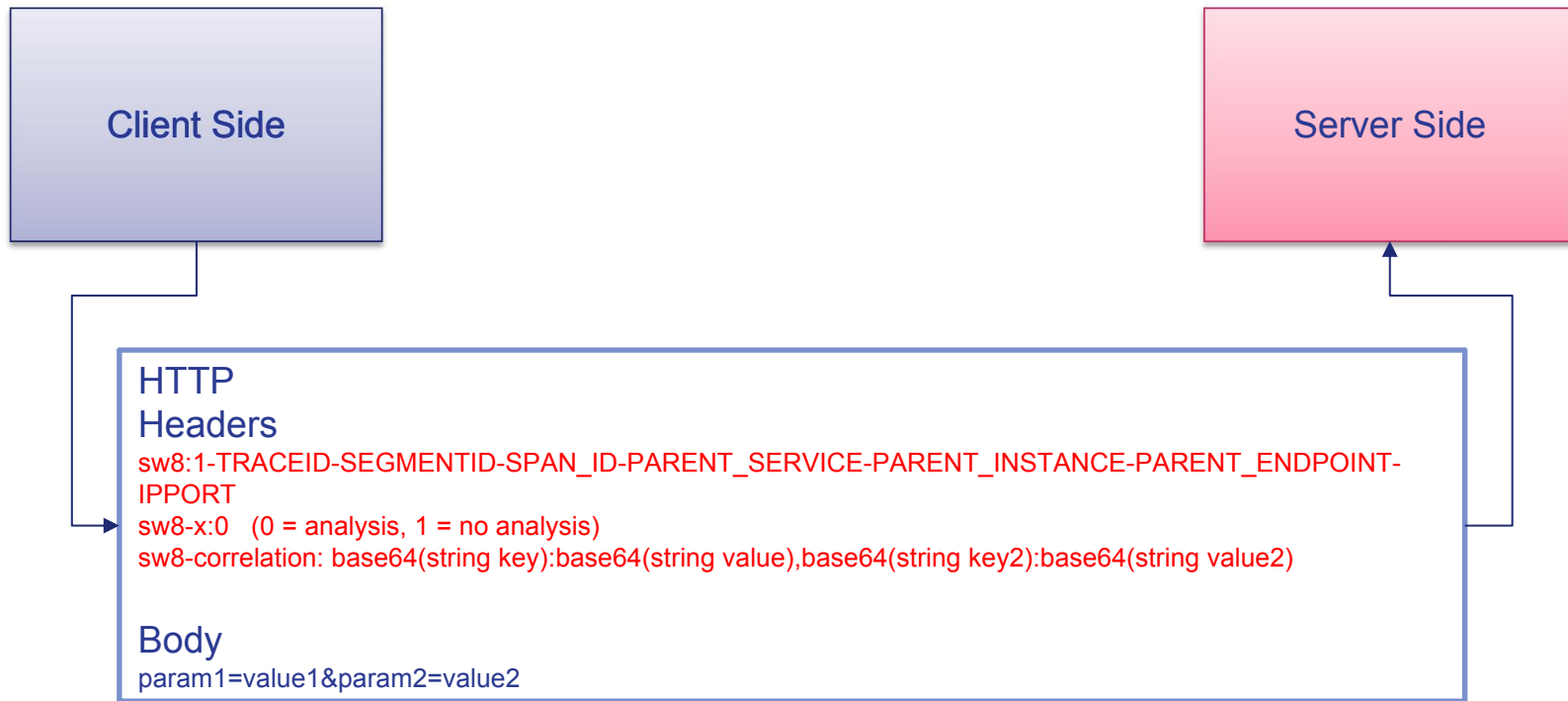  - Component
  - Layer
  - Tags
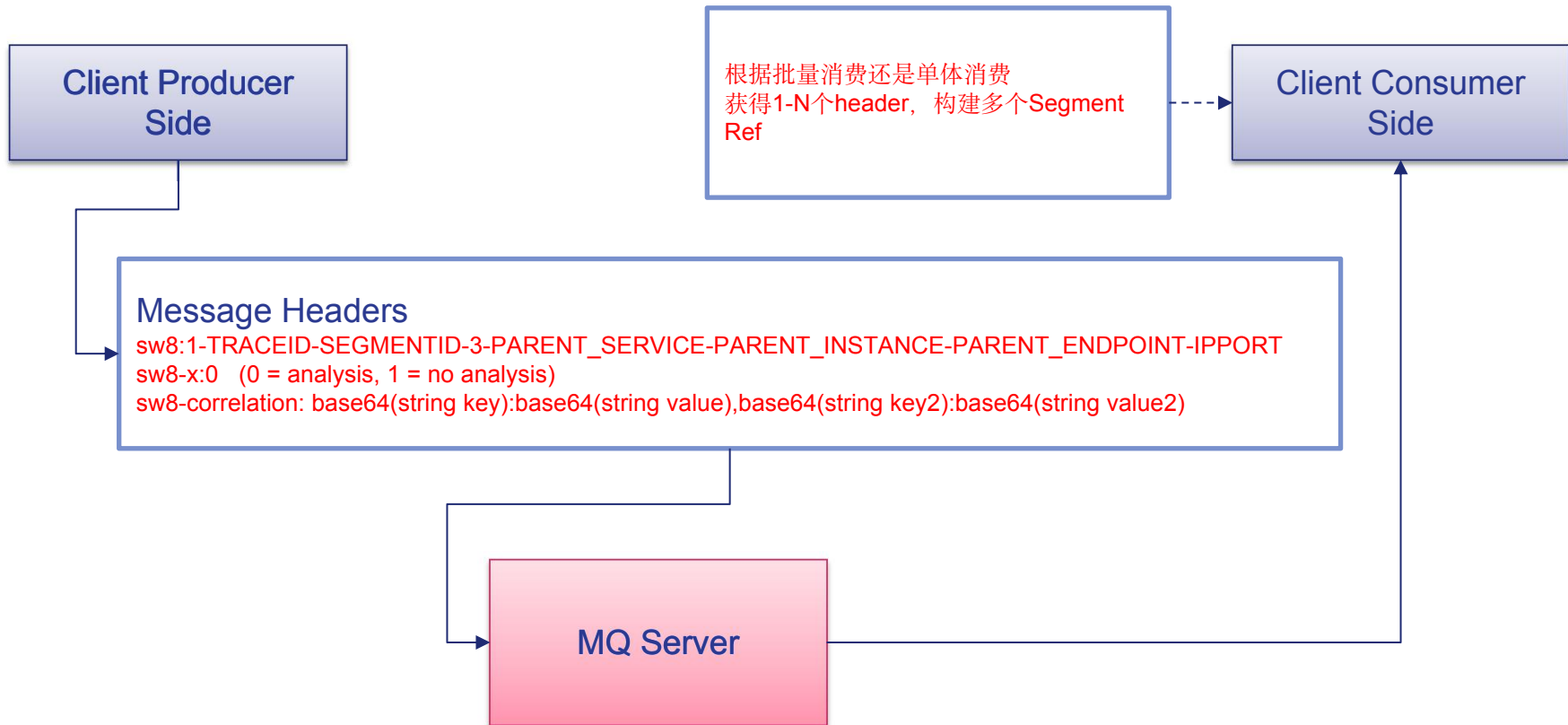  - Logs

- **Segment**
  - Span Collection

- **Segment Ref**
  - Trace Id
  - Parent Segment Id
  - Parent Span Id
  - Parent service and instance
  - Parent endpoint

# RPC Tracing

Client Side

Server Side

HTTP
Headers
sw8:1-TRACEID-SEGMENTID-SPAN_ID-PARENT_SERVICE-PARENT_INSTANCE-PARENT_ENDPOINT-IPPORT
sw8-x:0   (0 = analysis, 1 = no analysis)
sw8-correlation: base64(string key):base64(string value),base64(string key2):base64(string value2)

Body
param1=value1&param2=value2

# MQ Tracing

Client Producer Side

根据批量消费还是单体消费
获得1-N个header，构建多个Segment Ref

Client Consumer Side

Message Headers
sw8:1-TRACEID-SEGMENTID-3-PARENT_SERVICE-PARENT_INSTANCE-PARENT_ENDPOINT-IPPORT
sw8-x:0   (0 = analysis, 1 = no analysis)
sw8-correlation: base64(string key):base64(string value),base64(string key2):base64(string value2)

MQ Server

分布式调用栈

# 分布式追踪的两面性

- 有效的移除分布式造成的信息缺失

  - 如RPC故障

  - 路由以及熔断情况

  - 分布式访问参数

- 重新加强日志的能力，trace和log可以集成使用

- Trace流量可能高达数百T每天

  - 追踪探针的额外负载

  - 后端分析能力

  - 后端存储能力

  - 数据时间有效性（TTL）

- Trace只是一个有上下文的日志么?

- 通过Trace还能得到什么?

# 分布式追踪消耗 – Span Overload

- 对于100ms / 1000+QPS 级别的应用访问，

  大概造成10%的性能损耗（< 20个span）

- 方法级追踪 -> 分布式追踪

- 日志 -> 分布式追踪

- 100%采样 （探针 / 后端）

- 直接上报 / MQ / Local Collector

# Paper, STAM: Enhancing Topology Auto Detection For A Highly Distributed and Large-Scale Application System

- SkyWalking特有的拓扑分析方法

- Peer别名机制

- https://wu-sheng.github.io/STAM/
- 中文翻译 https://wu-sheng.github.io/STAM/README-cn

如何才能正确的定位问题

分布式追踪提供了数据

APM工具提供
聚合和
数据可视化



DATA

SORTED

ARRANGED

PRESENTED
VISUALLY

# 拓扑图的价值



- 理解系统
- 理论架构 -> 实际架构
- 领导需要
- 提供下钻点

# 监控的起点，metrics和告警

| Service Apdex | Service Avg Response Time ( ms ) | Successful Rate ( % ) | Service Load ( CPM – calls per minute ) |
|---|---|---|---|
| 0.32 |  | 100.00 | 209.00 |

| Service Apdex | Service Response Time Percentile ( ms ) | Successful Rate ( % ) | Service Load ( CPM – calls per minute ) |
|---|---|---|---|
|  |  |  |  |

## Entity name

Define the relation between scope and entity name.

- **Service**: Service name
- **Instance**: {Instance name} of {Service name}
- **Endpoint**: {Endpoint name} in {Service name}
- **Database**: Database service name
- **Service Relation**: {Source service name} to {Dest service name}
- **Instance Relation**: {Source instance name} of {Source service name} to {Dest instance name} of {Dest service name}
- **Endpoint Relation**: {Source endpoint name} in {Source Service name} to {Dest endpoint name} in {Dest service name}

# 下钻检查 – 实例 or Endpoint

| Service Instances Load（CPM – calls per minute） | |
|---|---|
| 104 | 612471e247804456983281b6cda70d64@192.168.252.12 |
| 104 | 77cca827da5943b1a1cd21bb81e8f18c@192.168.252.13 |

| Slow Service Instance（ms） | |
|---|---|
| 1869 | 612471e247804456983281b6cda70d64@192.168.252.12 |
| 1858 | 77cca827da5943b1a1cd21bb81e8f18c@192.168.252.13 |

| Service Instance Successful Rate（%） | |
|---|---|
| 100 | 612471e247804456983281b6cda70d64@192.168.252.12 |
| 100 | 77cca827da5943b1a1cd21bb81e8f18c@192.168.252.13 |

| Endpoint Load in Current Service（CPM – calls per minute） | |
|---|---|
| 209 | /projectA/{name} |

| Slow Endpoints in Current Service（ms） | |
|---|---|
| 1863 | /projectA/{name} |

| Successful Rate in Current Service（%） | |
|---|---|
| 100 | /projectA/{name} |

Endpoint Load

Endpoint Avg Response Time（ms）

Endpoint Response Time Percentile（ms）
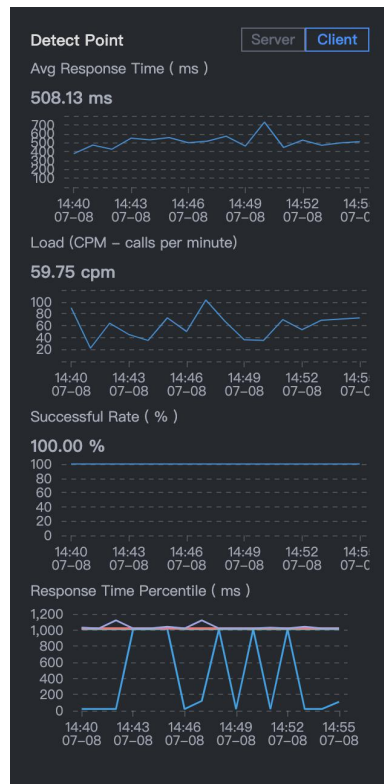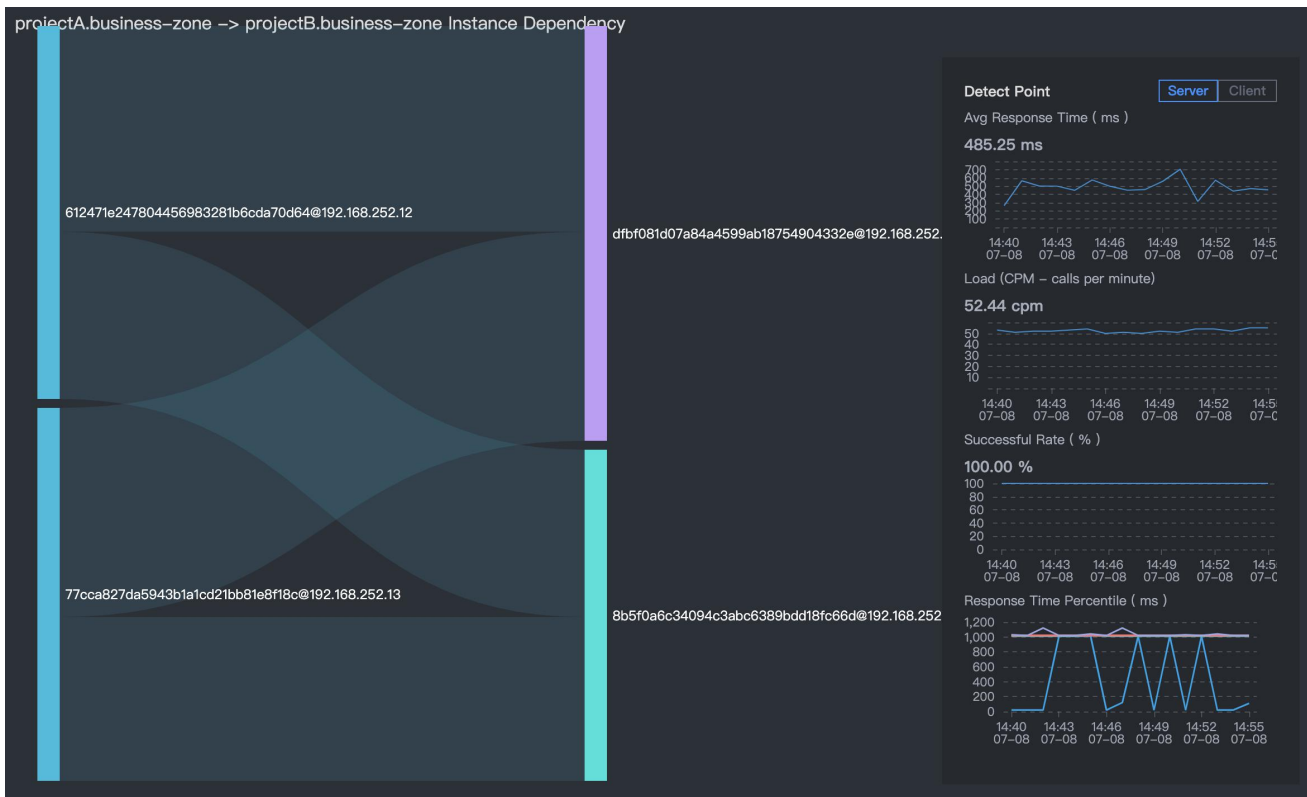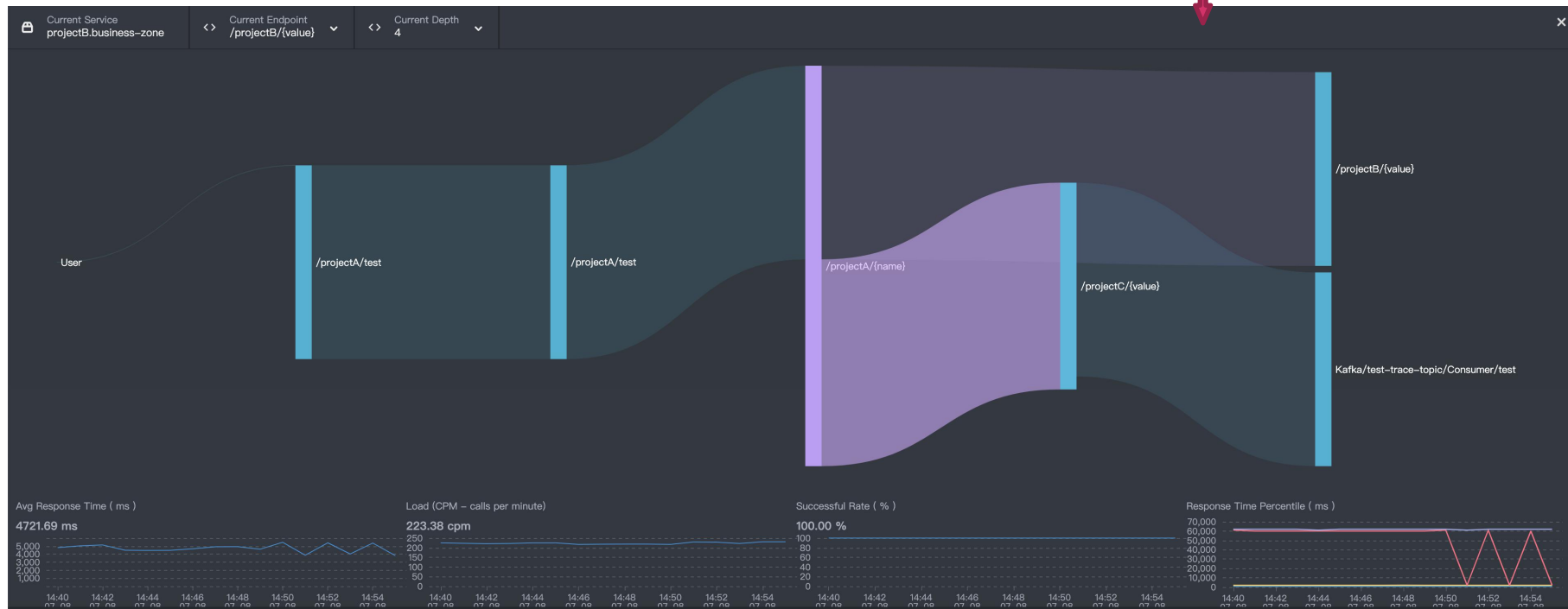● P50  ● P75  ● P90  ● P95  ● P99

Endpoint Successful Rate（%）

# 下钻检查 – 分布式依赖维度

# 下钻检查－实例级别依赖故障（限流、网络）

# 业务服务上下游依赖

# 明细级别探查 – Error Span

/projectC/{value}

9f3fa51b-78cc-44a4-b9f7-29112f31bd88 ▾

Start 2020-07-08 14:47:00  Duration 62715 ms  Spans 18

☰ List  🖧 Tree  ☷ Table

| | | | | | | |
|---|---|---|---|---|---|---|
| ⌄ /projectA/test | 2020-07-08 14:47:00 | 5056 | | 0 | Nginx | load balancer2.system |
| ⌄ /projectA/{name} | 2020-07-08 14:47:00 | 5056 | | 1 | SpringMVC | projectA.business-zone |
| ⌄ /projectB/test | 2020-07-08 14:47:00 | 28 | | 1 | SpringRestTempl... | projectA.business-zone |
| ⌄ /projectB/{value} | 2020-07-08 14:47:00 | 27 | | 0 | SpringMVC | projectB.business-zone |
| ⌄ org.skywalking.springcloud.test.projectb.dao.DatabaseOperateDao.saveUse... | 2020-07-08 14:47:00 | 0 | | 0 | Unknown | projectB.business-zone |
| H2/JDBI/PreparedStatement/execute | 2020-07-08 14:47:00 | 0 | | 0 | h2-jdbc-driver | projectB.business-zone |
| ⌄ selectUser | 2020-07-08 14:47:00 | 27 | | 1 | Unknown | projectB.business-zone |
| H2/JDBI/PreparedStatement/execute | 2020-07-08 14:47:00 | 26 | | 26 | h2-jdbc-driver | projectB.business-zone |
| ⌄ /projectC/{name} | 2020-07-08 14:47:00 | 5027 | | 0 | Feign | projectA.business-zone |
| ⌄ /projectC/{value} | 2020-07-08 14:47:05 | 24 | | 2 | SpringMVC | projectC.business-zone |
| / | 2020-07-08 14:47:05 | 21 | | 21 | HttpClient | projectC.business-zone |
| Kafka/test-trace-topic/Producer | 2020-07-08 14:47:05 | 1 | | 1 | kafka-producer | projectC.business-zone |
| ⌄ /projectC/{value} | 2020-07-08 14:47:00 | 62715 | | 1002 | SpringMVC | projectC.business-zone |
| / | 2020-07-08 14:47:01 | 1713 | | 1713 | HttpClient | projectC.business-zone |
| Kafka/test-trace-topic/Producer | 2020-07-08 14:47:03 | 60000 | | 60000 | kafka-producer | projectC.business-zone |

# 明细级别探查 – Slow Span

# 明细级别探查 – Access Overload

# 明细级别探查 – Dependency Depth

# 如果上述方法依赖失效？

Trace的局限性

# 基于线程快照的执行时间估算



- org.test.code.classA.methodA(ClassA.java, line:24) T2-T1 ms
- org.test.code.classA.methodB(ClassA.java, line:24) T2-T1 ms
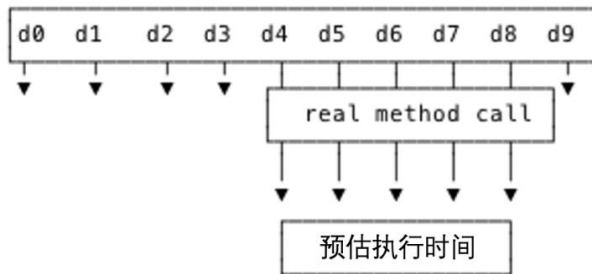- org.test.code.classC.methodB(ClassA.java, line:33) 0 ms
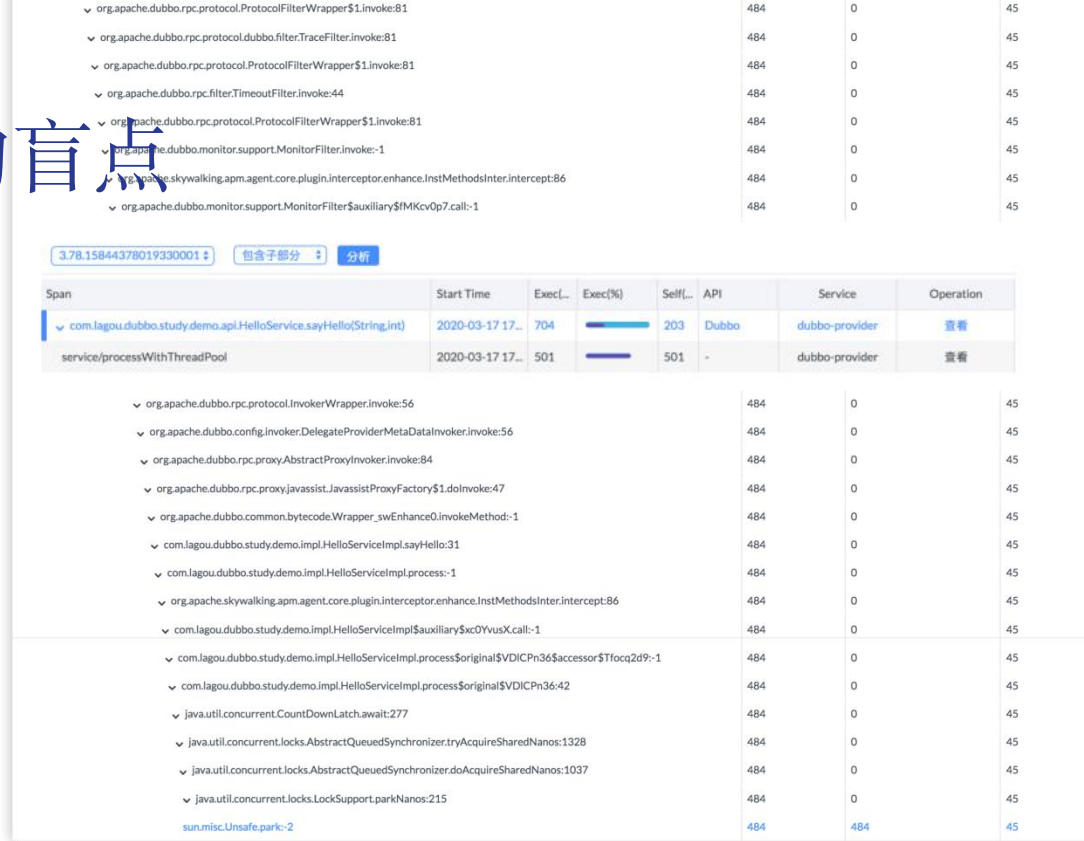
T1 Snapshot

- org.test.code.classA.methodA(ClassA.java, line:24)
  - org.test.code.classA.methodB(ClassA.java, line:24)
    - org.test.code.classC.methodB(ClassA.java, line:33)
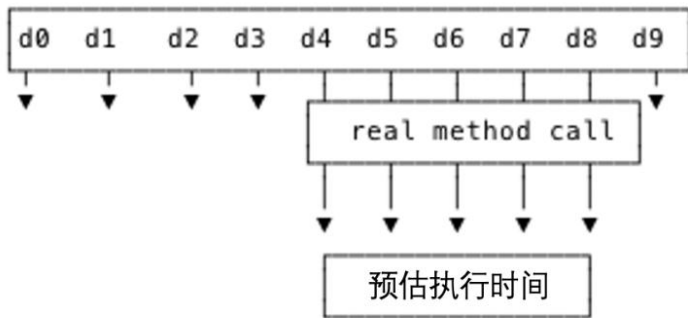
T2 Snapshot

- org.test.code.classA.methodA(ClassA.java, line:24)
  - org.test.code.classA.methodB(ClassA.java, line:35)
    - org.test.code.classD.methodB(ClassA.java, line:42)

# 性能剖析 – Trace的盲点





- InfoQ文章 https://www.infoq.cn/article/CWUOl1JA0EyXw0CxQ4Zm
- TheNewStack（英文）https://thenewstack.io/apache-skywalking-use-profiling-to-fix-the-blind-spot-of-distributed-tracing/

# Open Source Community

- 基于Trace的诊断分析, http://skywalking.apache.org/blog/2019-01-01-Understand-Trace.html
- 多语言

    ○ Java, https://github.com/apache/skywalking

    ○ Nginx Lua, https://github.com/apache/skywalking-nginx-lua

    ○ Python, https://github.com/apache/skywalking-python

    ○ PHP, https://github.com/SkyAPM/SkyAPM-php-sdk

    ○ .Net, https://github.com/SkyAPM/SkyAPM-dotnet

    ○ NodeJs, https://github.com/SkyAPM/SkyAPM-nodejs | https://github.com/apache/skywalking-nodejs

    ○ Golang, https://github.com/SkyAPM/go2sky

- 可以选择任何你喜欢的语言

# Apache Local Community



# Cloud Native Community