



百度在线数据仓库Palo

——开源架构解读及应用

百度大数据部 牟宇航

2017.12

QCon

全球软件开发大会

成为软件技术专家的 必经之路

[北京站] 2018

2018年4月20-22日 北京·国际会议中心

7折 购票中, 每张立减2040元

团购享受更多优惠



识别二维码了解更多



极客时间

重拾极客精神·提升技术认知

下载极客时间App

获取有声IT新闻、技术产品专栏，每日更新



扫一扫下载极客时间App

2018 Geekbang > InfoQ
极客邦 极客邦

AiCon

全球人工智能与机器学习技术大会

助力人工智能落地

2018.1.13 - 1.14 北京国际会议中心

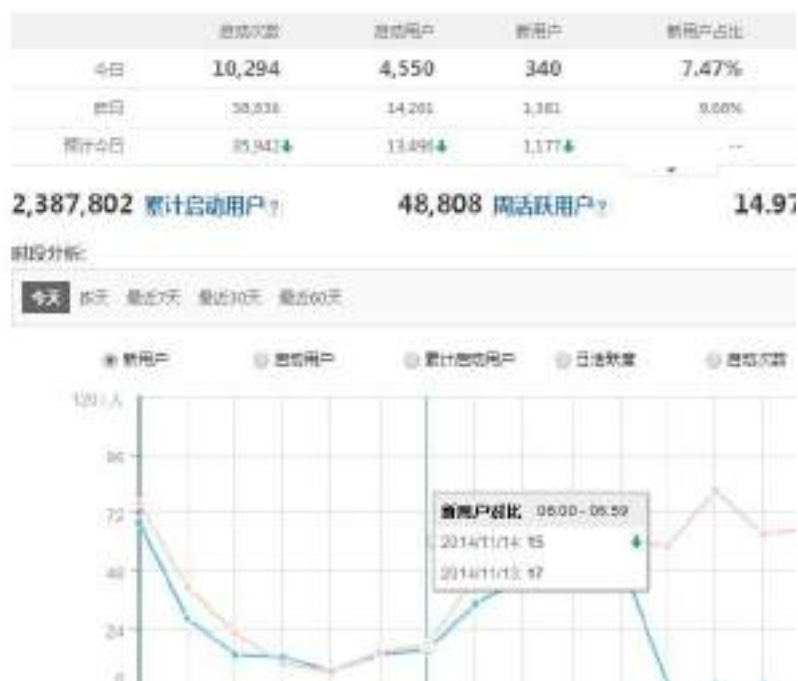


扫描关注大会官网

场景一

• 某在线报表业务

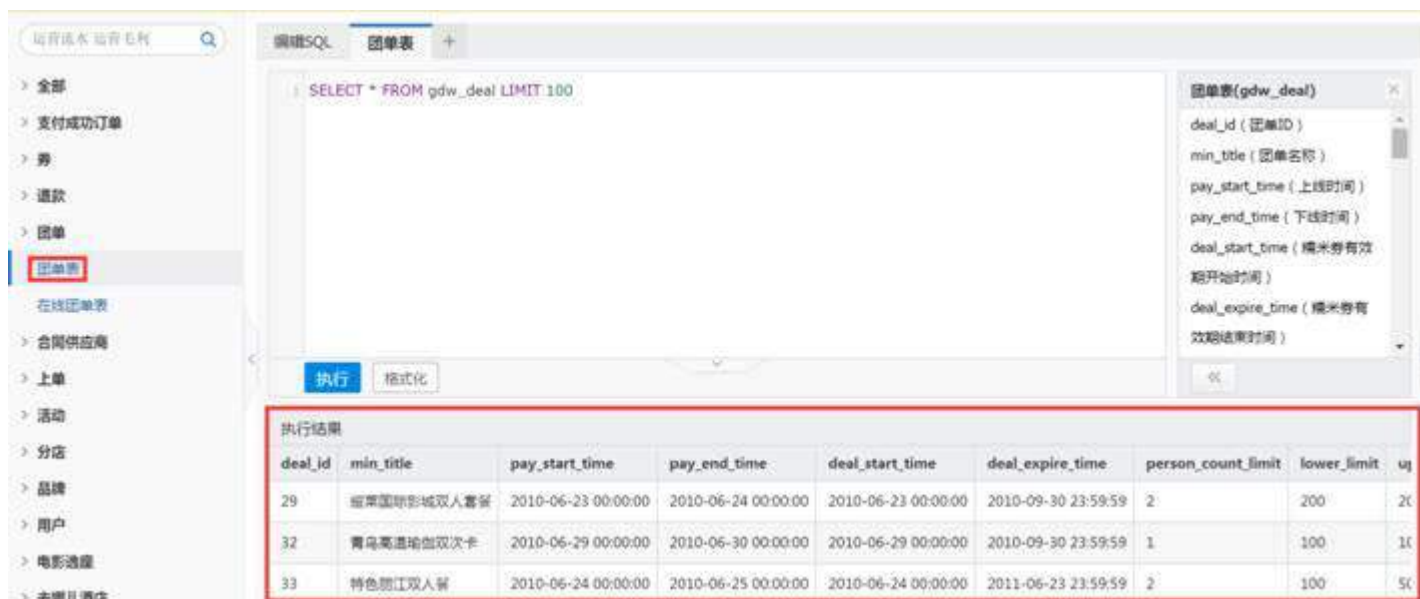
- 为网站站长提供流量分析，网站分析，受众分析等多种分析服务
- 300+表，数据清洗结构化后百TB+，单日增量1TB+
- 查询峰值QPS 2000+，日查询量千万级
- 一致性（会话内单调一致性、更新一致性）
- 导入5分钟一次
- 查询平均延时30+ms



场景二

• 某业务数据集市

- 集运营、业务分析、订单管理、会员管理、客户关系管理等数十个管理分析平台一体的综合数据平台
- 100+主题视图、10-100TB
- 标准SQL，Ad-Hoc（即席查询），秒级分析



The screenshot displays a data analysis tool interface. On the left is a navigation menu with categories like '全部', '支付成功订单', '券', '退款', '团单', '团单表', '在线团单表', '合同供应商', '上单', '活动', '分店', '品牌', '用户', '电影选座', and '去哪儿酒店'. The '团单表' (Group Deal Table) is selected. The main area shows a SQL editor with the query: `SELECT * FROM gdw_deal LIMIT 100`. Below the editor are buttons for '执行' (Execute) and '格式化' (Format). On the right, a '团单表(gdw_deal)' schema is shown with fields: deal_id (团单ID), min_title (团单名称), pay_start_time (上线时间), pay_end_time (下线时间), deal_start_time (糯米券有效期开始时间), and deal_expire_time (糯米券有效期结束时间). Below the editor is a table titled '执行结果' (Execution Results) with the following data:

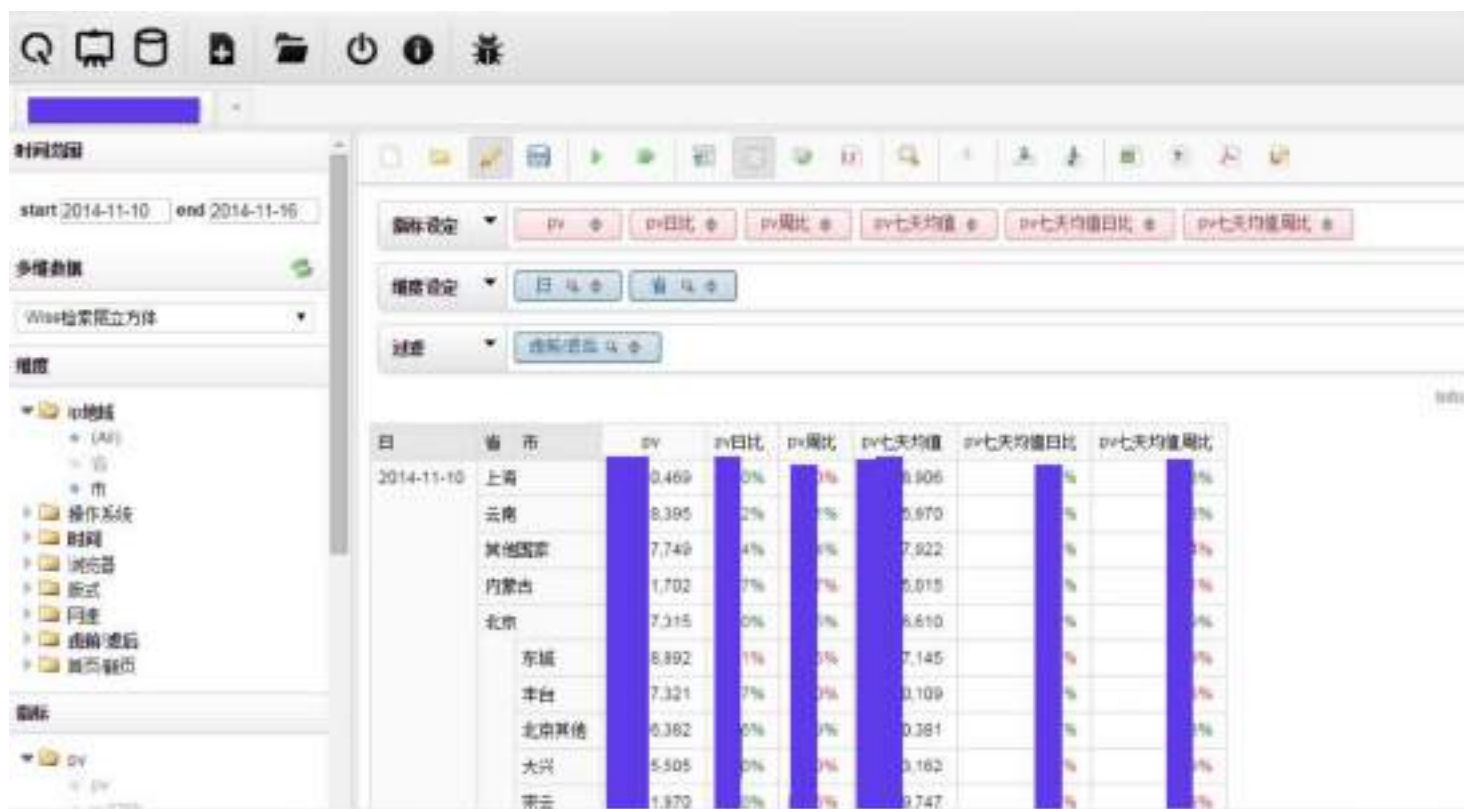
deal_id	min_title	pay_start_time	pay_end_time	deal_start_time	deal_expire_time	person_count_limit	lower_limit	uj
29	南京国际影城双人套餐	2010-06-23 00:00:00	2010-06-24 00:00:00	2010-06-23 00:00:00	2010-09-30 23:59:59	2	200	20
32	青岛高速瑜血双次卡	2010-06-29 00:00:00	2010-06-30 00:00:00	2010-06-29 00:00:00	2010-09-30 23:59:59	1	100	10
33	特色丽江双人餐	2010-06-24 00:00:00	2010-06-25 00:00:00	2010-06-24 00:00:00	2011-06-23 23:59:59	2	100	50



场景三

• 某在线多维分析平台

- 100+表，最大单表50+维度列、10+指标列，任意组合，秒级分析
- 10-100TB



场景

- 以前

- 报表 : Hadoop + MySQL
- 分析 : Hadoop + Hive

- 现在

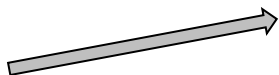
在线报表



在线多维分析



即席查询



在线数据仓库 → **Palo**

Palo在百度大数据技术栈的位置

App Client

(Mobile App, Browser App, Desktop App, Service API)

Web / App Server

(Business Logic)

Transaction Processing

Storage System
(NoSQL, Table, Object)

Database System
(OldSQL, NewSQL)

Analytical Processing

Online Data Serving

Stream Data Processing

Batch Data Processing

Data Bus / File System

(Bigpipe, Kafka / HDFS, AFS)

Palo在友商技术栈的应用

品友第三代数据分析平台的技术栈



在线数据仓库——OLAP

- **Online Analytical Processing**

- Online vs. Offline (Interactive vs. Batch)
- Analytical Processing vs. Transactional Processing

	OLTP	OLAP
面向应用	日常交易处理	明细查询，分析决策
访问模式	简单小事务，操作少量数据	复杂聚合查询，查询大量数据
数据	当前最新数据	历史数据
数据规模	GB	TB ~ PB
数据更新	实时更新	批量更新
数据组织	满足3NF	反范式，星型模型

OLAP-商业产品

产品	简介	技术特点	收购情况
Netezza	2000年在美国成立 Netezza TwinFin	<ul style="list-style-type: none"> ✓ 软硬一体机 ✓ 采用FPGA数据过滤代替索引 	2010年9月20日，IBM出资17.8亿美元收购
Greenplum	2003年在美国成立 Greenplum Database	<ul style="list-style-type: none"> ✓ 行存 + 列存 ✓ Shared-Nothing集群 	2010年7月6日，EMC出资3亿美元收购
Vertica	2005年在美国成立 Vertica Analytic Database	<ul style="list-style-type: none"> ✓ 列存 ✓ Shared-Nothing集群 	2011年2月，HP出资3.5亿美元收购
Aster Data	2005年在美国成立 nCluster	<ul style="list-style-type: none"> ✓ SQL-MapReduce ✓ Shared-Nothing集群 	2011年7月6日，Teradata出资2.63亿美元收购
ParAccel	2005年在美国成立 PADB	<ul style="list-style-type: none"> ✓ 列存 + 自适应压缩 ✓ Shared-Nothing集群 	2013年Actian出资1.5亿美元收购，Redshift宣称使用ParAccel

Vendor and Appliance	Memory (GB)	Total Cores	Compression	User Storage (TB, Compressed)	List Price
EMC Greenplum Data Computing Appliance	768	48	4 to 1	144	\$2,000,000
IBM PureData System for Analytics N1001-010	n/a	112	4 to 1	128	\$1,599,000
Microsoft SQL Server 2012 Parallel Data Warehouse ¹	2,304	144	5 to 1	340	\$1,569,970
Oracle Exadata Database Machine X3-2	2,048	128	10 to 1	450	\$13,580,000
Teradata Data Warehouse Appliance 2690	768	96	4 to 1	146	\$1,168,000

OLAP-开源社区



Apache Impala (incubating) is the open source, native analytic database for Apache Hadoop. Impala is shipped by Cloudera, MapR, Oracle, and Amazon.

Follow us on Twitter at [@ApacheImpala](https://twitter.com/ApacheImpala)



Lightning-fast cluster computing



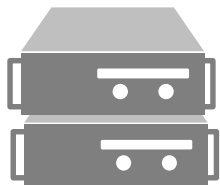
Druid is a high-performance, column-oriented, distributed data store.

百度OLAP功能需求

- 基本需求
 - high availability
 - Scalability
 - High Performance
- 特化需求
 - Consistent View
 - Monotonic Consistency
 - Heterogeneous Storage
 -

Palo定位

低成本



1/10 ~ 1/100 Cost

线性扩展



100~200节点 / 1000 TB

支持云化部署



高可用



99.9999 % Uptime

高查询性能



10W QPS/ 100GB/s

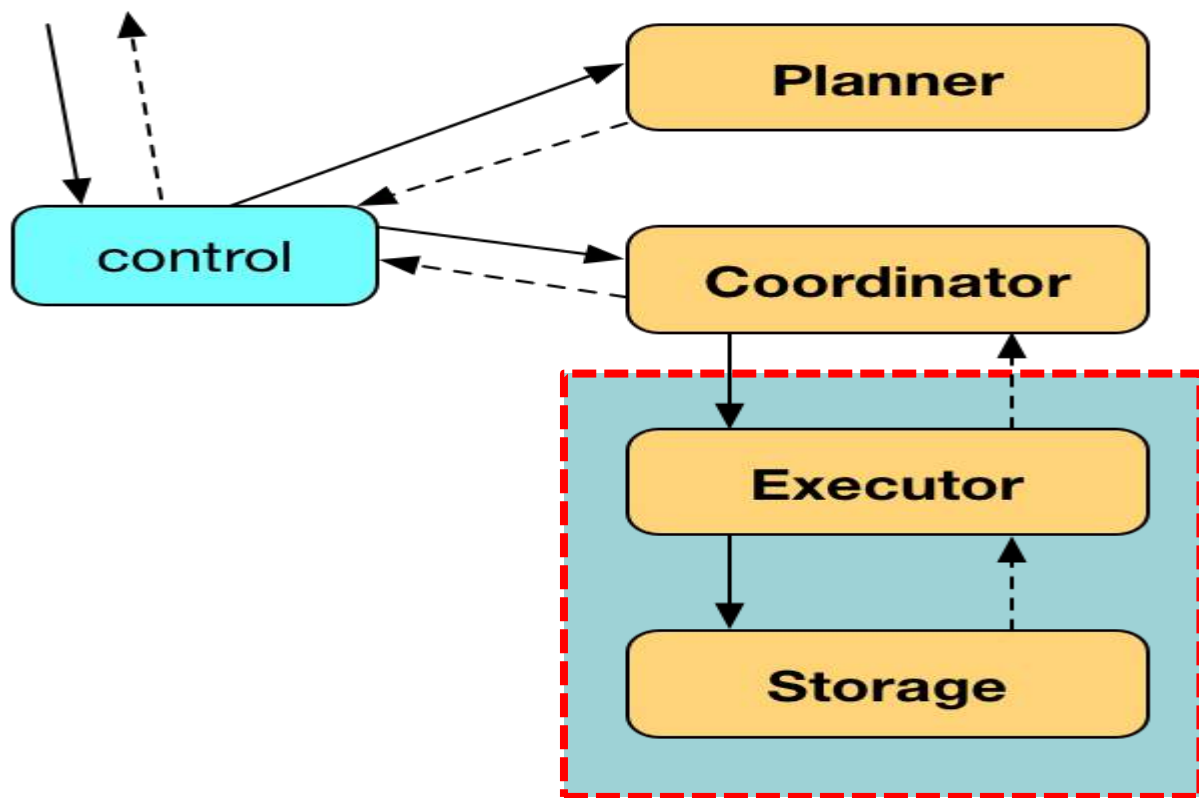
高加载性能



10 TB / Hour

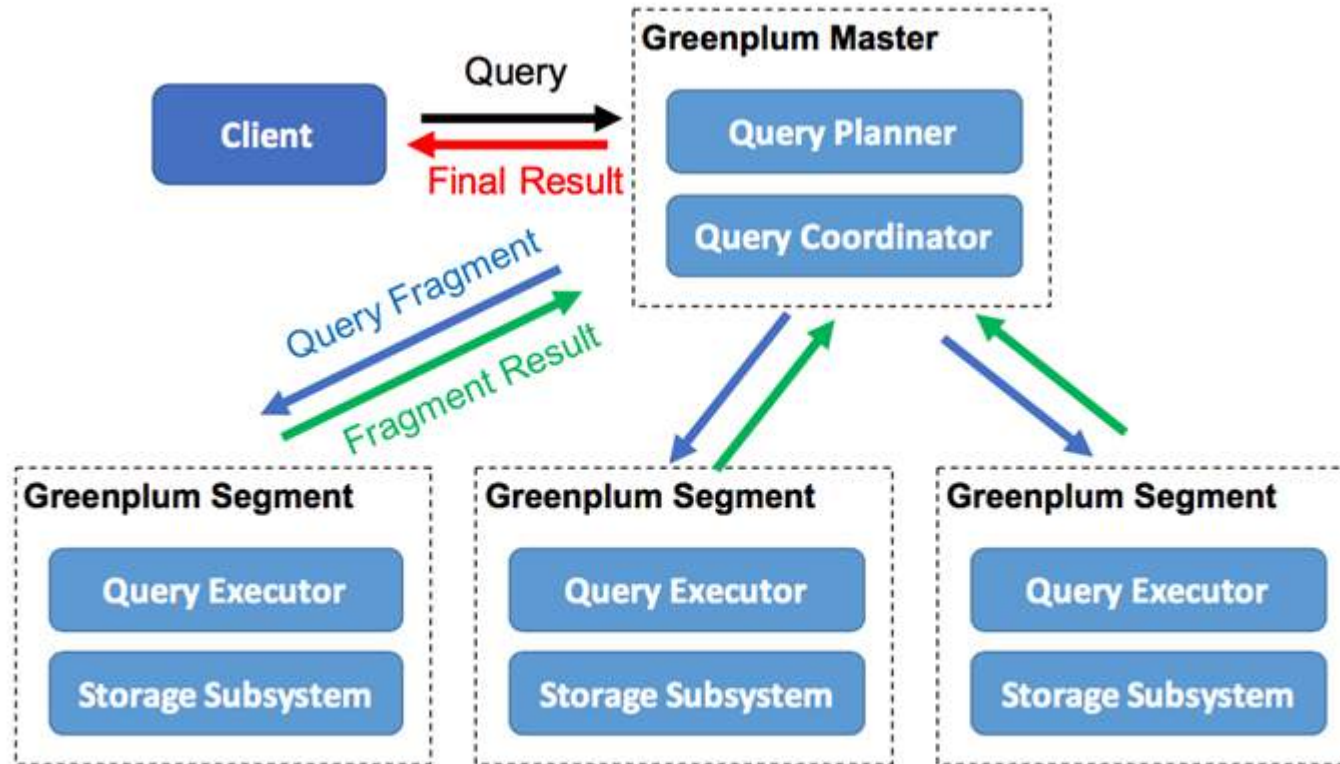
Palo整体架构

- MPP通用构件



中心架构 (Centric Architecture)

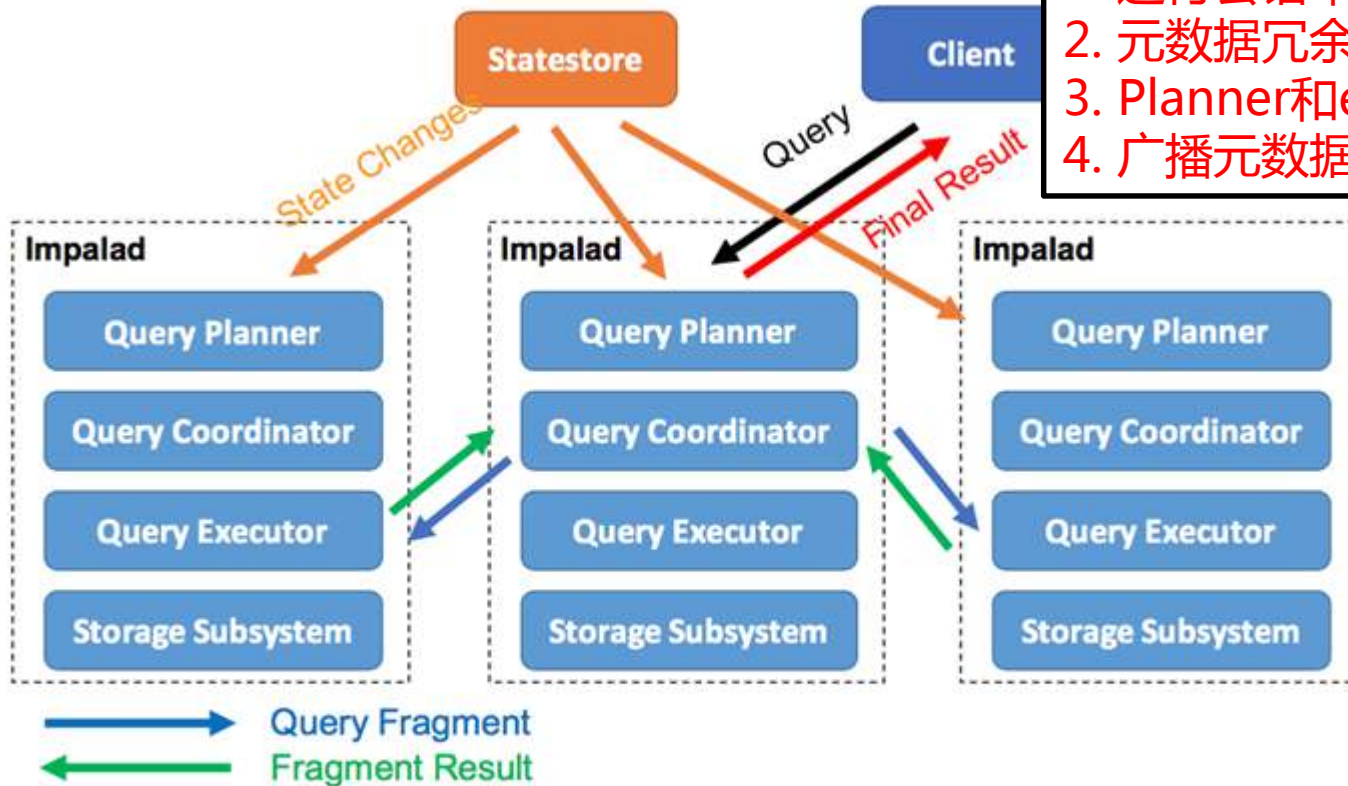
Master may be bottleneck



Centric architecture — Greenplum.

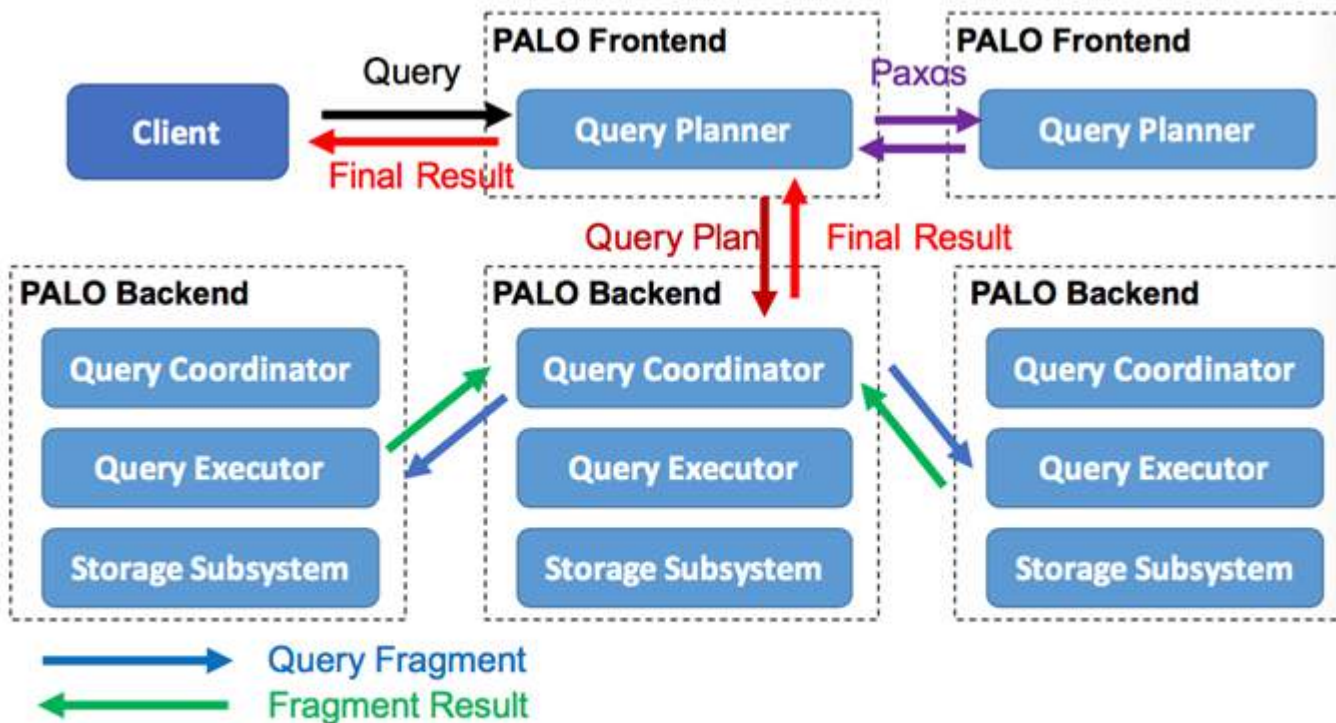
对称架构 (Symmetric Architecture)

1. 违背会话单调一致性
2. 元数据冗余
3. Planner和executor争内存
4. 广播元数据，可能有网络瓶颈



Symmetric architecture — Impala.

混合架构 (Hybrid Architecture)

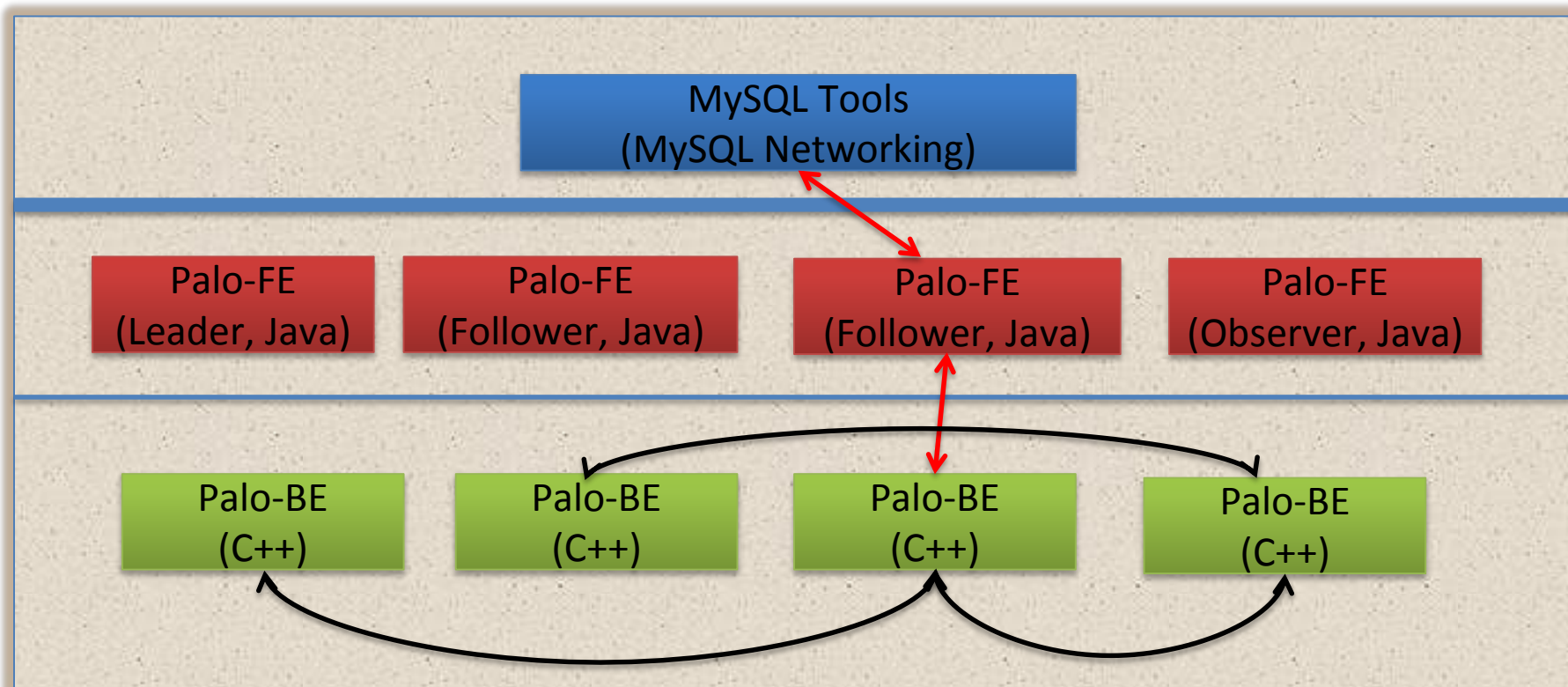


Hybrid architecture — PALO.

Palo整体架构

- **Hybrid Architecture**

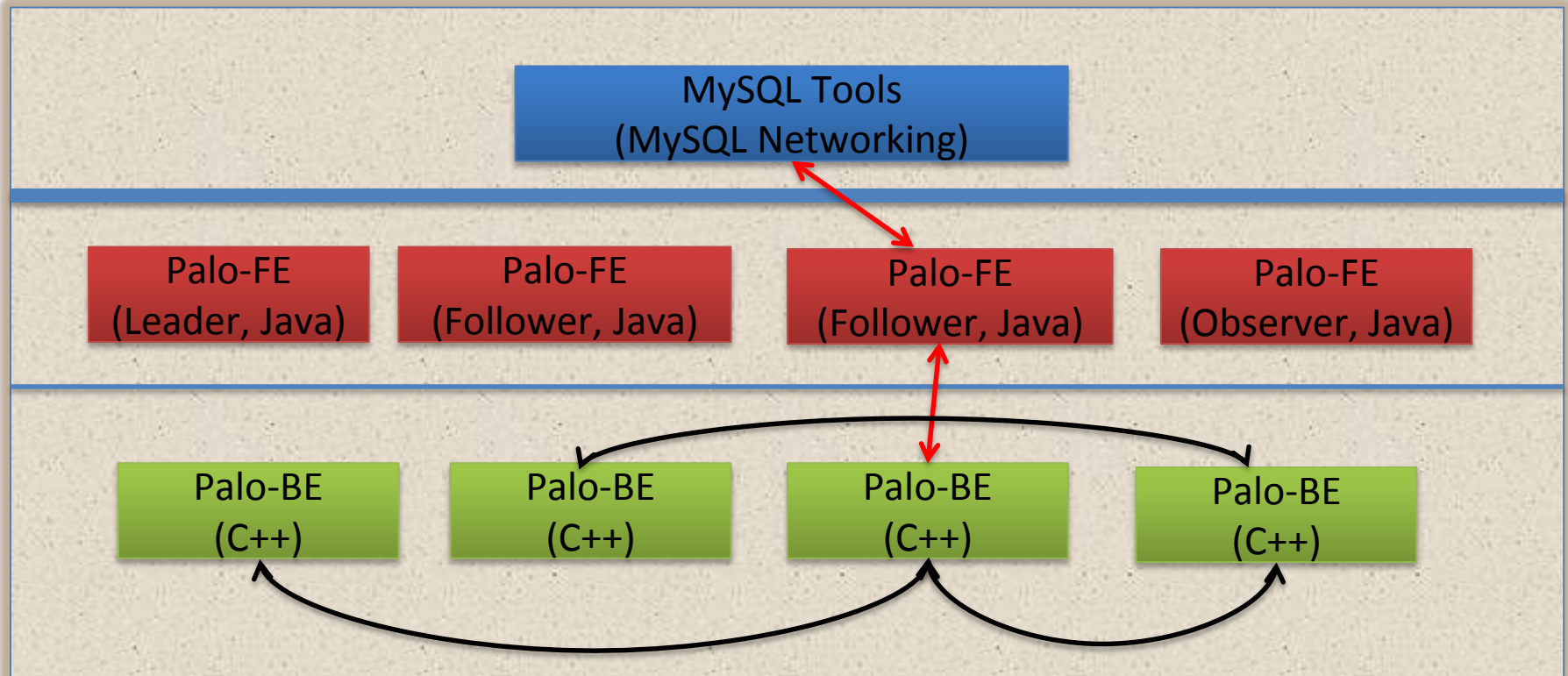
- Leader * 1 + Follower * 2 + Observer * n
- 高可靠；高QPS；可动态扩缩容；支持频繁的元数据更改



集成式系统

- **Integrated system VS layered method**

- 多表原子导入
- 进程内通信 VS 进程间通信
- 易于开发和调试



Palo使用

```
test@my-laptop:~$ mysql -h tc-inf-devop01.tc.baidu.com -P 8276 -u maruyue
Welcome to the MySQL monitor.  Commands end with ; or \g.
Your MySQL connection id is 0
Server version: 4.1.2 (Powered by Palo 2.0 Beta)
```

Copyright (c) 2000, 2013, Oracle and/or its affiliates. All rights reserved.

Oracle is a registered trademark of Oracle Corporation and/or its affiliates. Other names may be trademarks of their respective owners.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

```
mysql> show databases;
```

Database
demo
fc
information_schema
lbs
searchbox
test

6 rows in set (0.01 sec)

```
mysql> use test;
```

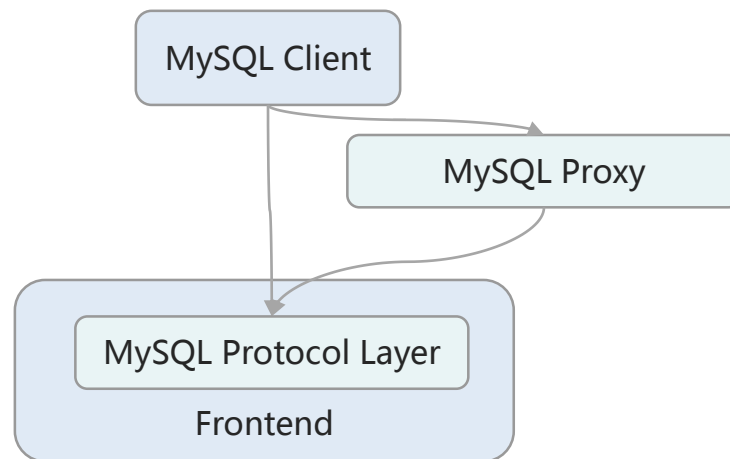
Reading table information for completion of table and column names
You can turn off this feature to get a quicker startup with -A

Database changed

```
mysql> show tables;
```

Tables_in_test
fc_cmatch_fact
tblDIM_pn
tblDIM_querytrade
tblDIM_region
tblDIM_wbws
tblDIM_wos
tblDIM_wpt

7 rows in set (0.01 sec)



- 轻量级客户端
- 与上层应用兼容容易
- 学习曲线平缓，方便用户上手使用
- 利用MySQL相关工具，比如MySQL Proxy

Palo使用

```
1  ./mysql -h PALO_FE_HOST -P PALO_FE_PORT -uYOUR_USERNAME -pYOUR_PASSWORD
2
3  CREATE DATABASE example_db;
4
5  USE example_db;
6
7  CREATE TABLE ps_stats_tbl (
8      siteid  INT,          DEFAULT '10',
9      day     DATETIME,
10     citycode SMALLINT,
11     username VARCHAR(32) DEFAULT '',
12     pv      BIGINT SUM DEFAULT '100'
13 ) DISTRIBUTED BY HASH(siteid) BUCKETS 32;
14
15 LOAD LABEL ps_stats_20150717 (
16     DATA INFILE("hdfs://host:port/ps_stats_data")
17     INTO TABLE ps_stats_tbl
18 );
19
20 SHOW LOAD WHERE LABEL = "ps_stats_20150717";
21
22 SELECT siteid, sum(pv) FROM ps_stats_tbl WHERE day = "2015-07-17" GROUP BY siteid;
23 +-----+-----+
24 | siteid | sum(pv) |
25 +-----+-----+
26 | 23143  | 114996  |
27 | 12345  | 318925  |
28 +-----+-----+
29 2 rows in set (0.02 sec)
```


Palo使用

ORACLE Business Intelligence

1.4 渠道统计与分析

总览

平台

(所有列位)

版式

(所有列位)

日期

介于 2014-11-07

2014-11-13

应用

重置



说明：用户量为各渠道用户量加和，并未做渠道间的去重处理。

说明标题

分析 - 编辑 - 导出

外渠道统计与内渠道总览

查询条件

版式：(所有列位)

平台：(所有列位)

日期	一	搜索 (过滤前)	搜索 (过滤后)	搜索 (过滤前)	搜索 (过滤后)
2014-11-07	RC	9,154	155	23,358	7,305
	其他	658	476	485	401
	其他	9,812	1,038	6,831	1,509
	内	3,429	666	2,533	1,981
	官	6,659	752	8,394	1,124
	客	9,167	1,337	7,940	1,264
	应	9,149	1,170	9,474	1,869
	方	2,117	199	4,806	1,559
	生	1,011	248	7,238	1,372



Palo使用

```
→ test_r R
R version 3.0.1 (2013-05-16) -- "Good Sport"
Copyright (C) 2013 The R Foundation for Statistical Computing
Platform: x86_64-pc-linux-gnu (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

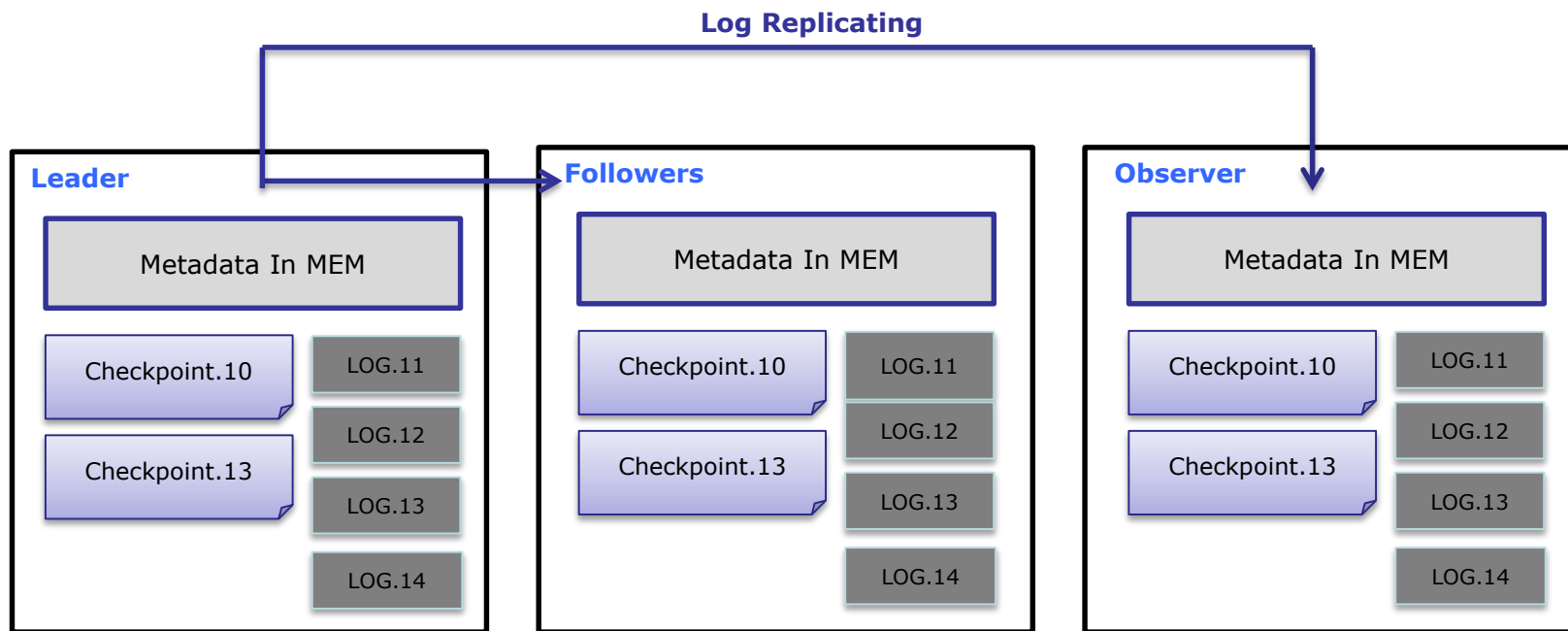
[Previously saved workspace restored]

> library(RMySQL)
Loading required package: DBI
> con <- dbConnect(MySQL(), user="root", password="123456", dbname="demo", host="tc-inf-devop01.tc.baidu.com", port=8276)
> dbListTables(con)
[1] "cumulative_detail_test" "fc_watch_fact"      "tblidm_pn"
[4] "tblidm_querytrade"     "tblidm_region"     "tblidm_mwms"
[7] "tblidm_wos"            "tblidm_wpt"        "ud_test"
> rs <- dbSendQuery(con, "select * from tblidm_region")
> d1 <- fetch(rs, n = 10)
> d1
  pid cid province  city
1  1  0  北京 北京其他
2  2  0  上海 上海其他
3  3  0  天津 天津其他
4  4  0  广东 广东其他
5  5  0  福建 福建其他
6  8  0  海南 海南其他
7  9  0  安徽 安徽其他
8  10 0  贵州 贵州其他
9  11 0  甘肃 甘肃其他
10 12 0  广西 广西其他
> |
```



元数据高可用

- Memory + Checkpoint + Journal
- 采用 Berkeley DB Java Edition , 类Paxos协议实现 ,

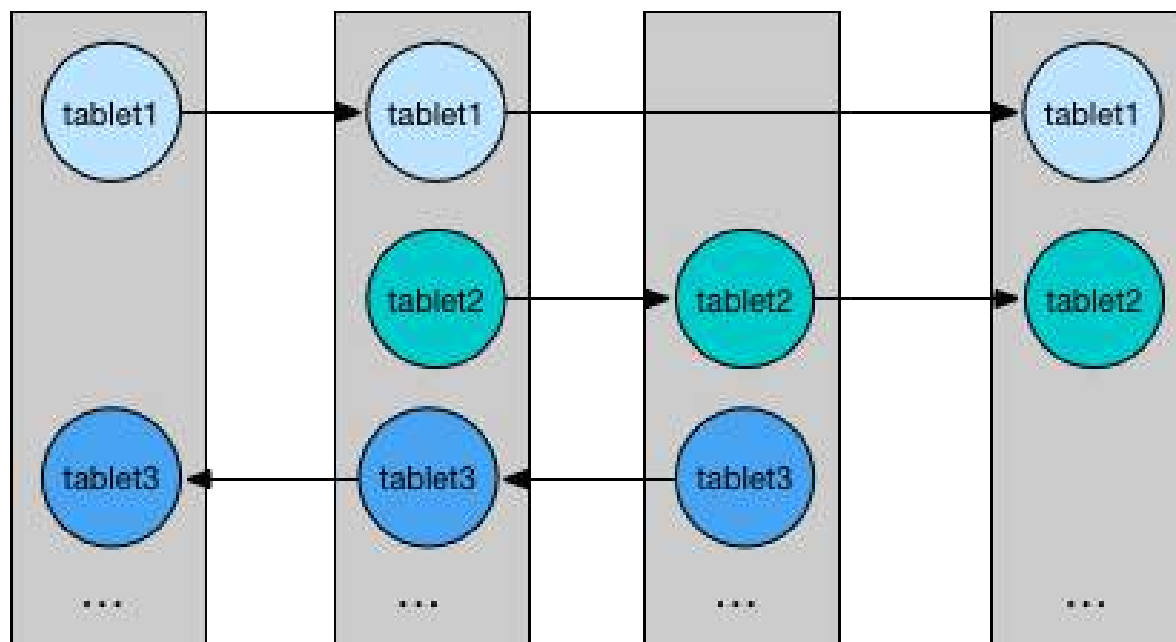


元数据一致性保证

- 会话内单调一致性
 - Master
 - Follower和Observer
- 多个会话间（尤其在失败重试时）
 - 提供sync命令来手动同步

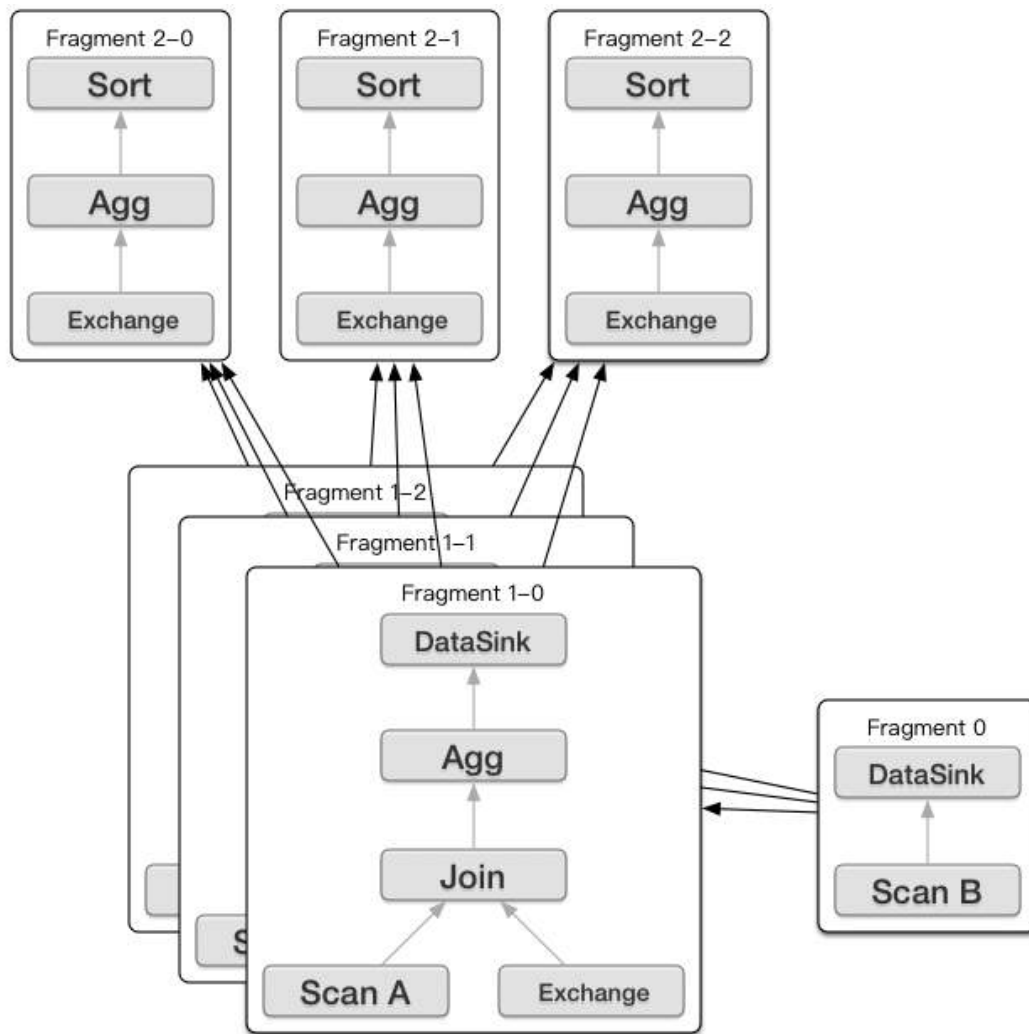
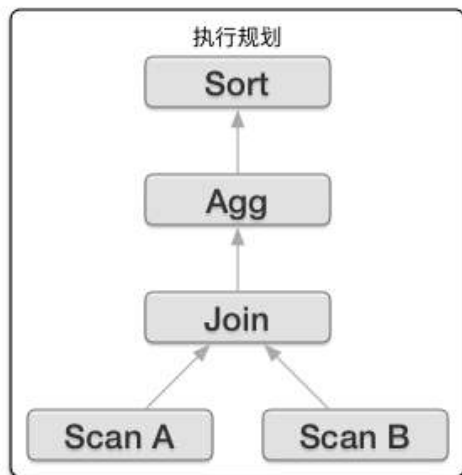
数据高可靠

- 默认三副本
- 自动均衡
- 自动补充



MPP执行

```
SELECT k1, SUM(v1)
FROM A, B
WHERE A.k2 = B.k2
GROUP BY k1
ORDER BY SUM(v1)
```

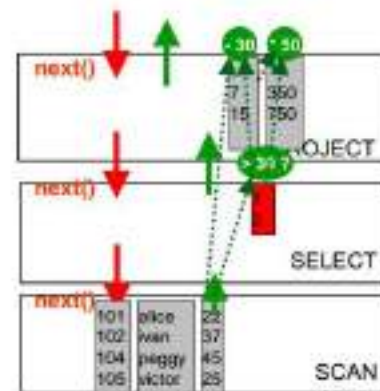
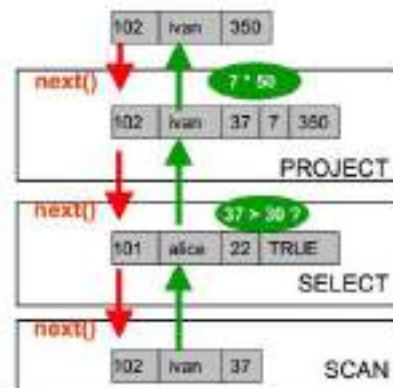


向量化执行&LLVM

• 向量化

- 行式执行引擎问题
 - 每行一次函数调用，打断CPU流水，不利于分支预测
 - 指令和数据cache miss
 - 编译器不友好，不利于循环展开，SIMD
- 设计思想
 - 单条处理到批量处理
 - 行式处理转化为列式处理
- 效果：star-schema测试整体提升3~4倍

```
SELECT id, name  
      (age-30)*50 AS bonus  
FROM   employee  
WHERE  age > 30
```



• LLVM

- 运行时代码生成
- 大型ad-hoc查询可提升5倍以上

存储格式-列存

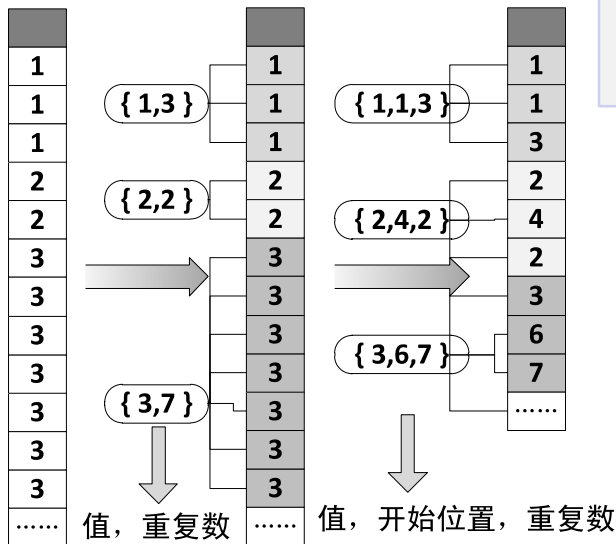
SSN	Name	Age	Addr	City	St
101258787	SMITH	68	986 FIRST ST	JUNO	AL
892375862	CHIN	37	15137 MAIN ST	FOONUN	CA
318370701	HANDU	12	42 JUNE ST	CHICAGO	IL

- ✓ 数据按列存储，每一列单独存放
- ✓ 只访问查询涉及的列，大量降低I/O
- ✓ 数据类型一致，方便压缩
- ✓ 数据包建索引，数据即索引

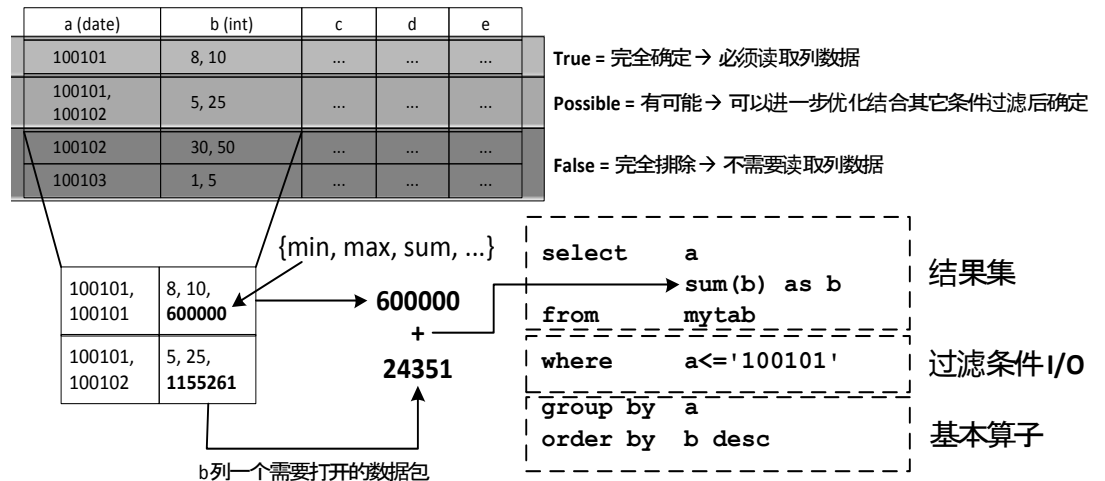
181258787 | 892375862 | 318370701 | 408240180 | 270568310 | 251340875 | 317048351 | 770006228 | 277332171 | 405124589 | 735085947 | 30738201

Block 1

- ✓ Palo存储引擎利用原始过滤条件以及min、max和sum智能索引技术将数据集查询范围尽可能地缩小，可以大大减少I/O，提升查询性能



数值类型的行程编码 (RLE) 压缩



数据模型-Schema

- **Google Mesa模型**

- 维度 (Key列) , 指标 (Value列)

- **Key列有序存储**

- 查询快速定位

- **全Key全局唯一**

- 相同Key的行, 其Value列自动合并(SUM,MIN,MAX,REPLACE)

Time	Id	Country	Clicks	Cost
2016/12/31	1	US	10	32
2017/01/01	2	UK	40	20
2017/01/01	2	US	150	80

数据模型-预聚合

	Time	Id	Country	Clicks	Cost
Base	2016/12/31	1	US	10	32
	2017/01/01	2	UK	40	20
	2017/01/01	2	US	150	80

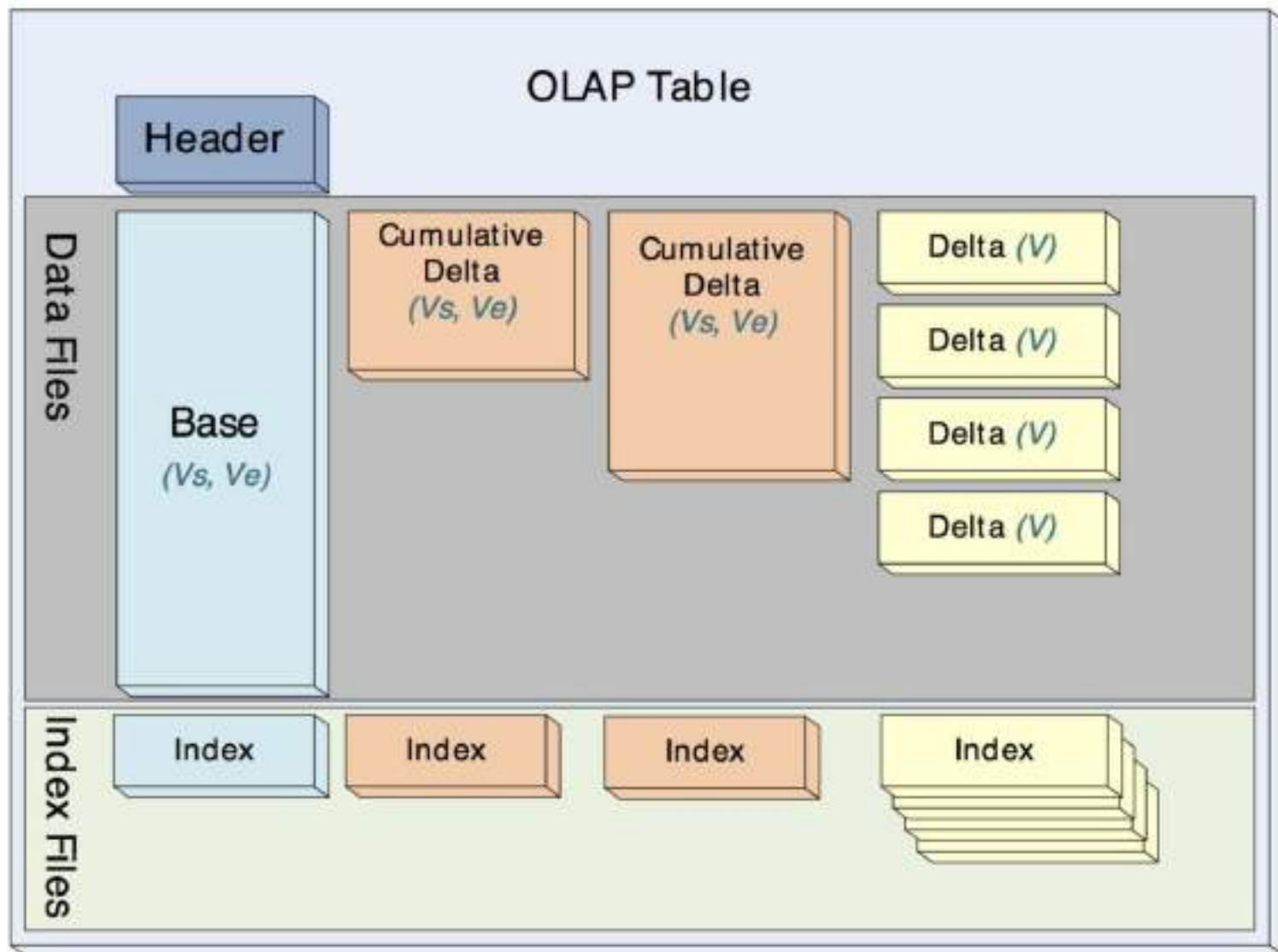
+

	Time	Id	Country	Clicks	Cost
Delta	2017/01/01	1	US	5	3
	2017/01/01	2	UK	60	30
	2017/01/01	2	US	50	20

	Time	Id	Country	Clicks	Cost
New Base	2016/12/31	1	US	10	32
	2017/01/01	1	US	+5	+3
	2017/01/01	2	UK	40+60	20+30
	2017/01/01	2	US	150+50	80+20



数据模型-多版本



数据模型-更多存储模型

- **聚合模型的缺点**

- 不易理解，某些Olap场景没有聚合需求
- 读放大
 - value列的过滤条件无法下推
 - count一个列会造成读取所有列
- 较多key列时，排序本身可能成为瓶颈

- **更丰富的存储模型选择**

- DUPLICATED KEY
- UNIQUE KEY
- AGGREGATE KEY

物化视图(rollup)

- 以空间换时间

- Base表中列的子集
- key列重新排序

- 查询时自动选择

- 导入原子生效

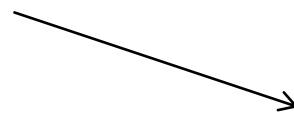
时间	Id	省份	pv
2017.01.01	1	北京	10
2017.01.01	2	天津	30
2017.01.02	1	北京	20
2017.01.02	2	北京	40

Base表



按指定列排序

Id	时间	省份	pv
1	2017.01.01	北京	10
1	2017.01.02	北京	20
2	2017.01.01	天津	30
2	2017.01.02	北京	40



聚合表

Id	pv
1	30
2	70

两层分区 & 分级存储

- **两层分区**

- 方便新旧数据分离，使用不同的存储介质（新数据SSD，历史数据SATA）
- 减少了大量历史数据不必要的重复合并，节省了大量的IO和CPU开销
- 简化了表的扩容，shard调整

- **分级存储**

- 用户可以指定数据放到SSD上或者SATA盘上，也支持根据TTL将冷数据从SSD迁移到SATA上，高效利用SSD提高查询性能

```
1 CREATE TABLE example_tbl (  
2     k1 DATE,  
3     k2 INT,  
4     v1 VARCHAR(2048) REPLACE,  
5 ) PARTITION BY RANGE (k1) (  
6     PARTITION p1 VALUES LESS THAN ("2014-01-01")  
7         properties ("storage_media"="ssd", "storage_cooldown"="2015-06-01 10:00:00"),  
8     PARTITION p2 VALUES LESS THAN ("2014-06-01")  
9         properties ("storage_media"="ssd"),  
10    PARTITION p3 VALUES LESS THAN ("2014-12-01")  
11        properties ("storage_media"="hdd"),  
12 ) DISTRIBUTED BY HASH(k2) BUCKETS 32;  
13
```

Online Schema Change

- 变更期间不停服，用户业务上层不需感知
- 加列、减列，修改列类型等
- 详见 help alter table

数据导入

- 按批导入
- 异步
 - Show load查看状态
- Label机制
 - 防止数据被重复导入

数据导入

- **Broker**

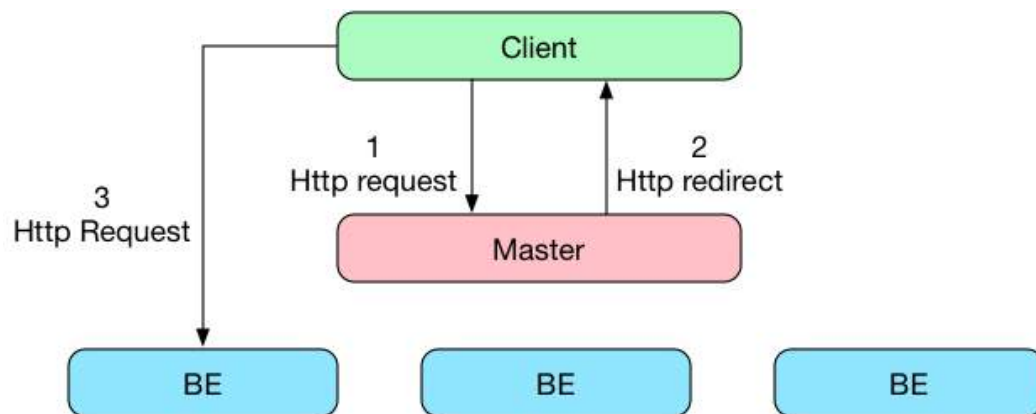
- 支持直接从HDFS、Baidu BOS上进行数据读取。
- 每个物理机部署broker进程，提供并发读取能力。
- 可以访问外部表。

```
LOAD LABEL example_db.label1(  
  DATA INFILE("hdfs://dpp.cluster.com:54310/user/palo/data/input/file")  
  INTO TABLE `my_table`  
)  
WITH BROKER hdfs ("username"="hdfs_user", "password"="hdfs_password")  
PROPERTIES(  
  "timeout"="3600",  
  "max_filter_ratio"="0.1"  
);
```


数据导入

- **MINI BATCH**

- 使用http即可导入，减少客户端对其它组件的依赖
- 实现了多表导入的事务提交



```
-- BATCH DATA LOADING --  
LOAD LABEL ps_stats_20150717 (  
  DATA INFILE("hdfs://host:port/input/ps_stats_data")  
  INTO TABLE ps_stats_tbl  
);
```

```
-- Mini-BATCH DATA LOADING --  
curl -u username,password -T ./input/ps_stats_data http://fe.host:port/api/db1/ps_stats_tbl/_load?label=ps_stats_20150717
```

资源隔离

• 问题

- 多用户影响
- 单用户多任务影响

• 解决

- 线程级cgroup
- 两级资源组织

```
mysql> show resource;
```

User	Resource type	Value
root	CPU_SHARE	1000
zw	CPU_SHARE	1000

```
mysql> show quota;
```

User	Group	Quota
root	high	800
root	low	100
root	normal	400

3 rows in set (0.01 sec)

```
-- cgroup.clone_children
-- cgroup.event_control
-- cgroup.procs
-- cpu.cfs_period_us
-- cpu.cfs_quota_us
-- cpu.shares
-- cpu.stat
-- cpuacct.stat
-- cpuacct.usage
-- cpuacct.usage_percpu
-- notify_on_release
-- root
-- tasks
-- test_user
-- yiguoiei
```

```
-- cgroup.clone_children
-- cgroup.event_control
-- cgroup.procs
-- cpu.cfs_period_us
-- cpu.cfs_quota_us
-- cpu.shares
-- cpu.stat
-- cpuacct.stat
-- cpuacct.usage
-- cpuacct.usage_percpu
-- high
-- cgroup.clone_children
-- cgroup.event_control
-- cgroup.procs
-- cpu.cfs_period_us
-- cpu.cfs_quota_us
-- cpu.shares
-- cpu.stat
-- cpuacct.stat
-- cpuacct.usage
-- cpuacct.usage_percpu
-- notify_on_release
-- tasks
-- low
-- cgroup.clone_children
-- cgroup.event_control
-- cgroup.procs
-- cpu.cfs_period_us
-- cpu.cfs_quota_us
-- cpu.shares
-- cpu.stat
-- cpuacct.stat
-- cpuacct.usage
```

部分SQL增强

- **支持窗口函数**

- 计算排名、同环比问题变得简单高效
- 同一部门内，薪水最高的三位？

rank() over (partition by dept order by salary)

- **类Proc机制**

- Show proc " /XX"

- **HELP**

- 在线查看帮助文档
- Web help
- Help xxx

联系方式

- GitHub :

<https://github.com/baidu/palo>

- 百度云 :

<https://cloud.baidu.com/product/palo.html>

- 邮件 :

palo-rd@baidu.com

- Palo开源讨论群 :

加我的微信，备注“加入Palo技术讨论群”

[myh13161636186](https://www.weixin.qq.com/wxa/myh13161636186)





谢谢