

---

# 大数据应用

OLAP & 数据挖掘

---

# Agenda

OLAP概念

OLAP&OLTP

kylin介绍

# 数据仓库

一个面向主题的、集成的、时变的、非易失性的数据集合，对信息处理提供支持。

- 面向主题的
- 集成的
- 时变的
- 非易失性的

# 操作数据库系统 vs 数据仓库

操作数据库系统：OLTP

数据仓库：OLAP

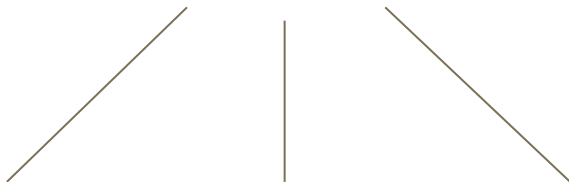
# 数据处理类型

OLAP

OLTP

Online Analytical Processing

Online Transaction Processing



ROLAP

MOLAP

HOLAP

Relational OLAP

Multi dimensional OLAP

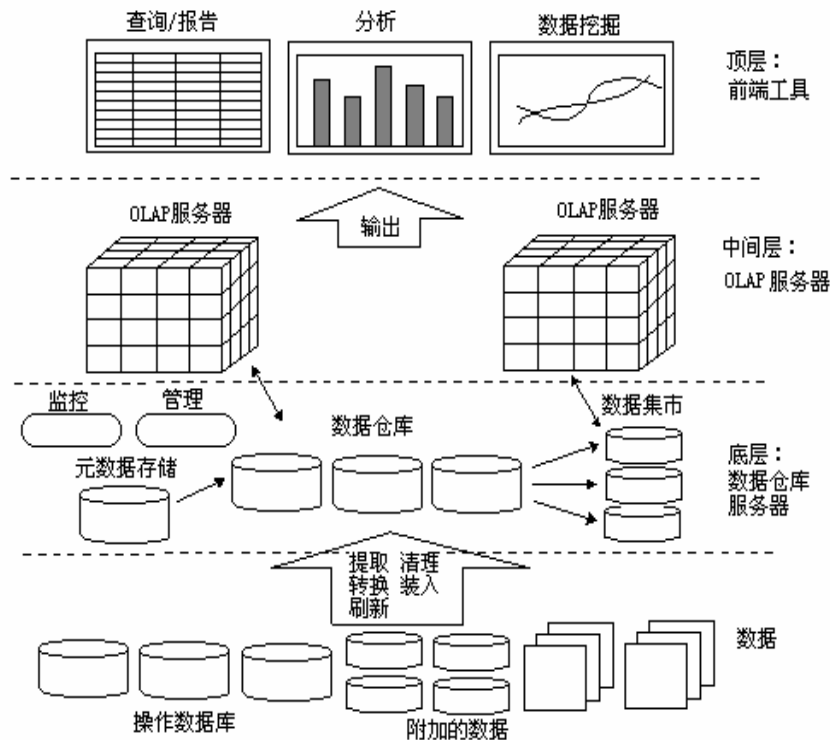
Hybrid dimensional OLAP

# OLTP vs OLAP

特征	OLTP	OLAP
特性	操作处理	信息处理
面向	事务	分析
用户	办事员、DBA、数据库专业人员	知识工人（如经理、主管、分析人员）
功能	日常操作	长期信息需求、决策支持
DB 设计	基于 E-R，面向应用	星形/雪花、面向主题
数据	当前的、确保最新	历史的、跨时间维护
汇总	原始的、高度详细	汇总的、统一的
视图	详细、一般关系	汇总的、多维的
工作单元	短的、简单事务	复杂查询
访问	读/写	大多为读
关注	数据进入	信息输出
操作	主码上索引/散列	大量扫描
访问记录数量	数十	数百万
用户数	数千	数百
DB 规模	GB 到高达 GB	≥TB
优先	高性能、高可用性	高灵活性、终端用户自治
度量	事务吞吐量	查询吞吐量、响应时间

# 三层数据仓库结构

- 数据仓库服务器
- OLAP服务器
- 前端客户层



# 数据仓库建模－数据立方体

维度 & 事实

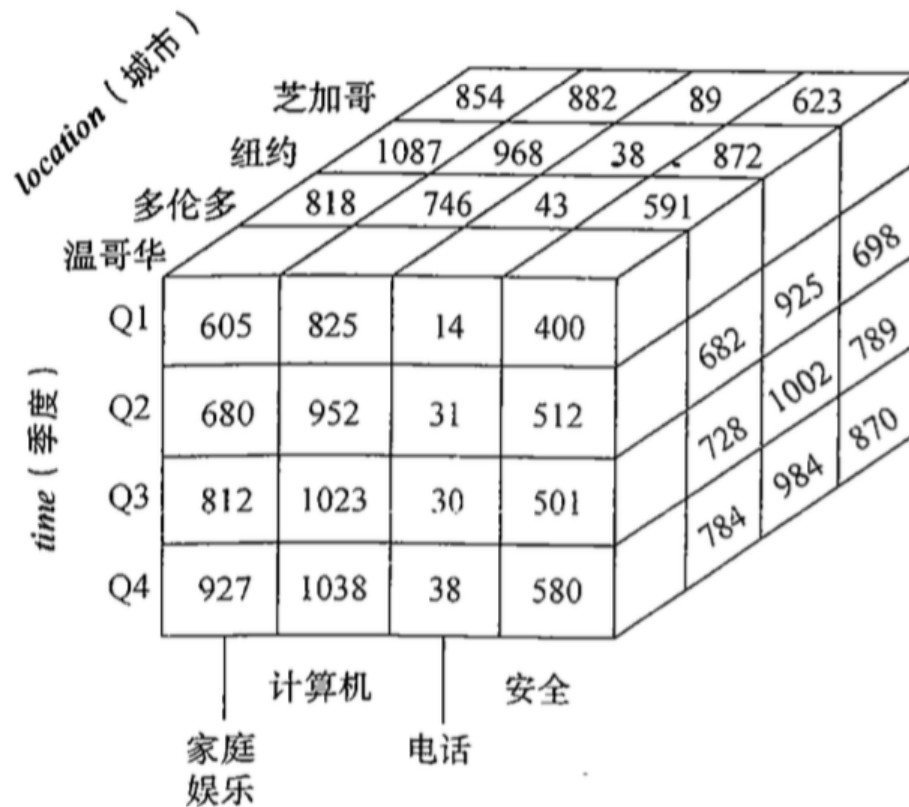
<i>location</i> = “温哥华”				
<i>time</i> (季度)	<i>item</i> (类型)			
	家庭娱乐	计算机	电话	安全
Q1	605	825	14	400
Q2	680	952	31	512
Q3	812	1023	30	501
Q4	927	1038	38	580



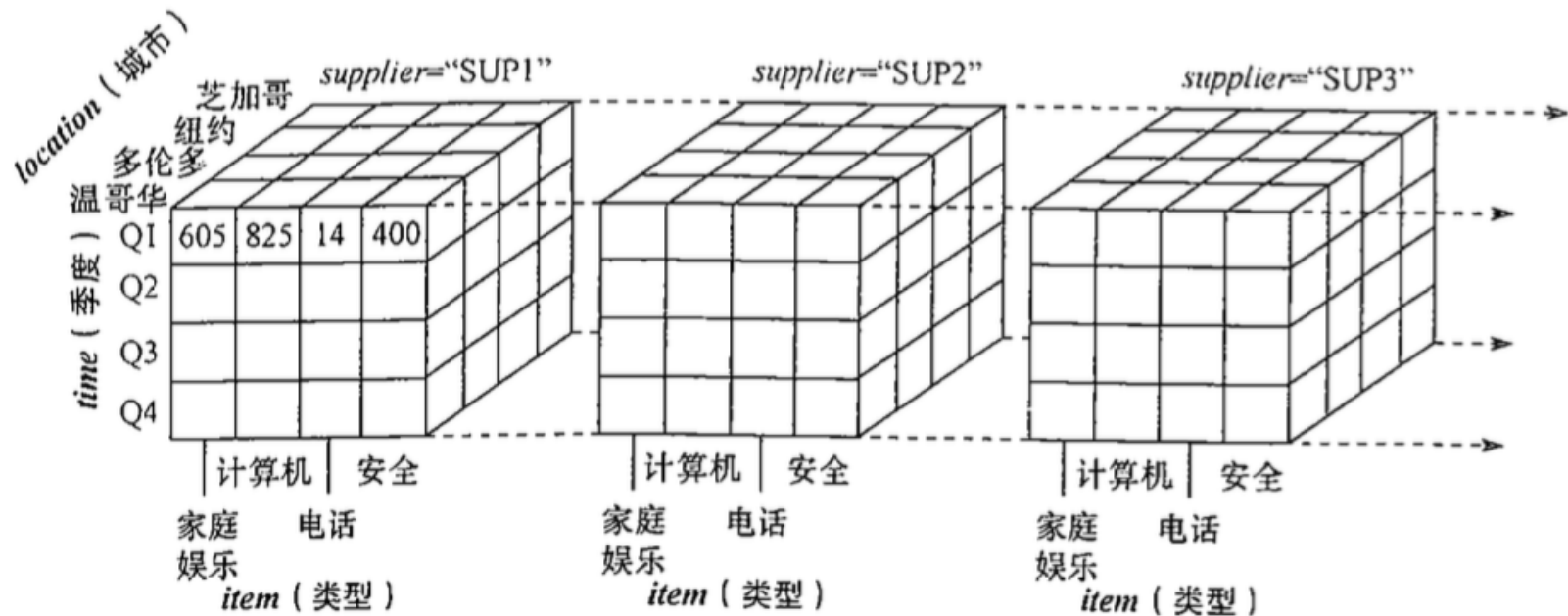
# 3D

<i>time</i>	<i>location</i> = “芝加哥” <i>item</i>				<i>location</i> = “纽约” <i>item</i>				<i>location</i> = “多伦多” <i>item</i>				<i>location</i> = “温哥华” <i>item</i>			
	家庭 娱乐	计算 机	电话	安全	家庭 娱乐	计算 机	电话	安全	家庭 娱乐	计算 机	电话	安全	家庭 娱乐	计算 机	电话	安全
Q1	854	882	89	623	1087	968	38	872	819	746	43	591	605	825	14	400
Q2	943	890	64	698	1130	1024	41	925	894	769	52	682	680	952	31	512
Q3	1032	924	59	789	1034	1048	45	1002	940	795	58	728	812	1023	30	501
Q4	1129	992	63	870	1142	1091	54	984	978	864	59	784	927	1038	38	580

# 3D

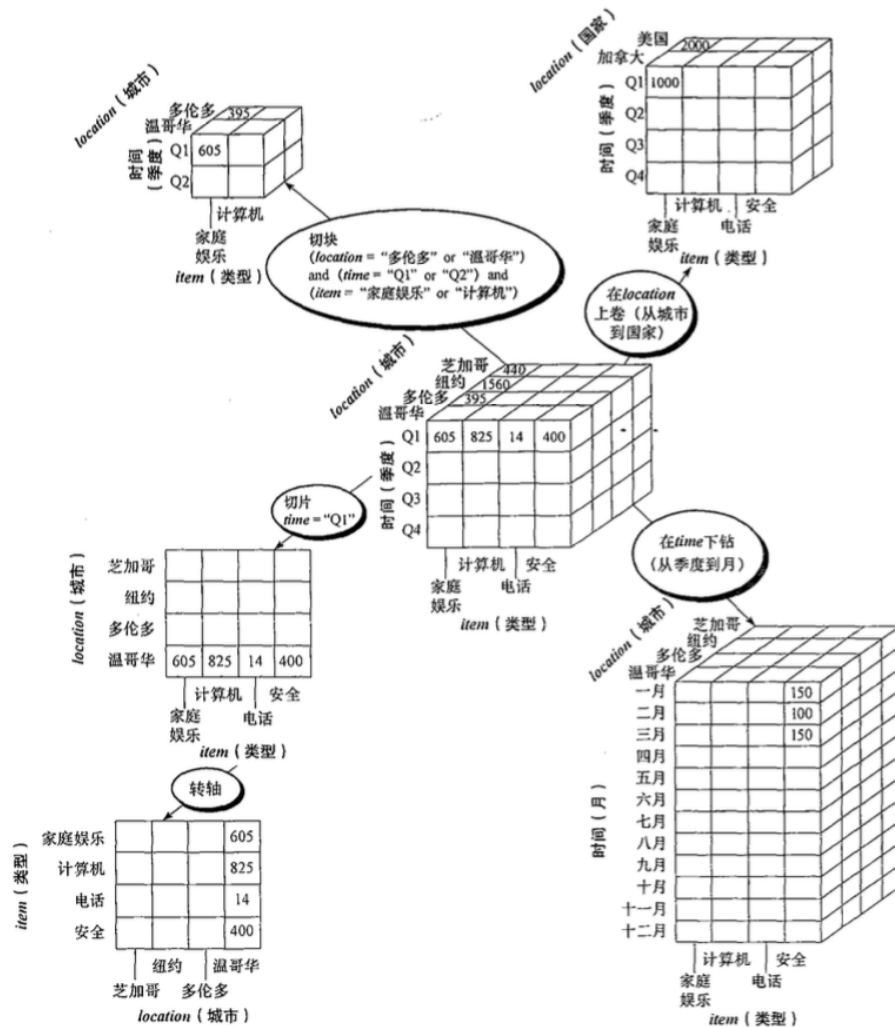


# 4D



# OLAP 操作

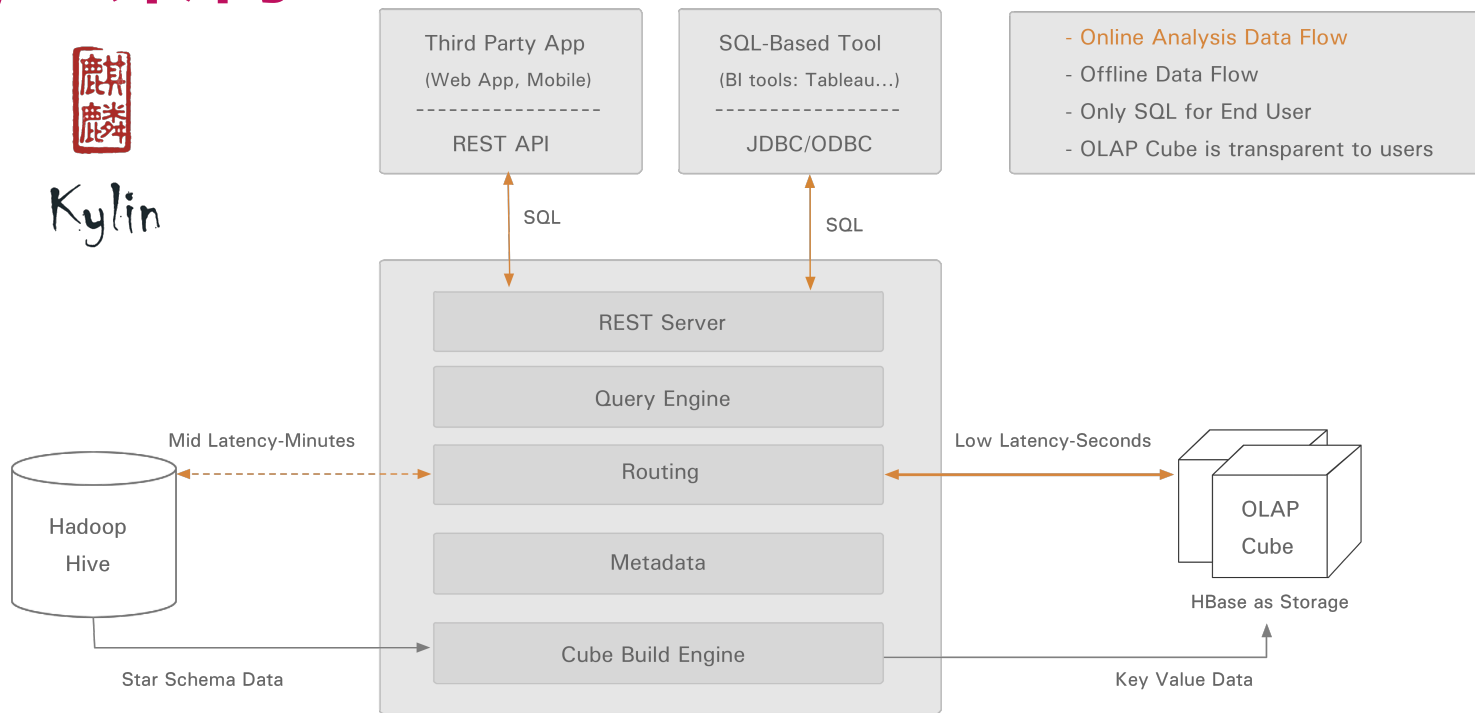
- 上卷
- 下钻
- 切片 / 切块
- 转轴



# Kylin架构



Kylin



# Cube

Dimension 维度

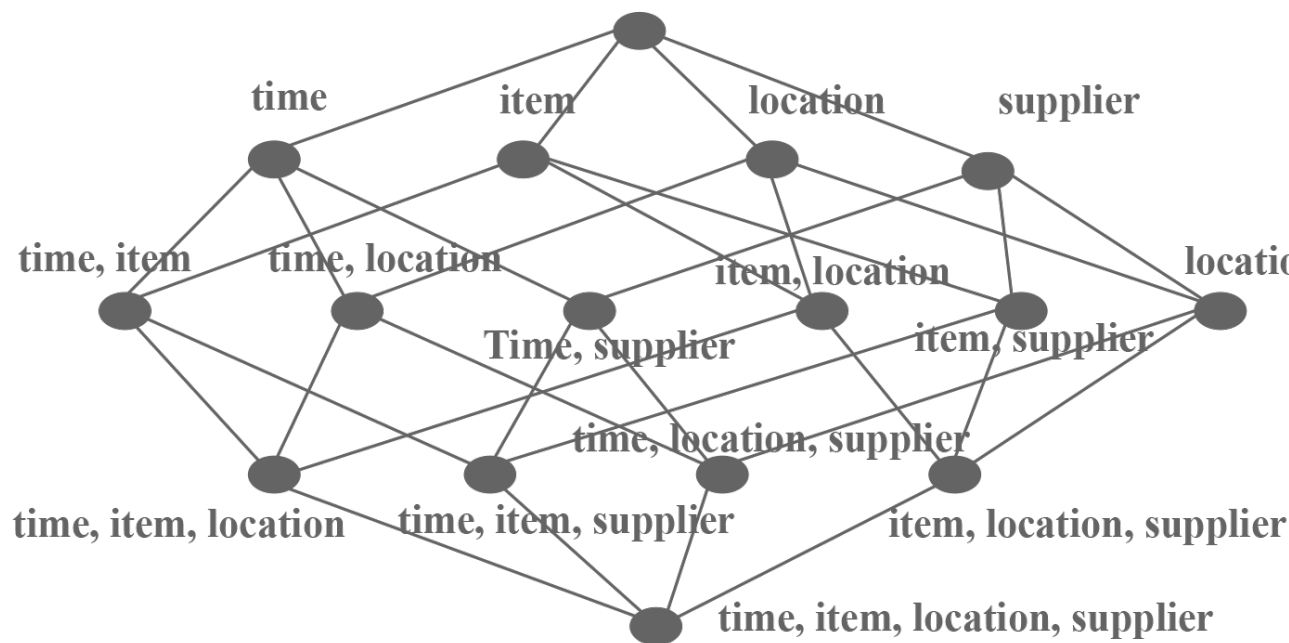
Measure 度量

Member 成员

Level 层次

A	B	C	D	v
A1	B1	C1	D1	2
A1	B2	C1	D1	3
A2	B1	C1	D1	5
A3	B1	C1	D1	6
A3	B2	C1	D1	8

# Cube 构建



0-D(apex) cuboid

1-D cuboids

2-D cuboids

3-D cuboids

4-D(base) cuboid



# 聚合

A	B	C	D	v
A1	B1	C1	D1	2
A1	B2	C1	D1	3
A2	B1	C1	D1	5
A3	B1	C1	D1	6
A3	B2	C1	D1	8

base cuboid

$\langle A1, B1, C1, D1 \rangle, 2$   
 $\langle A1, B2, C1, D1 \rangle, 3$   
 $\langle A2, B1, C1, D1 \rangle, 5$   
 $\langle A3, B1, C1, D1 \rangle, 6$   
 $\langle A3, B2, C1, D1 \rangle, 8$

$A1, B1, C1, 2$   
 $A1, B1, D1, 2$   
 $A1, C1, D1, 2$   
 **$B1, C1, D1, 2$**

$A2, B1, C1, 5$   
 $A2, B1, D1, 5$   
 $A2, C1, D1, 5$   
 **$B1, C1, D1, 5$**

$A3, B1, C1, 6$   
 $A3, B1, D1, 6$   
 $A3, C1, D1, 6$   
 **$B1, C1, D1, 6$**

sum(v)

next cuboid

$\langle A1, B1, C1 \rangle, 2$   
 $\langle A1, B1, D1 \rangle, 2$   
 $\langle A1, C1, D1 \rangle, 2 + 3$   
 **$\langle B1, C1, D1 \rangle, 2 + 5 + 6$**

.....

# kyin dimension

- Mandatory dimension cuts cuboid combinations by half.

Normal Dimensions

A	B	C
A	B	-
-	B	C
A	-	C
A	-	-
-	B	-
-	-	C
-	-	-



A is Mandatory

A	B	C
A	B	-
A	-	C
A	-	-

# hierarchy

- Hierarchy dimension reduces combination from  $2^N$  to  $N+1$ .

Normal Dimensions

A	B	C
A	B	-
-	B	C
A	-	C
A	-	-
-	B	-
-	-	C
-	-	-



A->B->C is Hierarchy

A	B	C
A	B	-
A	-	-
-	-	-

# Derived

- Derived dimension reduces combination from  $2^N$  to 2 at the cost of extra runtime aggregation.

Normal Dimensions

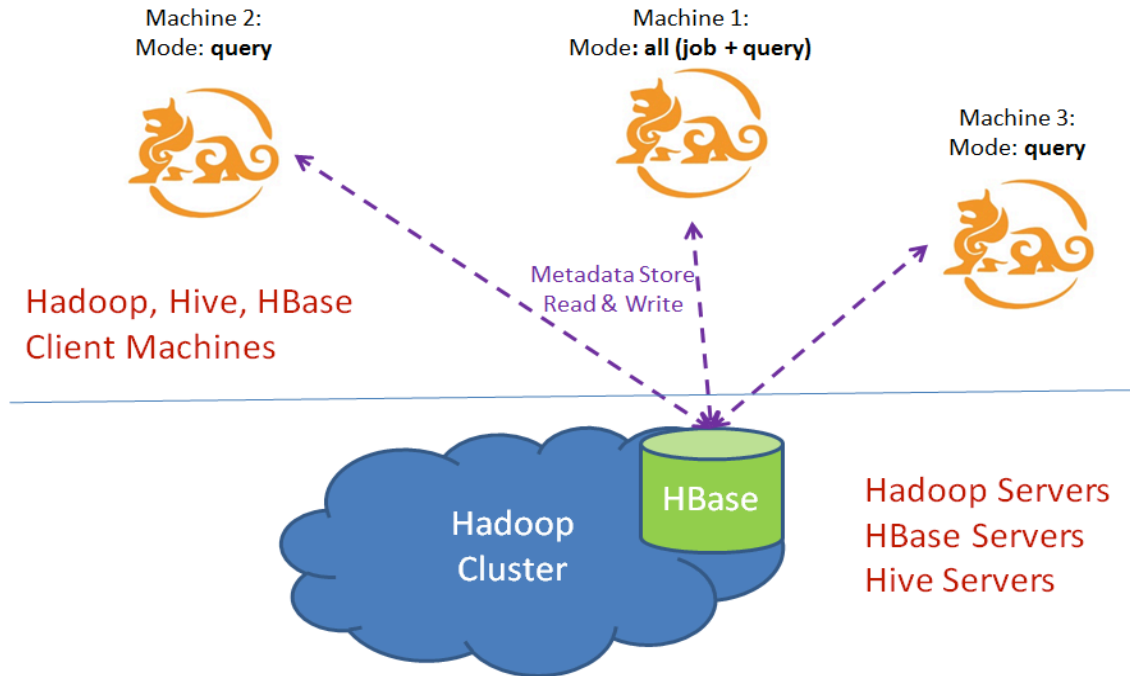
A	B	C
A	B	-
-	B	C
A	-	C
A	-	-
-	B	-
-	-	C
-	-	-



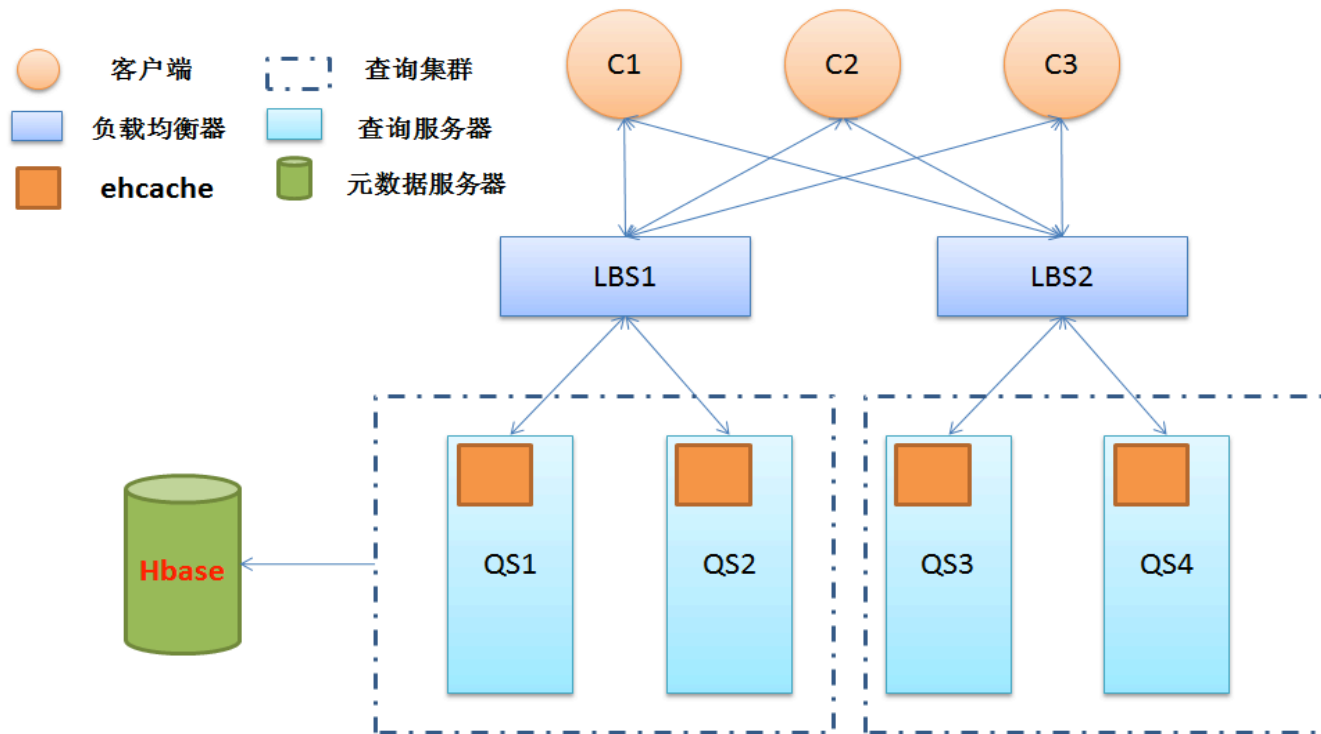
A, B, C are Derived by ID

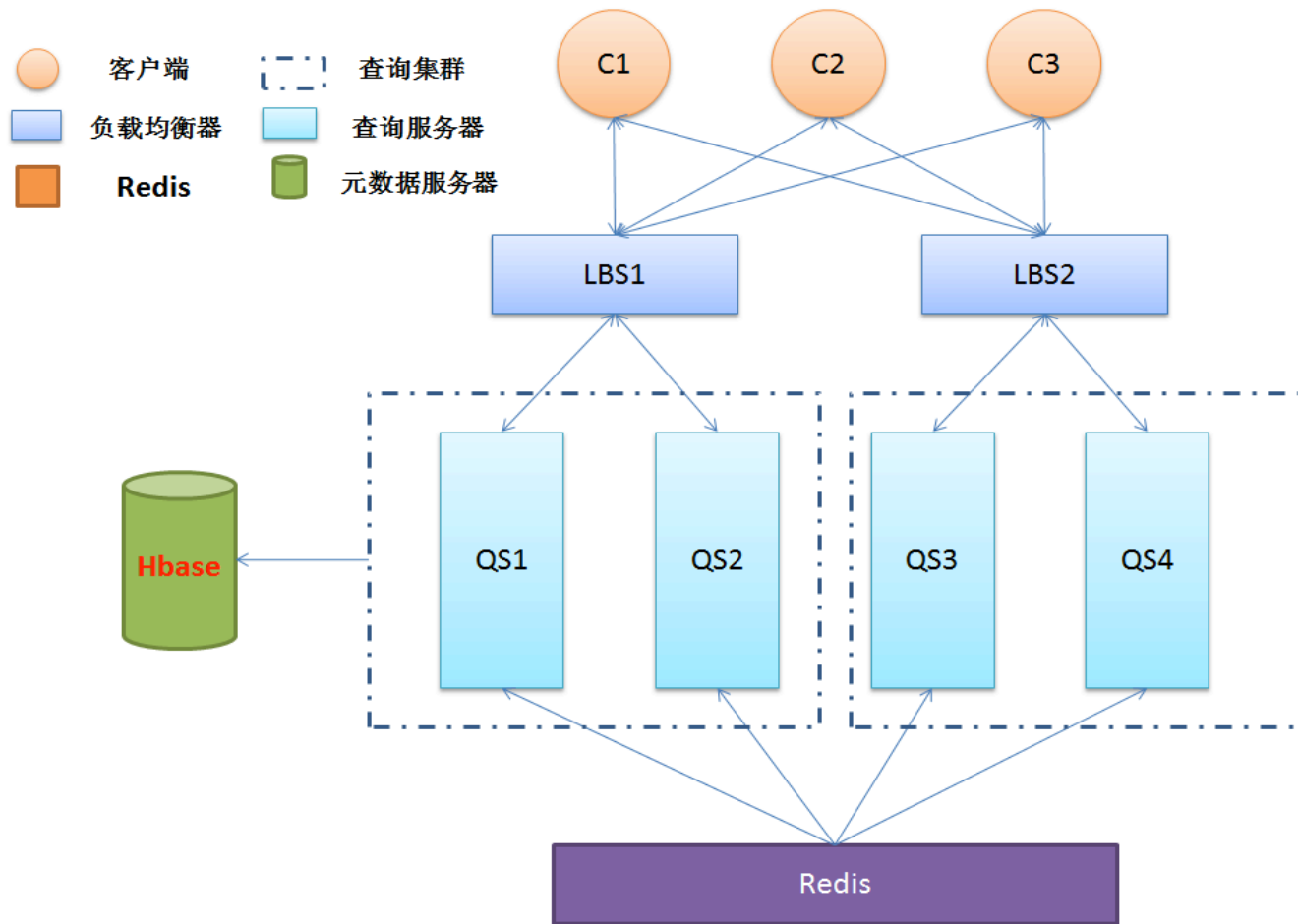
ID
-

# kylin 设计概览



# Kylin 集群架构

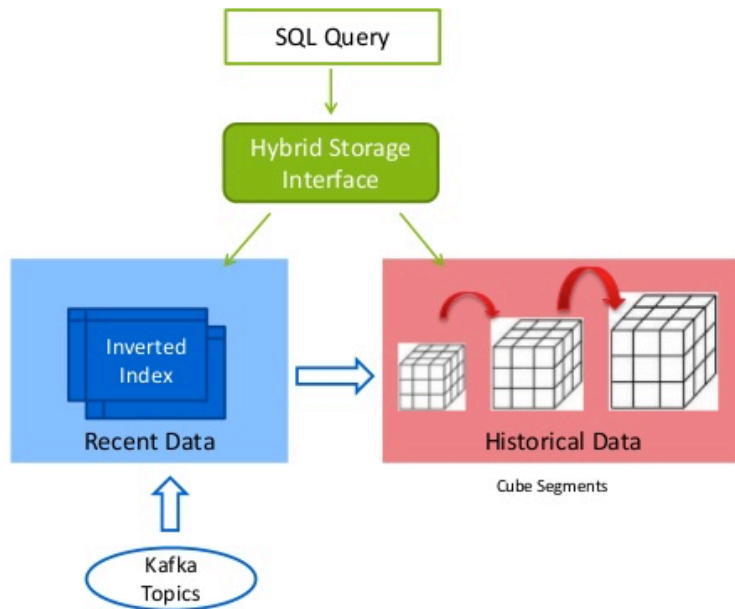






# kylin segment merge

## Optimizing Cube Builds: Streaming



# OLAP projects



MOLAP

kylin      druid



ROLAP

presto, imala,  
hive, sparkSQL

# OLAP 比较

## kylin

预计算，适合业务清楚，分析场景明确。  
要避免维度灾难。

## presto

没有使用MapReduce，大部分场景下比hive快一个数量级，  
所有的处理都在内存中完成

## druid

实时处理时序数据，索引按照时间分片，查询也是按照时间线去路由索引。  
保证数据实时写入，但查询目前只支持部分SQL，适合用于工业大数据