



GOPHER CHINA 2020

中国 上海 / 2020-11.21-22

探探分布式存储的实践

彭亮



Agenda

- Why tantan db
- What is tantan db
- Golang
- Roadmap

Why do we need tantan db?

GOPHER CHINA 2020

中国 上海 / 2020-11.21-22

tantan

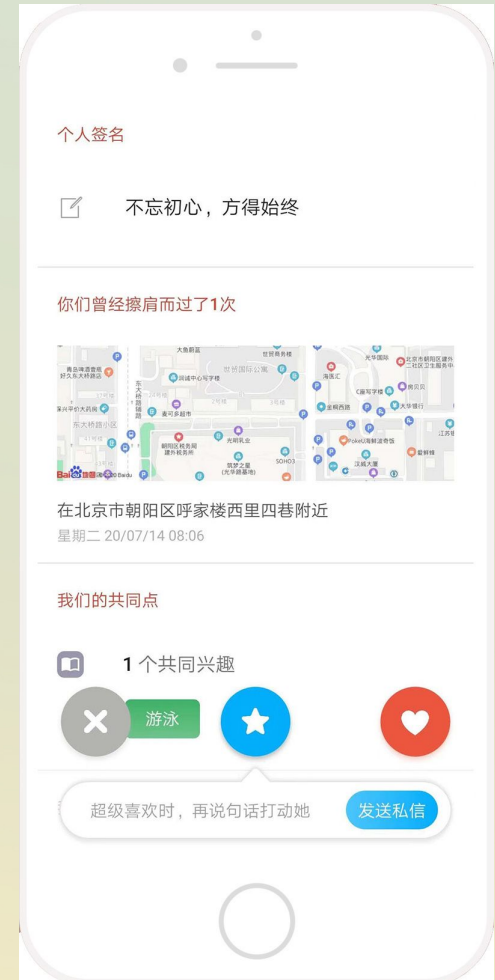
左滑无感，右滑喜欢



破冰利器，附近动态



擦肩而过，回眸一笑



tantan

业务特性

- 大数据量
- 数据快速增长
- 低时延

数据特性

- 分区
- 聚集性
Clustered

需求

架构师

- 分布式
- 可用性
- 扩展性
- 定制化

使用者

- SQL
- BASE vs ACID
- PACELC

DBA

- 存储成本
- 运维友好
- 容灾能力

开源方案

NoSQL

- SQL限制
- 存储成本
- 定制化

NewSQL

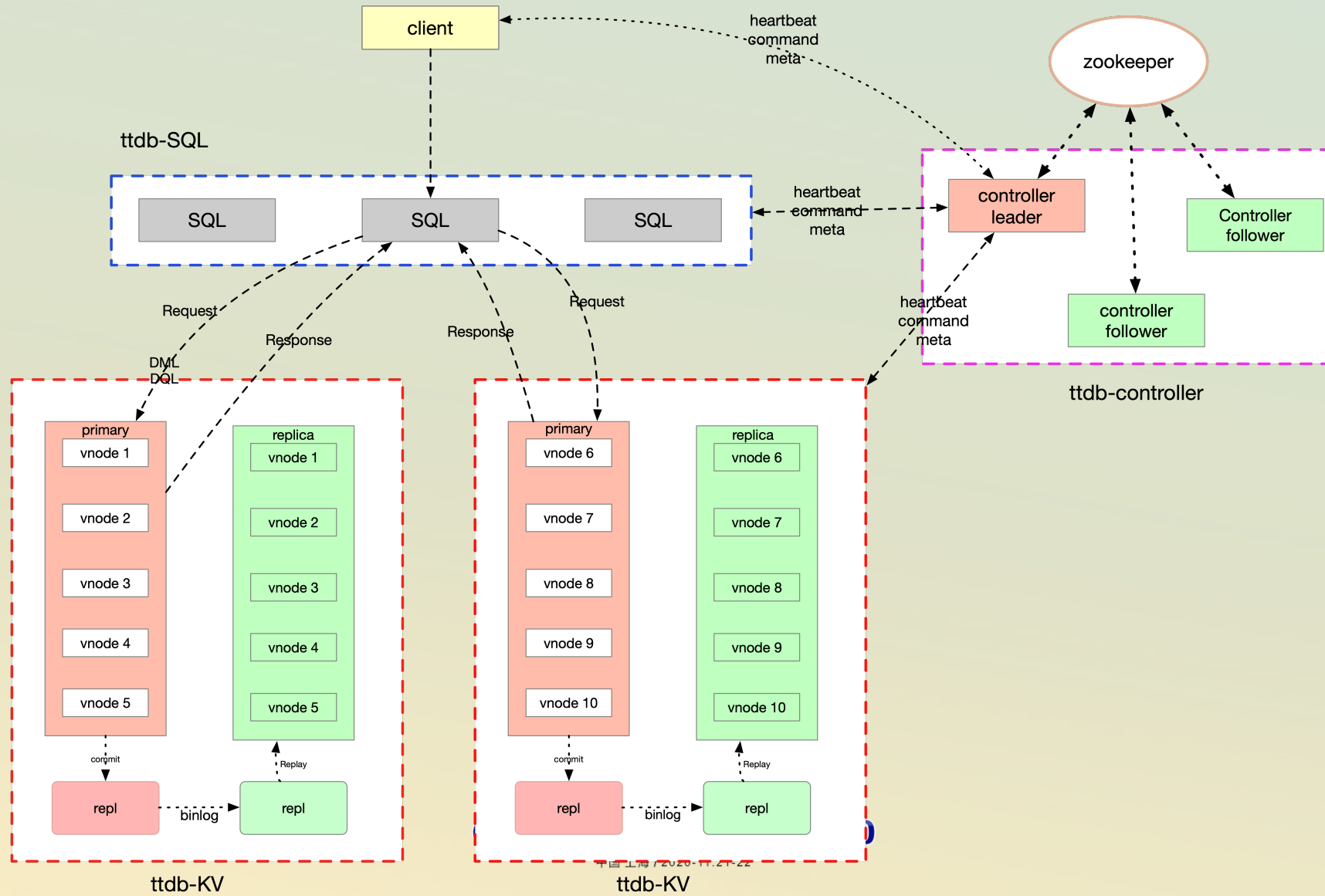
- 存储成本
- 强一致性
- ACID
- 延迟
- 定制化

What is tantan db(ttldb)?

GOPHER CHINA 2020

中国 上海 / 2020-11.21-22

架构

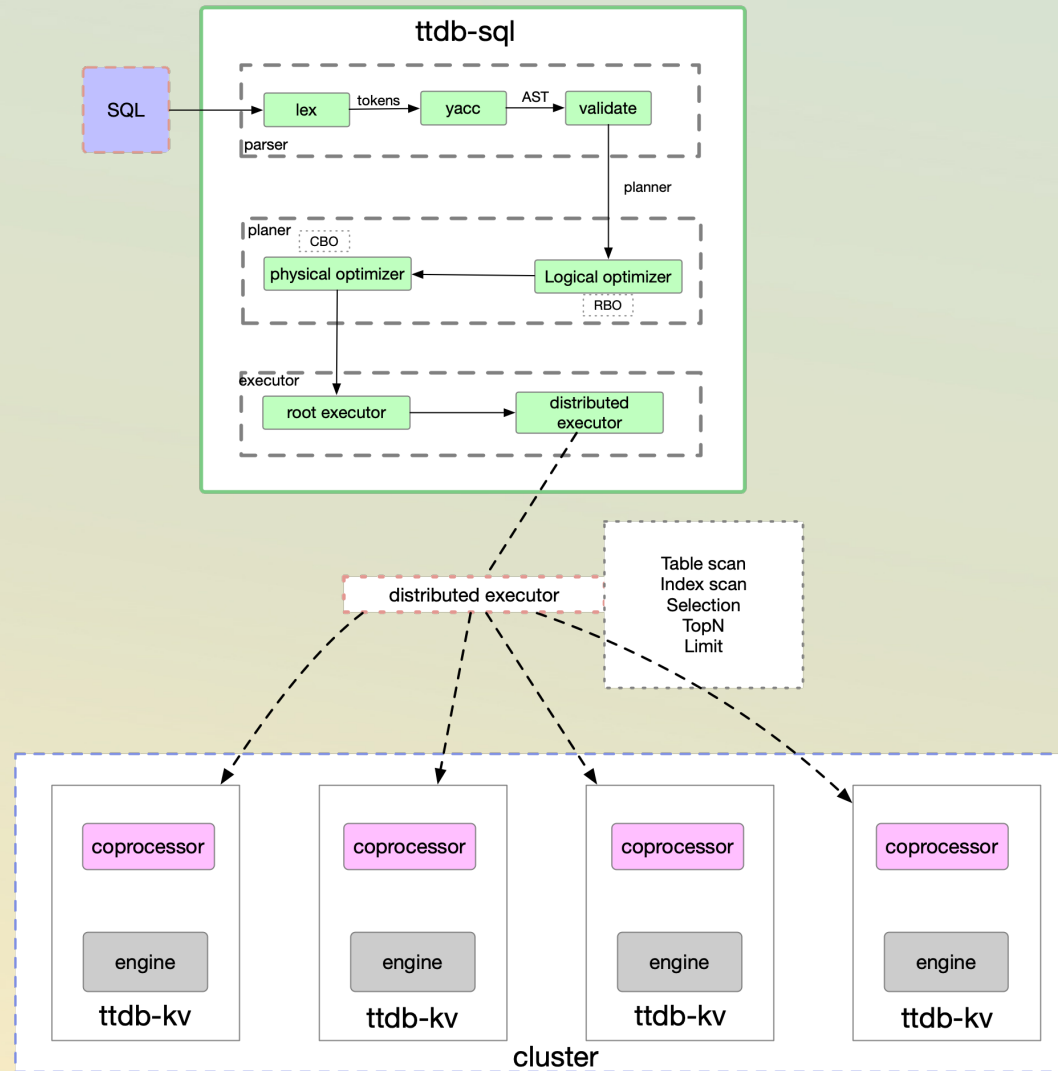


SQL

GOPHER CHINA 2020

中国 上海 / 2020-11.21-22

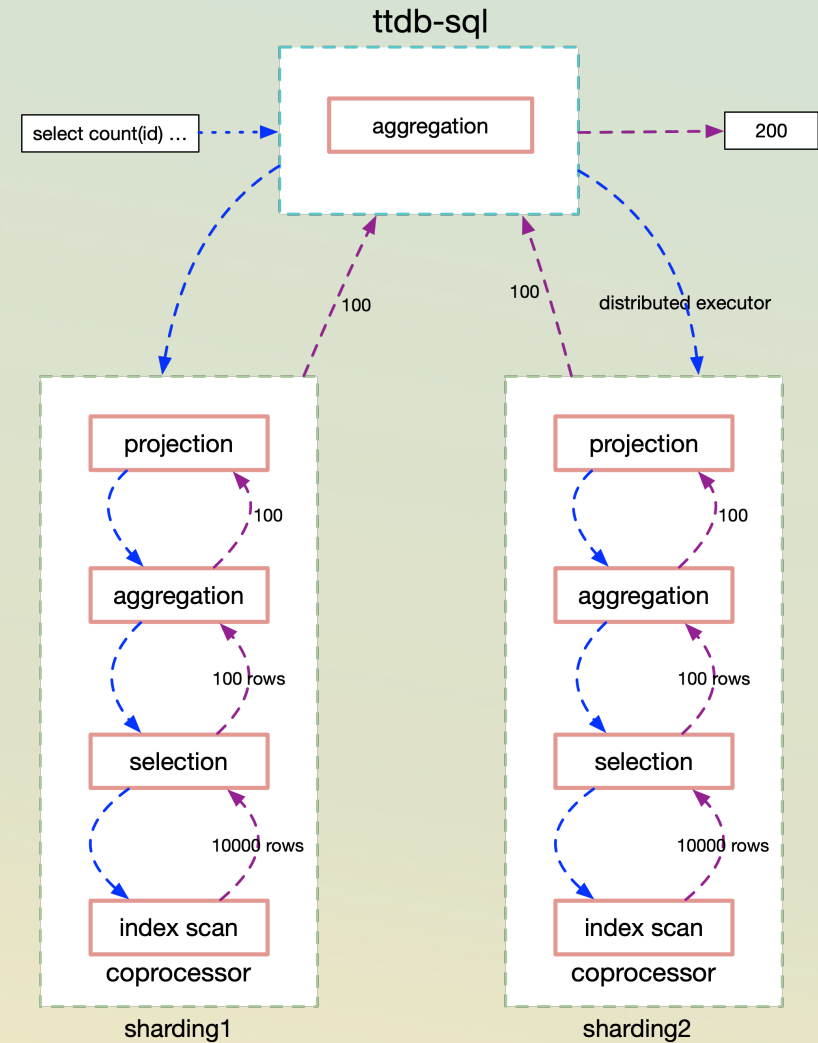
SQL 执行



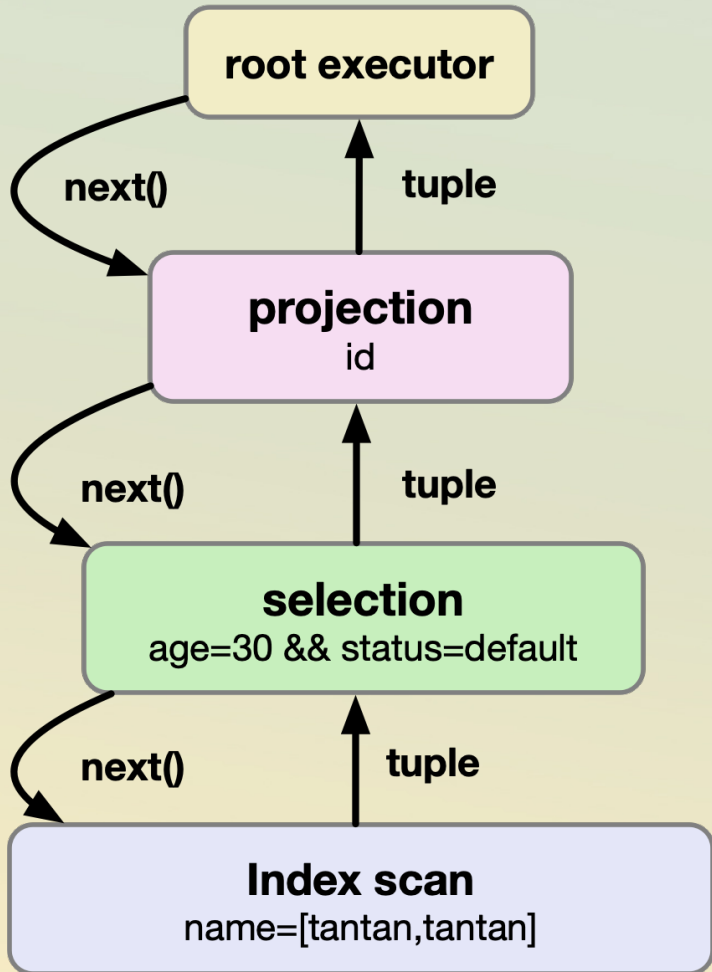
SQL 优化器

Rule based optimizer

- 列裁剪(prune columns)
- 谓词下推(push down predicate)
- 聚合下推(push down aggregation)
- topN下推 (push down topN)



Volcano model executor



Index scan

id	name	age	status
5	ken	30	default
6	tantan	30	default
7	tantan	30	default
8	tantan	30	banned
9	tantan	31	default

name=tantan

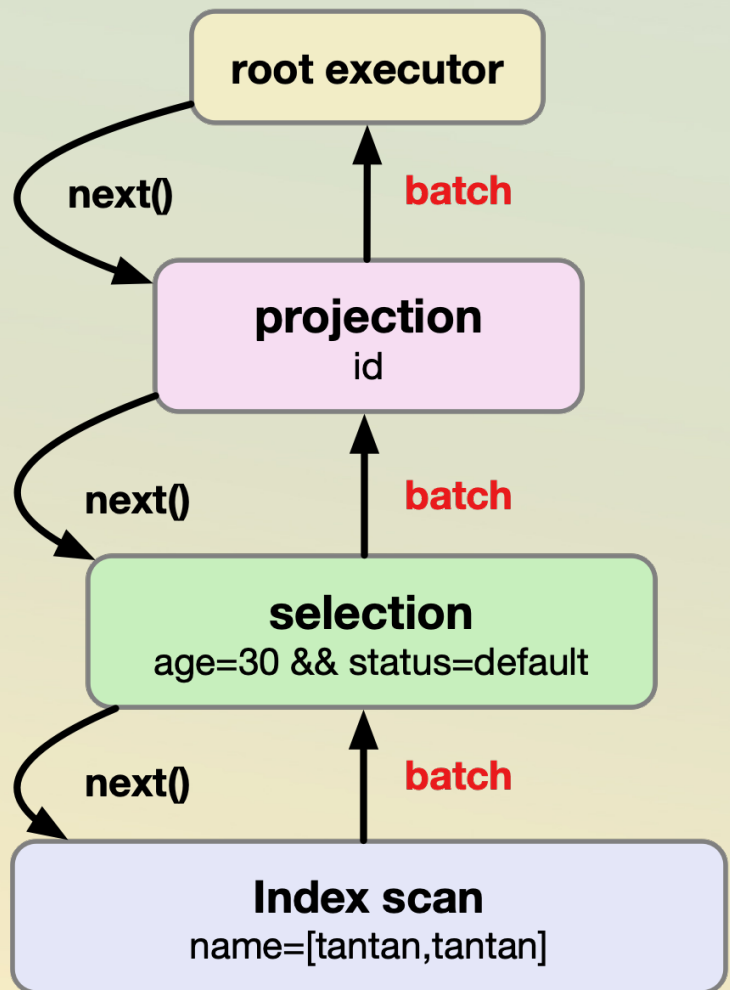
selection

	id	age	status
row	6	30	default
row	7	30	default
row	8	30	banned
row	9	31	default

age=30

status=default

Vectorization model executor



Index scan

id	name	age	status
5	ken	30	default
6	tantan	30	default
7	tantan	30	default
8	tantan	30	banned
9	tantan	31	default

name=tantan

selection

column	column	column
id	age	status
6	30	default
7	30	default
8	30	banned
9	31	default

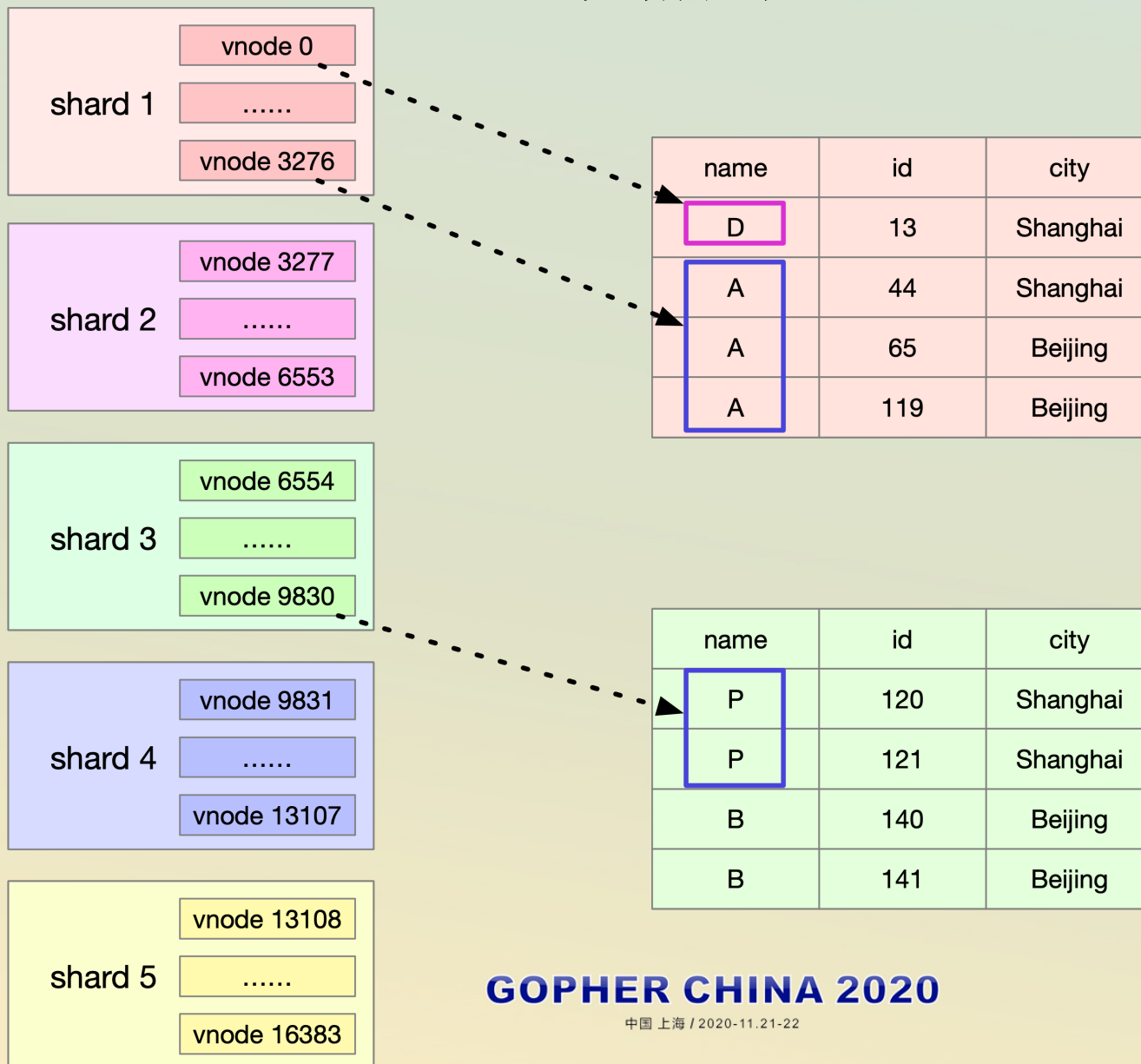
age=30 status=default

数据分片

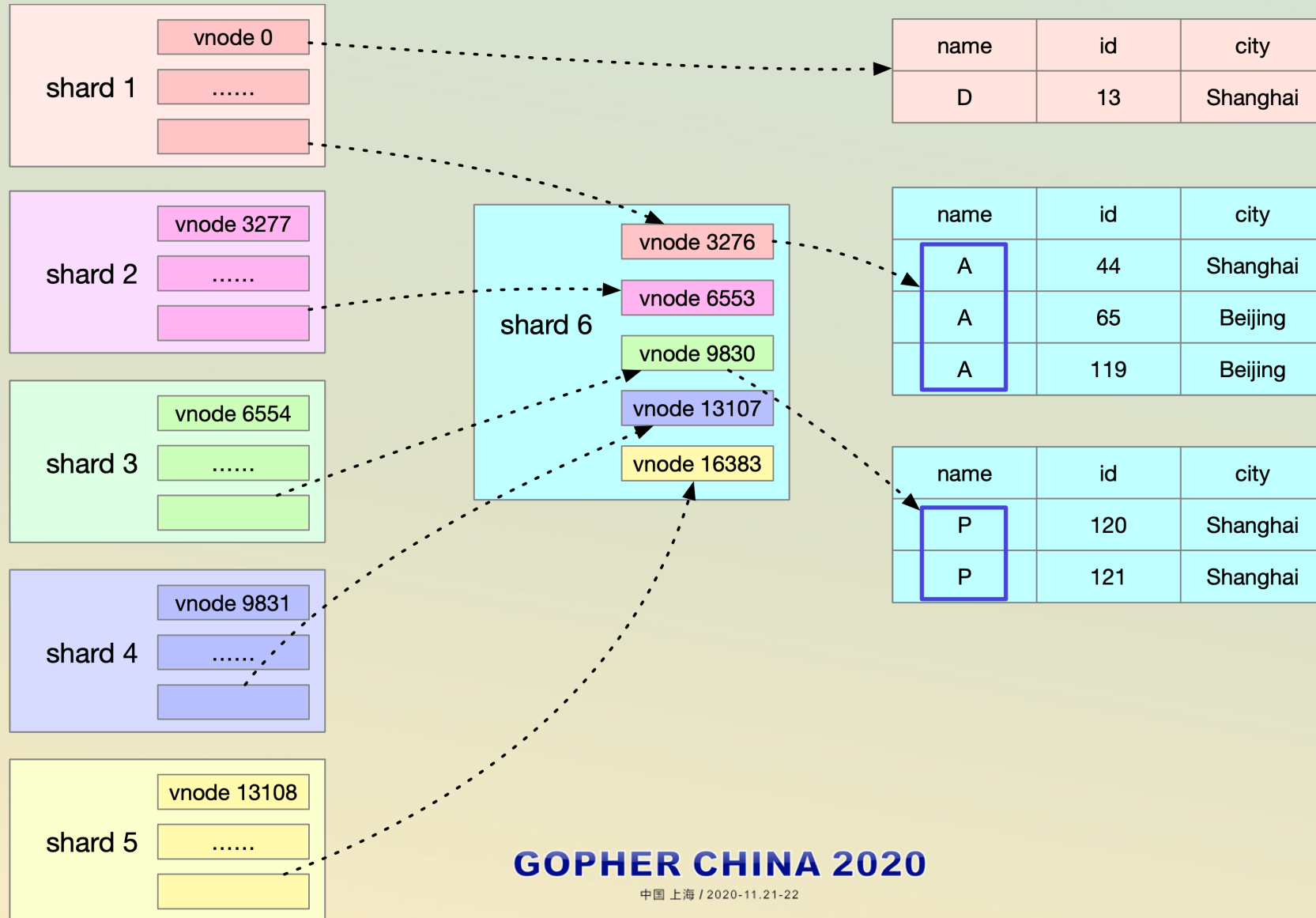
GOPHER CHINA 2020

中国 上海 / 2020-11.21-22

数据分片



扩容



高可用

GOPHER CHINA 2020

中国 上海 / 2020-11.21-22

故障检测

中心化

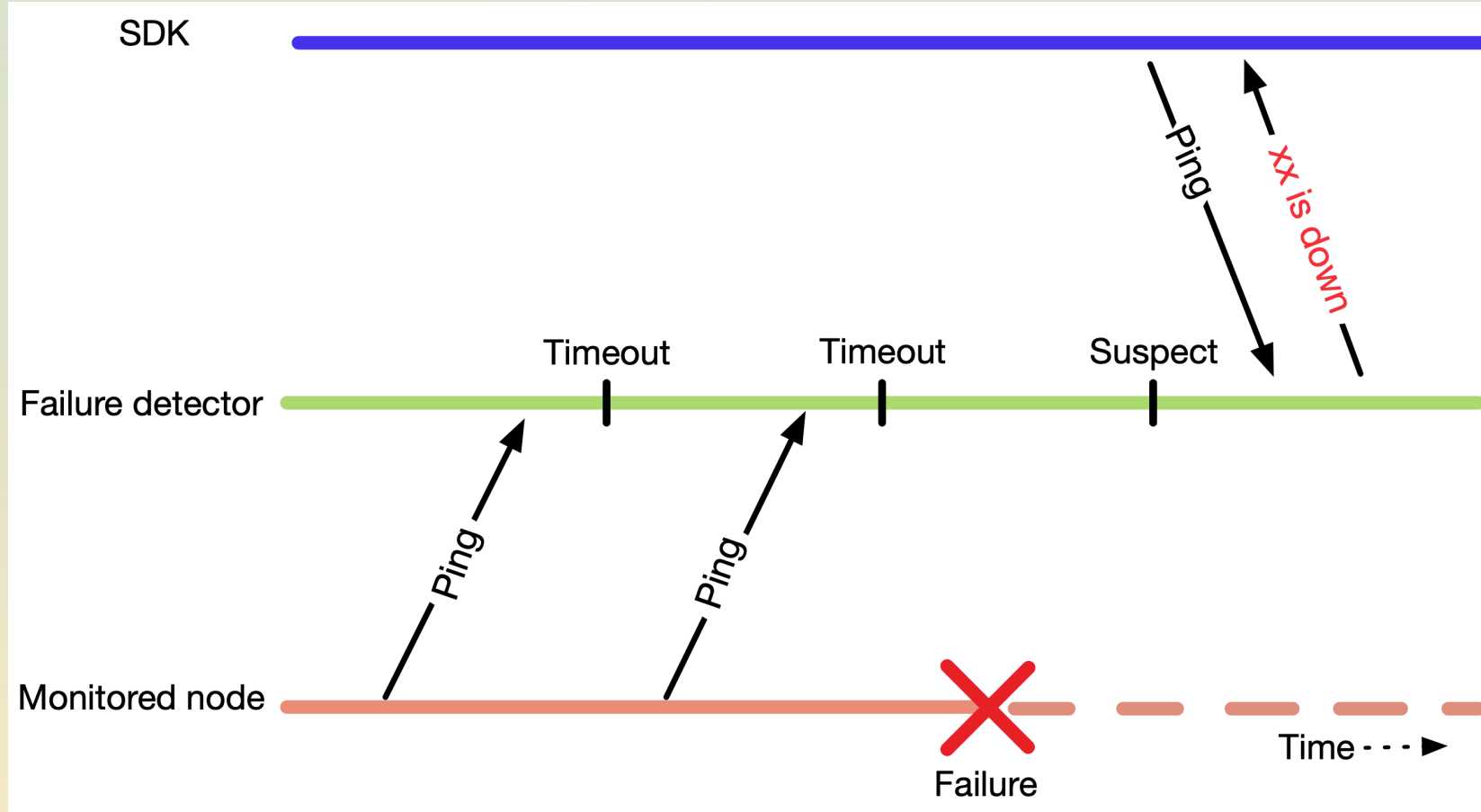
- Timeout(binary or accrual)
vs false positive
- Network partition

无中心化

Gossip and Failure Detection

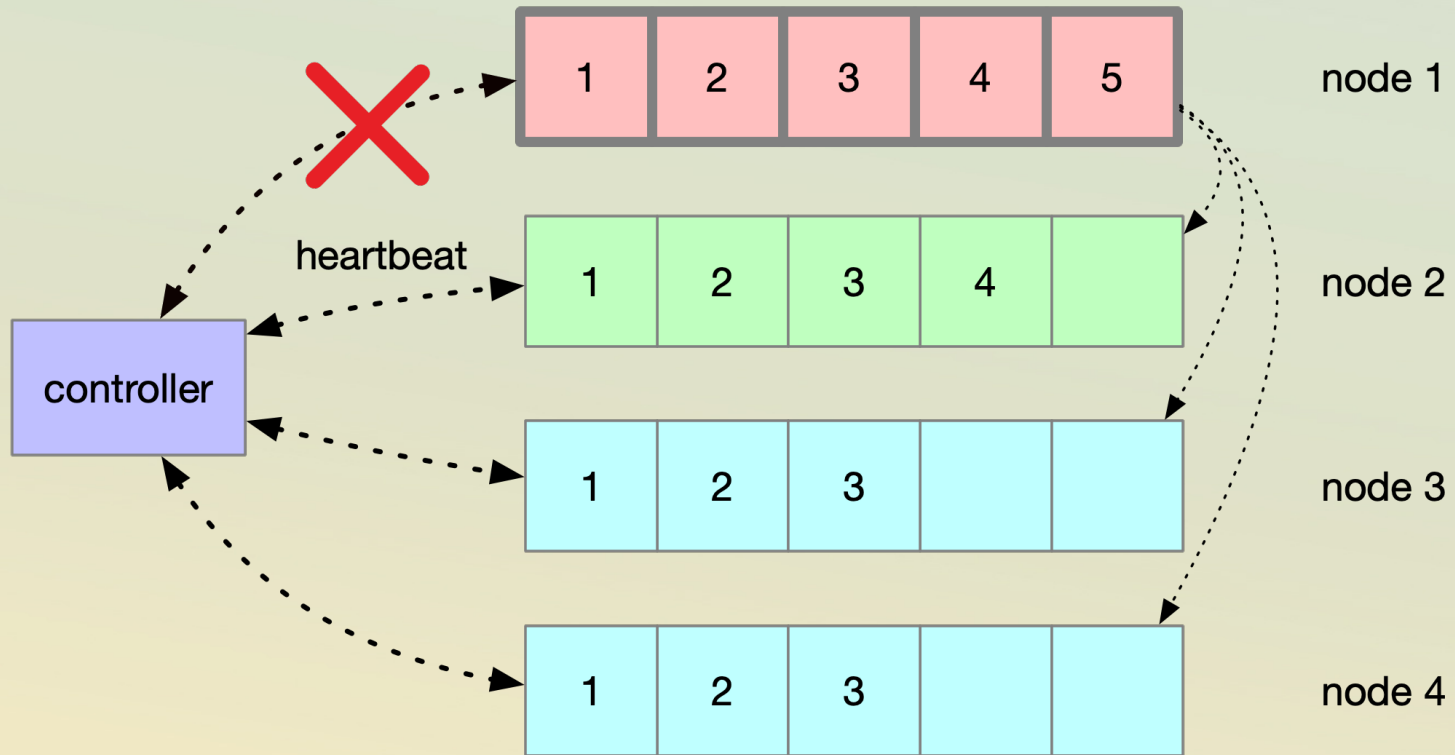
- P2P
- Outsourced
- Convergence

故障检测

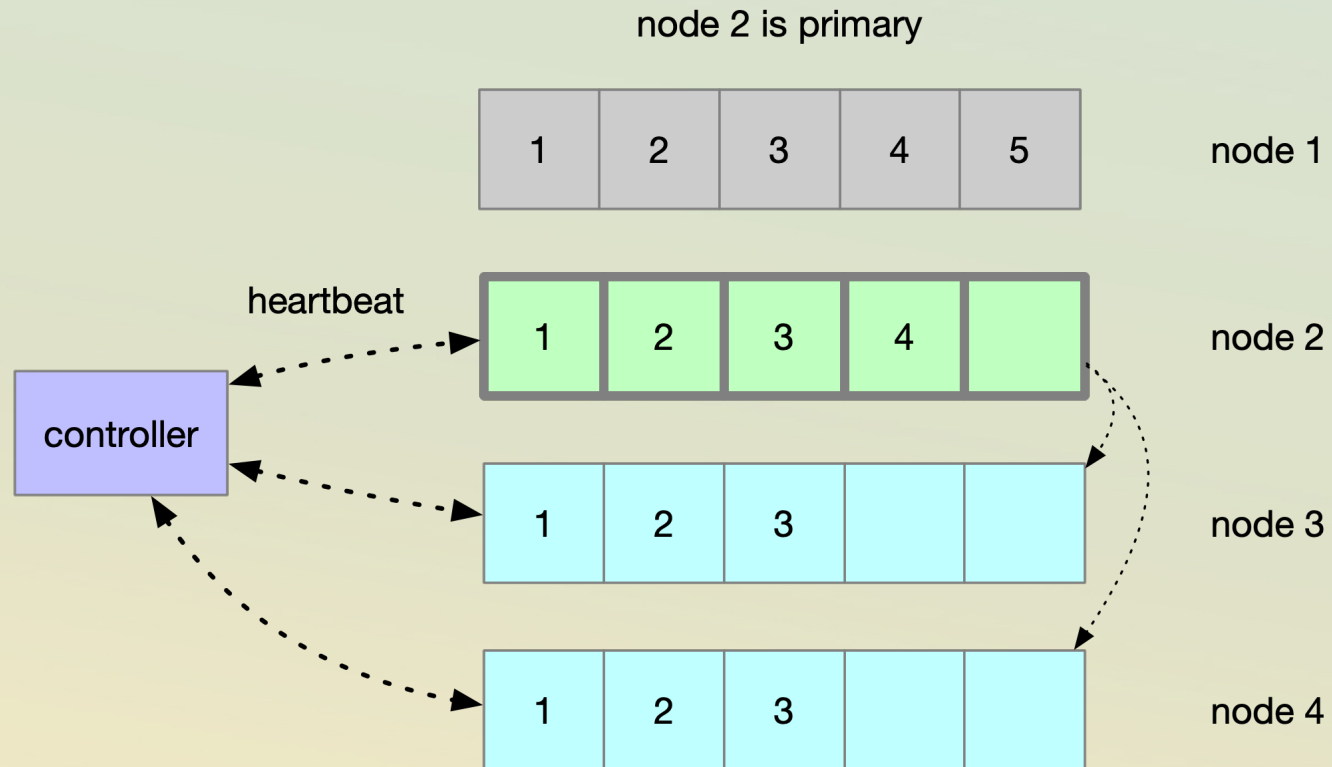


故障恢复

node 1 is down



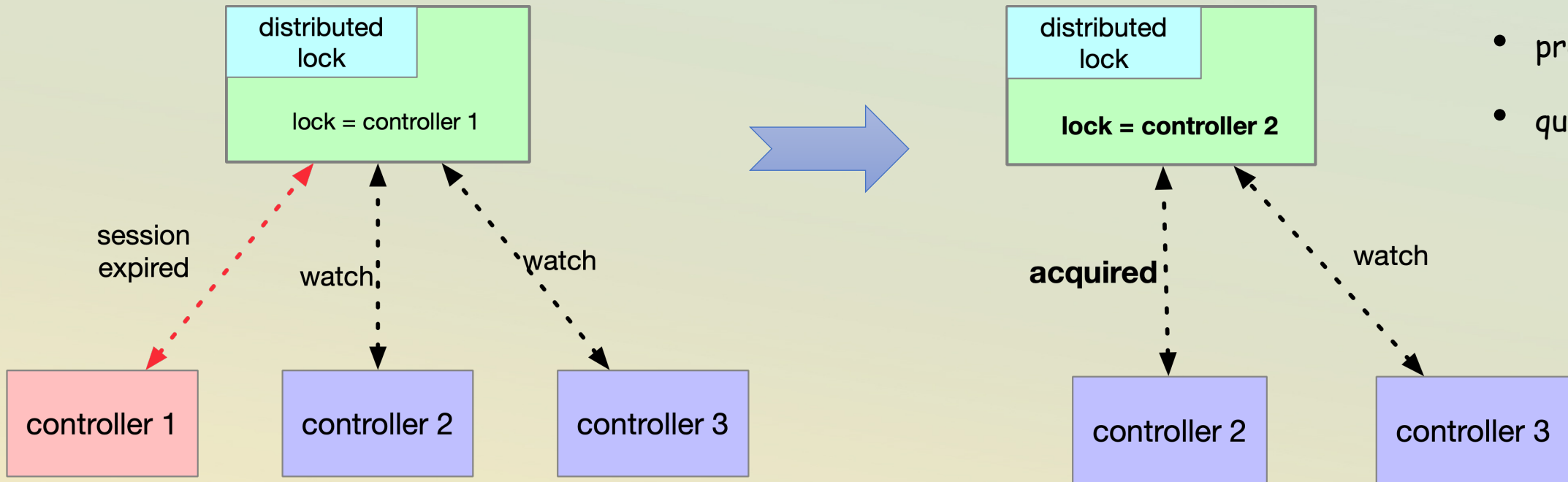
故障恢复



Stateful

- up-to-date

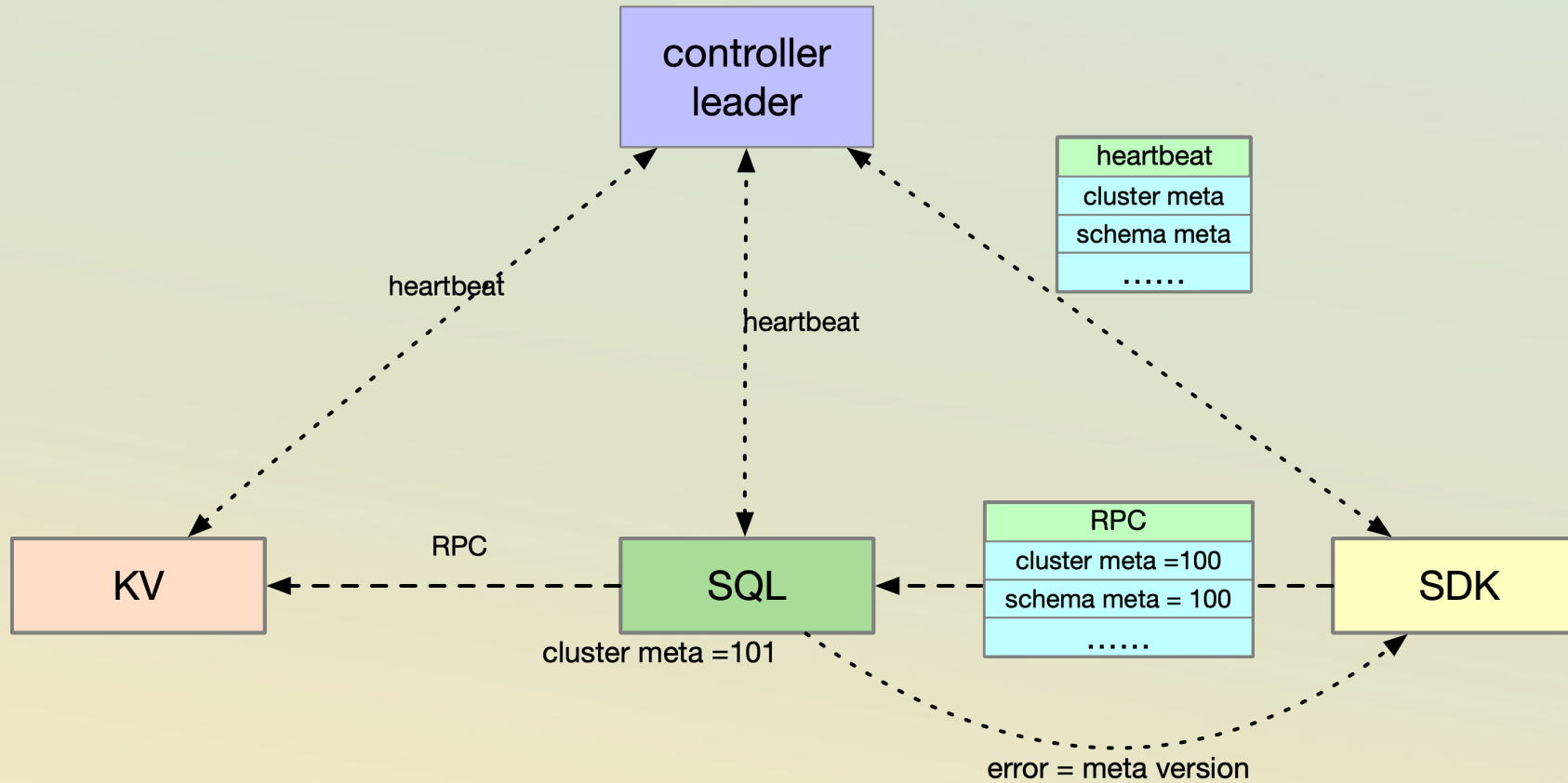
故障恢复



Stateless

- preemptive
- queue

重新配置



故障恢复了?

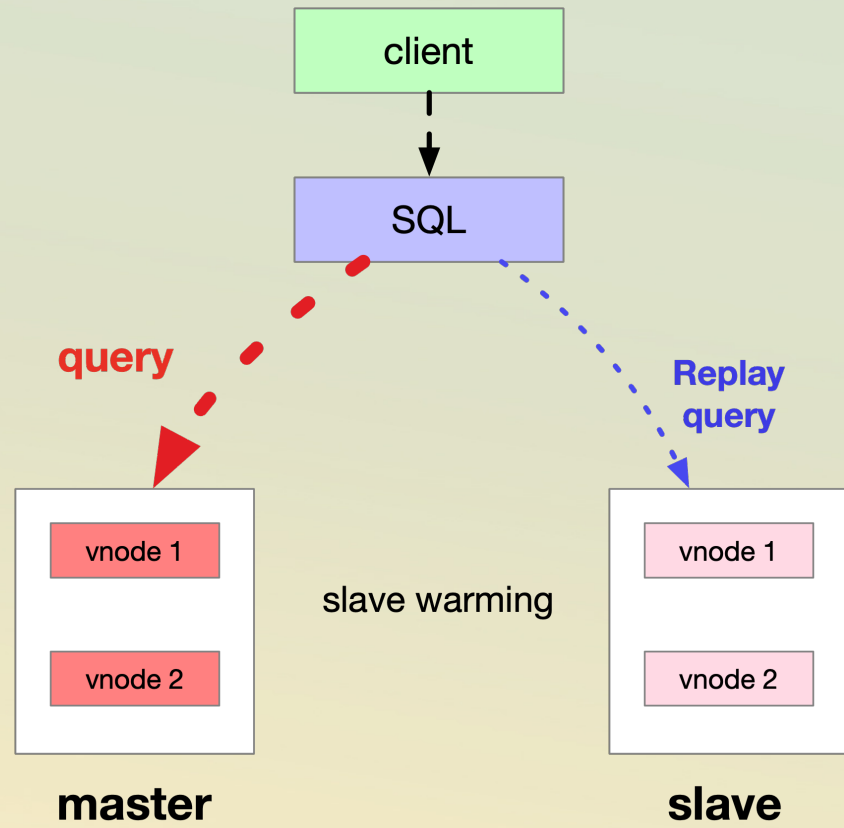
cache miss



GOPHER CHINA 2020

中国 上海 / 2020-11.21-22

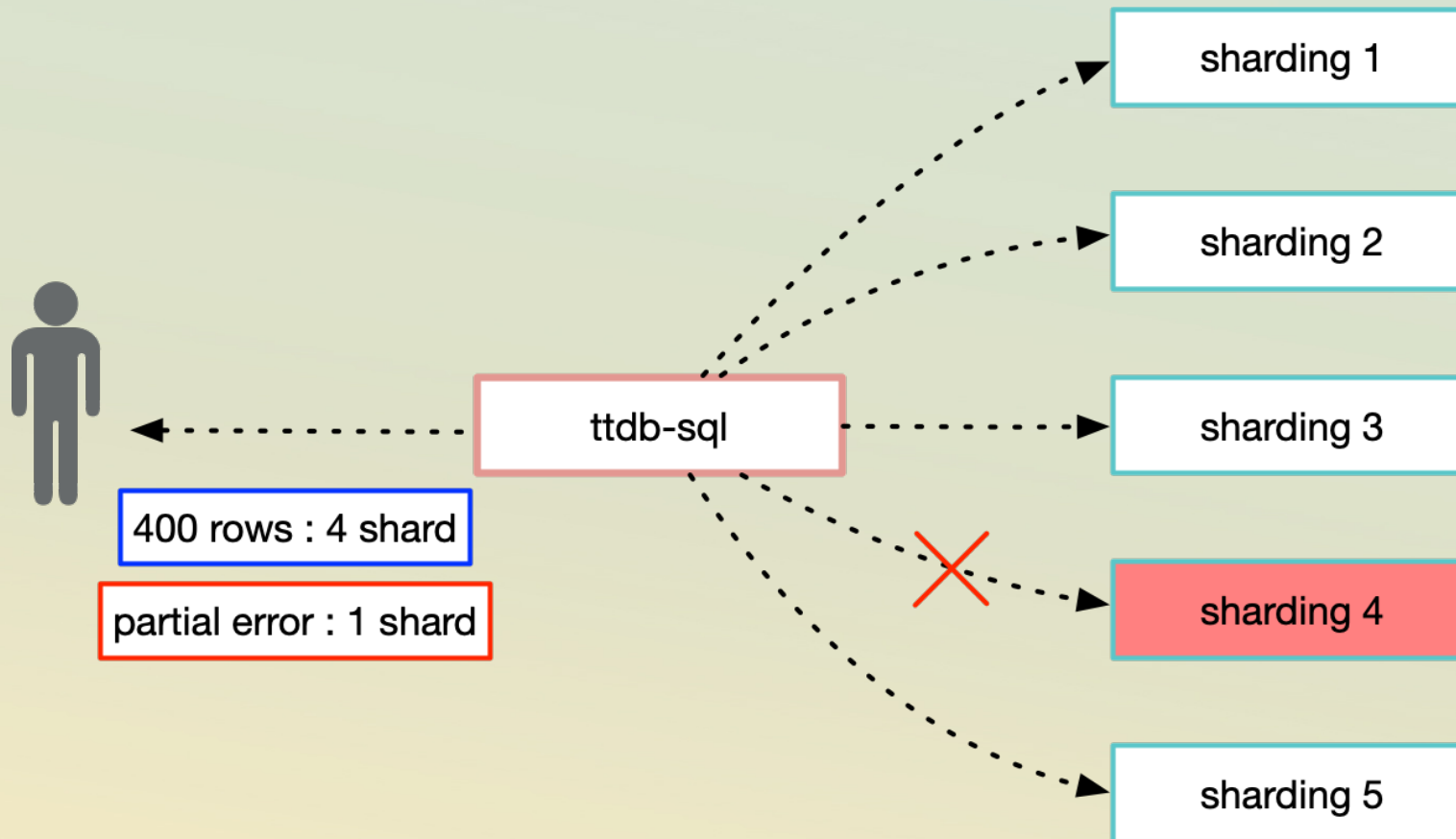
故障恢复了?



cache hit



结果与产出 一致性与可用性



业务应用

GOPHER CHINA 2020

中国 上海 / 2020-11.21-22

擦肩而过

user_id	other_user_id	location	time
user_A	user_B	142.654, 71.892	2020-01-25 11:00:00
user_A	user_C	142.617, 71.903	2020-01-25 11:05:00
user_A	user_D	142.598, 72.061	2020-01-25 11:07:00

```
ttdb> explain select user_id, other_user_id from location.passby where user_id=1 and other_user_id in (2,3,4) and count > 2;
plan          type          info
-----
RootReader    root
  Projection  dist          location.passby.user_id, location.passby.other_user_id
    Selection dist          gt(location.passby.count, 2)
      TableScan dist          table:passby, range:[1 2,1 2], [1 3,1 3], [1 4,1 4]

ttdb>
```

擦肩而过



selection latency	P99	P90	Avg
ttdb	3.4ms	2.1ms	960us
cassandra	19ms	8.6ms	4ms

	1TB存储的数据行数
ttdb	41M
cassandra	7.8M

Golang

GOPHER CHINA 2020

中国 上海 / 2020-11.21-22

Why golang

- 学习成本
- 开发效率
- 标准库
- Goroutine
- GC

Goroutine 限制和泄漏

context, timeout, channel, select

- <https://github.com/uber-go/goleak>
- `runtime.NumGoroutine()`
- `pprof/goroutine`

GC 优化

减少对象分配

- Reuse
- Stack
- Preallocation

降低scan成本

- Pointers

Roadmap

GOPHER CHINA 2020

中国 上海 / 2020-11.21-22

Todo

- 开源 (2021)
- 物理优化
- 同步复制
- Range 分区
- 生态



GOPHER CHINA 2020

中国 上海 / 2020-11.21-22



Thanks

