# Contents
# 目录

FLINK
FORWARD

相关背景

Background

01

# 计算架构

Computing Architecture

FLINK FORWARD

| User | Third-party Platform | Smart Resource |
| --- | --- | --- |
| Streaming Compute Platform | | Trace & Debug |
| Flink SQL | Flink Streaming Flink Batch | |
| YARN | | Dashboard |

# 计算规模

Scale of Computing

1s+ Yarn Cluster
1w+ Machines

100s+ Users

1k+ Flink Streaming
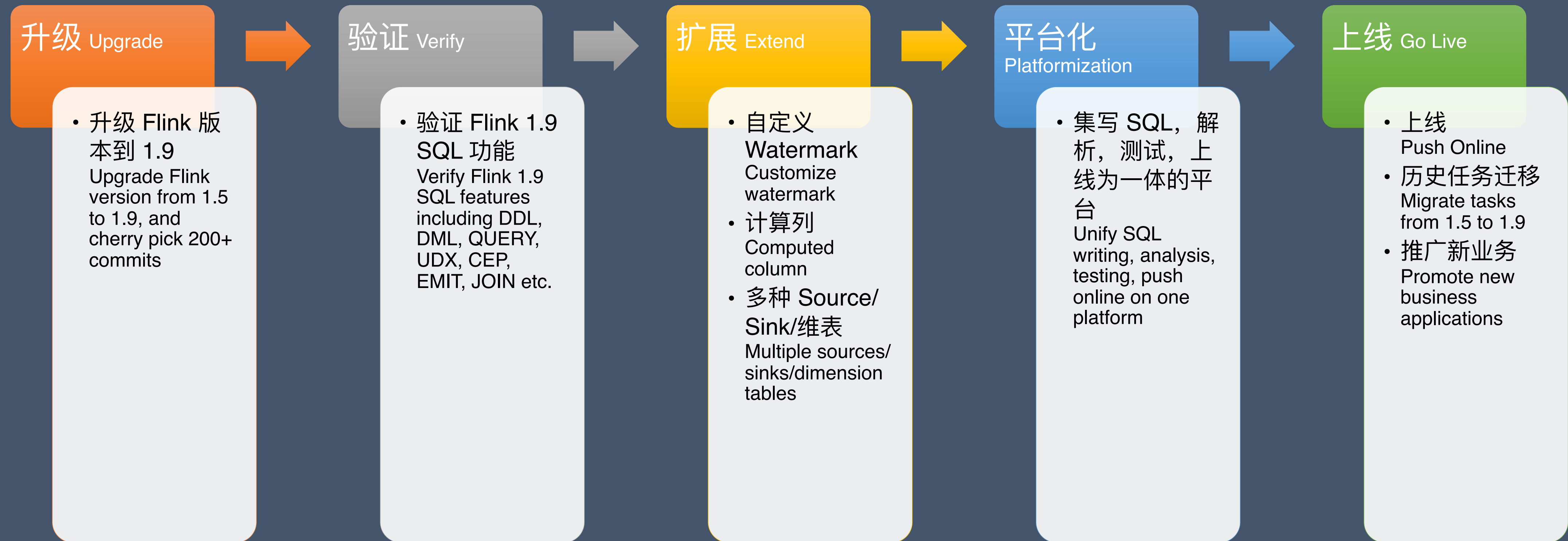1w+ Flink Batch
100s+ of Flink SQL

# 核心关注

Main Focuses

- 负载均衡 Load Balancing
- 环境隔离 Isolation
- 作业故障率 Failure Rate
- CP 成功率 Checkpoint Success Rate
- 双机房容灾 Disaster Tolerance

- 平台化 Platformization
- 便捷 Debug Easy Debug

可靠
Reliability

易用
Usability

引擎能力
Power of Engine

- 状态计算 Stateful Computations
- 有且仅有 Exactly-Once
- 窗口 Windowing
- SQL SQL
- 批流统一 Unified Batch & Streaming

FLINK FORWARD

# Flink SQL 的应用和扩展

Flink SQL Applications and Extensions

**02**

# Flink 1.9 SQL 上线

Flink 1.9 SQL Goes Live

**升级 Upgrade** → **验证 Verify** → **扩展 Extend** → **平台化 Platformization** → **上线 Go Live**

- 升级 Flink 版本到 1.9
  Upgrade Flink version from 1.5 to 1.9, and cherry pick 200+ commits

- 验证 Flink 1.9 SQL 功能
  Verify Flink 1.9 SQL features including DDL, DML, QUERY, UDX, CEP, EMIT, JOIN etc.

- 自定义 Watermark
  Customize watermark
- 计算列
  Computed column
- 多种 Source/Sink/维表
  Multiple sources/sinks/dimension tables

- 集写 SQL，解析，测试，上线为一体的平台
  Unify SQL writing, analysis, testing, push online on one platform

- 上线
  Push Online
- 历史任务迁移
  Migrate tasks from 1.5 to 1.9
- 推广新业务
  Promote new business applications

# Flink SQL 扩展

Flink SQL Extensions

**FLINK FORWARD**

Source/Sink: Kafka, Redis, Mysql, ES, ClickHouse, Stdin/Stdout

Dimension Table: Redis/MySQL

Format: JSON, Protobuf, Binlog

Create Function/View/Resource

Customize  Watermark

Computed Column

# Task 调度优化

Task Scheduling Optimizations

**03**

# YARN Gang 调度器

➢ 一次分配请求返回所有或者没有资源

All or nothing resources for an allocation request

➢ 为应用调度而不是为节点调度

Schedule for application instead of node

- 低延迟 (每次请求毫秒延迟)

  Low latency (RT in milliseconds per request)
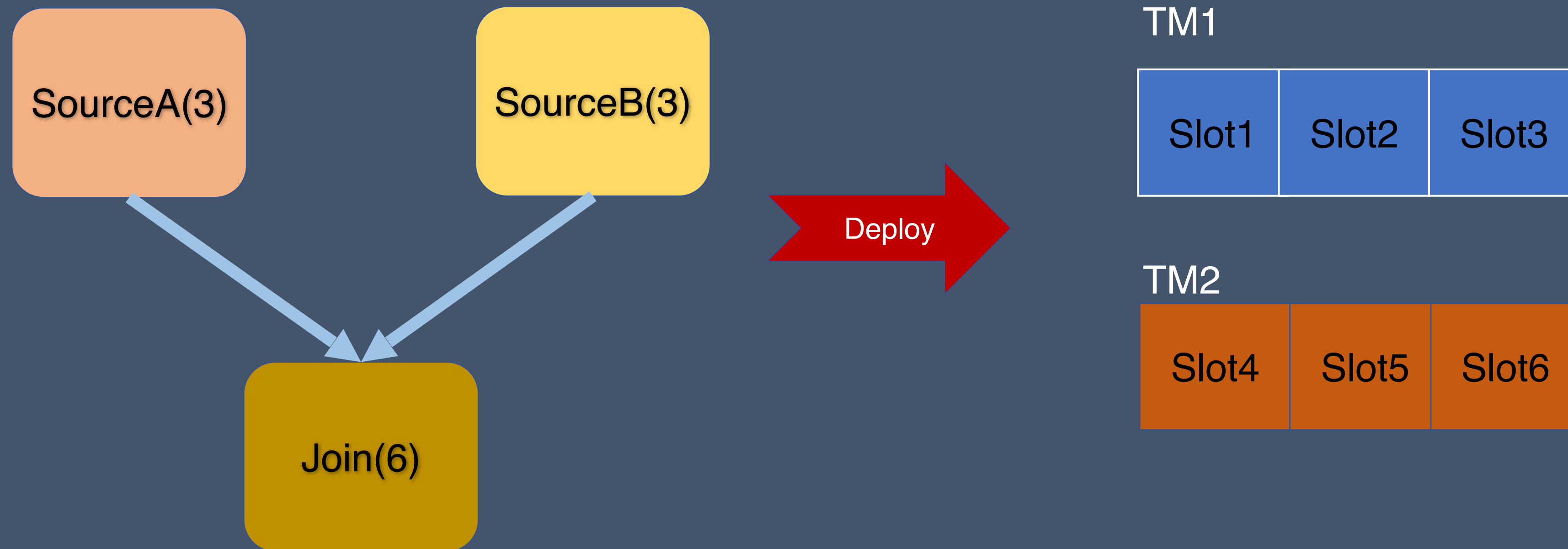
- 全局视野

  Global view

  ◆ 强约束，例如：属性，负载

    Hard constraints, e.g., Attributes, Load avg…

  ◆ 弱约束，例如：属性，负载，Container 打散，Quota 平均

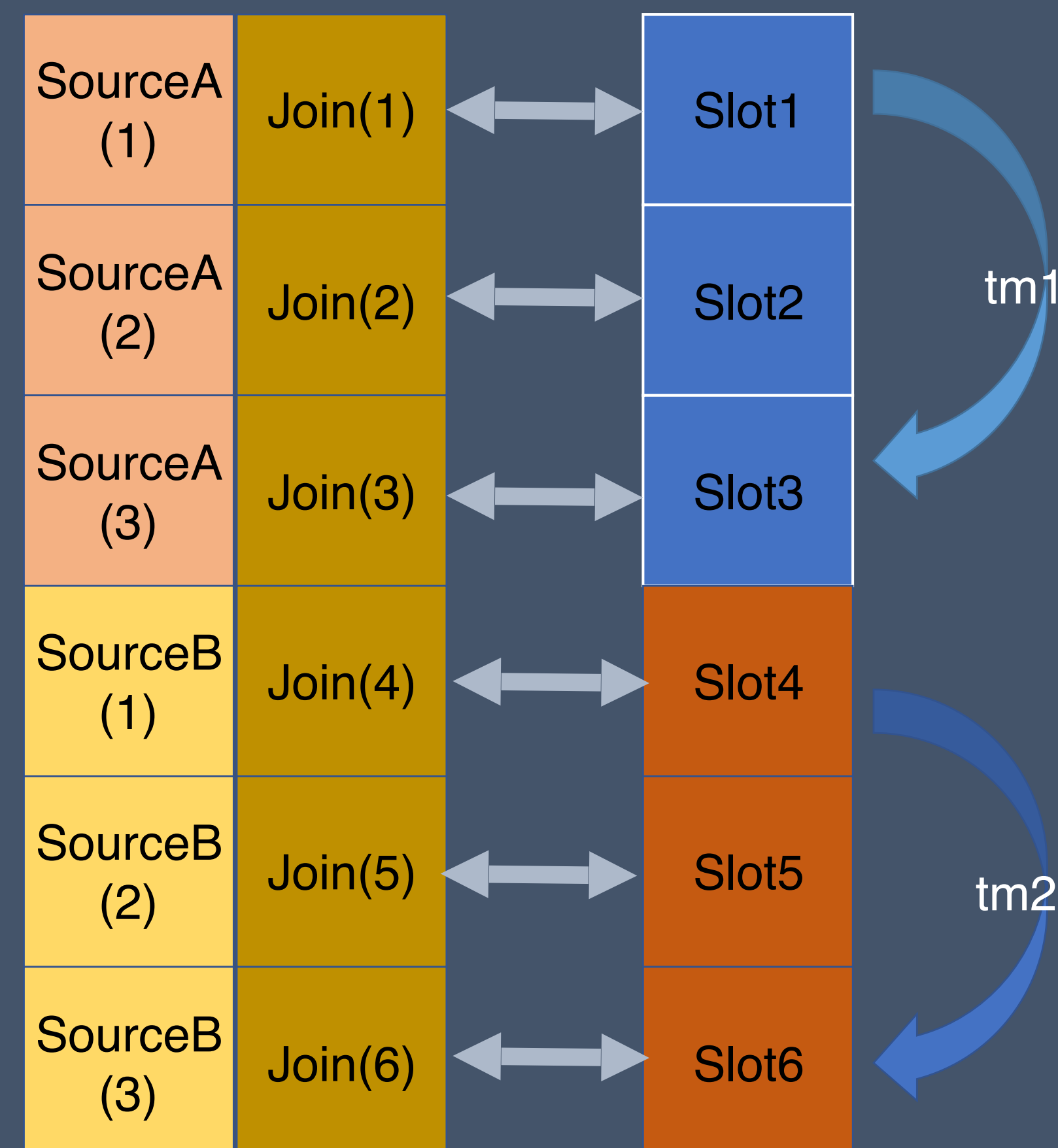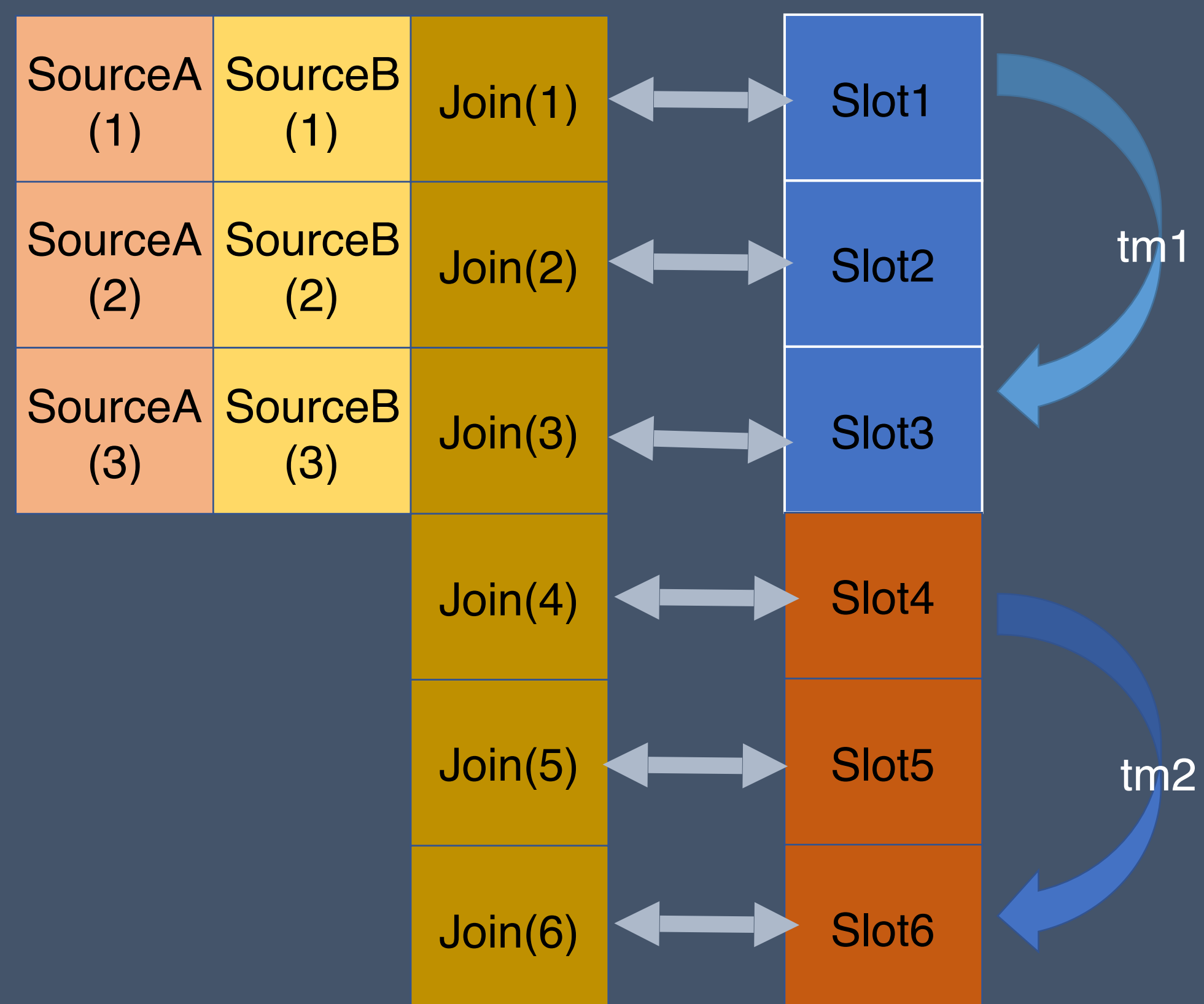    Soft constraints, e.g., Attributes, Load avg, Container Discretization, Quota avg, …

# Task 均衡调度

Task Load-Balanced Scheduling

SourceA(3)

SourceB(3)

Join(6)

Deploy

TM1

| Slot1 | Slot2 | Slot3 |

TM2

| Slot4 | Slot5 | Slot6 |

# Task 均衡调度
## Task Load-Balanced Scheduling

# Flink Docker 化

Flink Dockerization

04

# Flink Docker 化

Flink Dockeralization

使用 Python udf，包依赖管理需求强烈
To use python udf, users need dependency management

磁盘空间隔离，/tmp，日志目录
Disk space management, /tmp, log dir

镜像服务是 P2P 分发模式，加快 container 初始化速度
Image distribution service uses P2P, which accelerates the container localization

# Flink Docker 化

Flink Dockeralization

YARN 支持通过
Docker 启动 JM/TM
YARN supports starting JM and
TM using docker

剔除 Flink 提交过程中
对 HDFS 的依赖，加快
初始化速度
Remove HDFS dependencies in the
submission process and accelerate
the container localization

# Flink 其他优化

Other Flink Optimizations

**05**

# Checkpoint 使用优化

Checkpoint Usage Optimization

Checkpoint 不能跨应用，当作业重启，无法默认从上一个应用的 Checkpoint 中恢复

Checkpoint does not cross applications. When a job restarts, it can not recover from the checkpoint of the previous invocation.

作业重启默认从 Checkpoint 恢复

When a job restarts, it now by default recovers from the checkpoint of the previous invocation.

# JM 异步并发启动 TM

JobManager Starts TaskManager Asynchronously and Concurrently

串行启动 TM，速度慢
Sequentially starting TMs takes long time

慢节点
Slow nodes

异步接口启动
Start TMs asynchronously

多线程并发
Multi-threaded concurrency

超时重试，黑名单机制
Timeout retry and blacklist mechanism

# Client/Cluster 类加载机制不同

Different Class Loading Mechanisms between Client and Cluster



Client 端类加载机制与 Cluster 端不同
The class loader used by the client is different from the default one used by the cluster

修改 client 端支持 child-first/parent-first 类加载机制可配置
The class loader used by the client supports configurable child-first/parent-first class loading mechanisms.

FLINK FORWARD

# JM 资源可配置

Configurable JobManager Resources



jm vcore 默认为 1，jm mem
与 tm mem 保持一致

The vcore limit of the job manager is 1 by
default and the memory limit of the job
manager is same as the task manager's

jm vcore, jm mem 支持单独配置

The vcore and memory limits of the job manager can be
configured independently

FLINK
FORWARD

# Flink 在字节跳动未来展望

Outlook of Flink at ByteDance

06

# Q&A

# 致谢

Thanks

感谢数据开发套件、计算架构团队的小伙伴，以上
分享内容是我们共同努力的成果

Thanks my colleagues of the data development kit team and the compute infrastructure team @ ByteDance, the above content is the result of our joint efforts