

# CEDAR: A Distributed Database for Scalable OLTP



Weining Qian  
wnqian@dase.ecnu.edu.cn

華東師範大學



# The Internet economy

2

- Content/traffic => money
- O2O :  
Online => Offline, or  
Offline => Online
- **O2O of mission-critical apps: 互联网+ (Internet+)**
- OLTP is inevitable

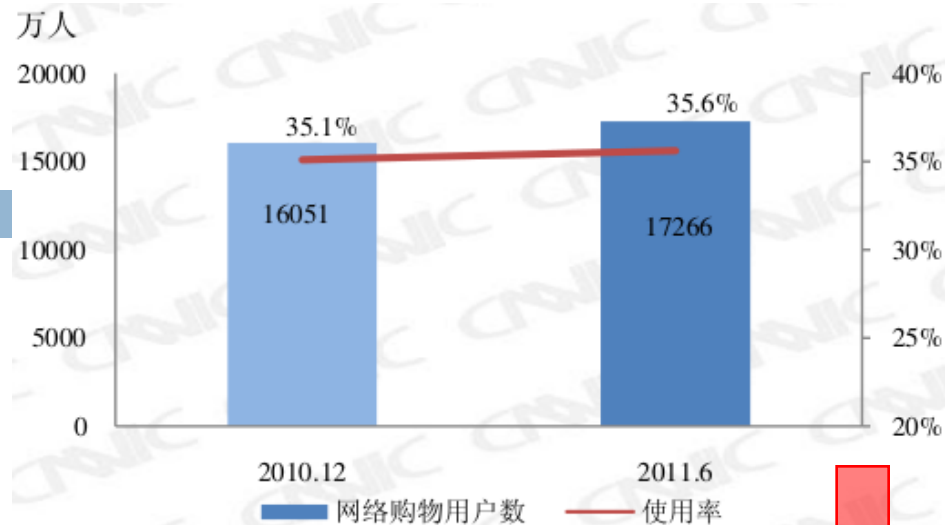
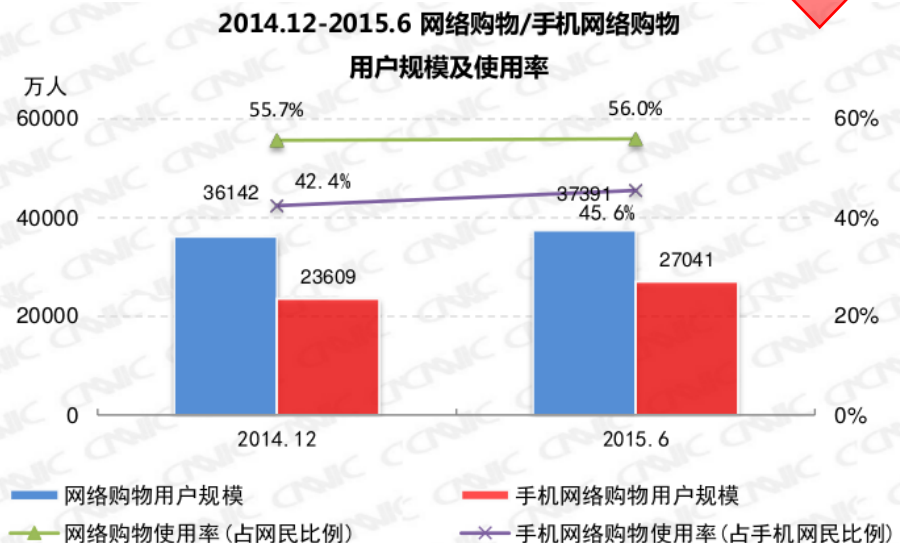


图 20 2010.12-2011.6 网络购物用户数及使用率



来源：CNIC 中国互联网络发展状况统计调查

2015.6

# Phenomenal applications

3

- Phenomenal - very remarkable; extraordinary.



- 180,000 tps in year 2016

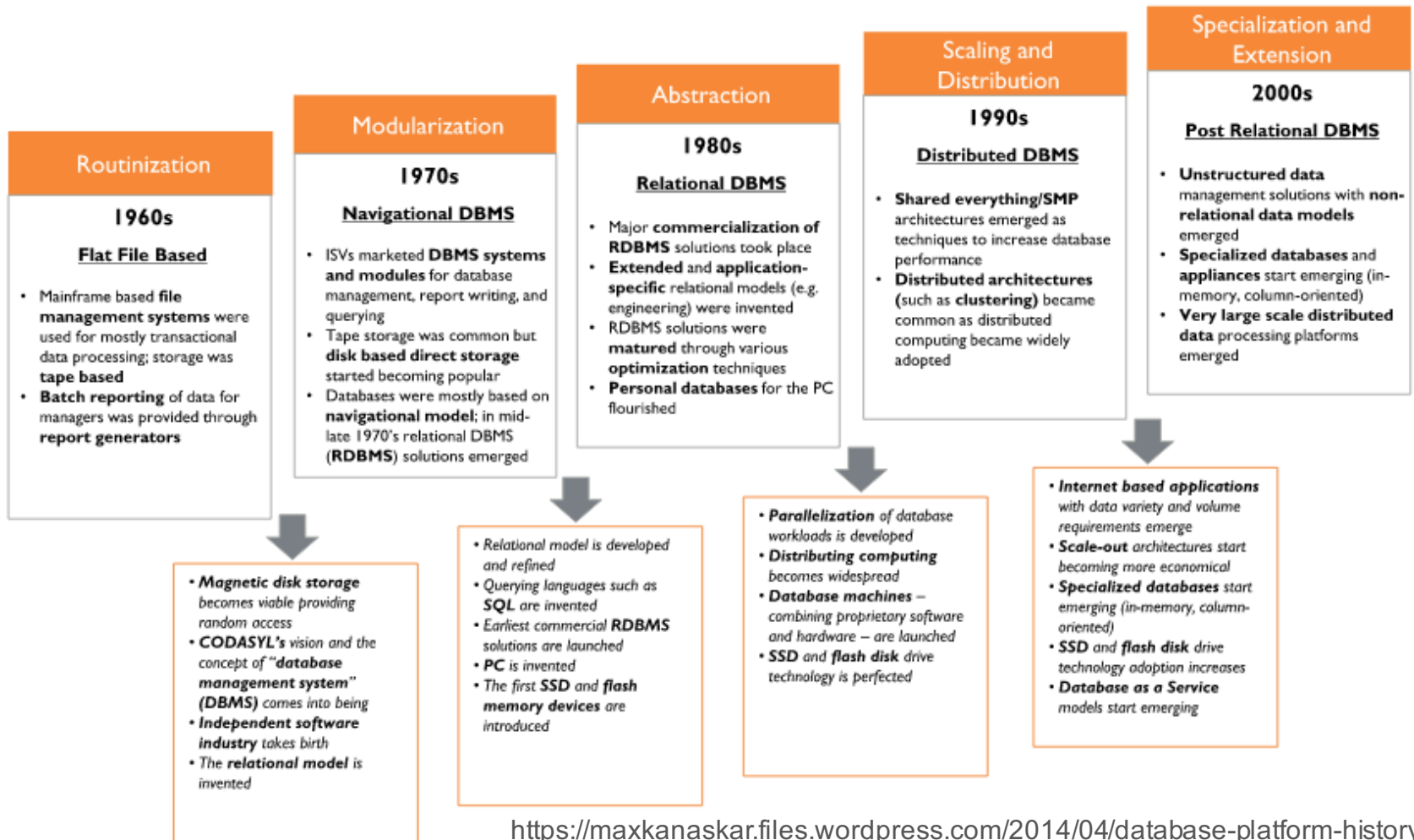
# Phenomenal **is common**

4

- 12306 during Spring Festival
- Black Friday promotion/second kill, ...
- The pressure is on **backend (transaction | payment) systems**
- Be inevitable and day-to-day more common
  
- **All pressure will finally go to mission critical systems**
- **Essentially a High-throughput, scalable transaction processing problem**

# A brief history of DBMS

5



# One size fits all => One size fits none!

6



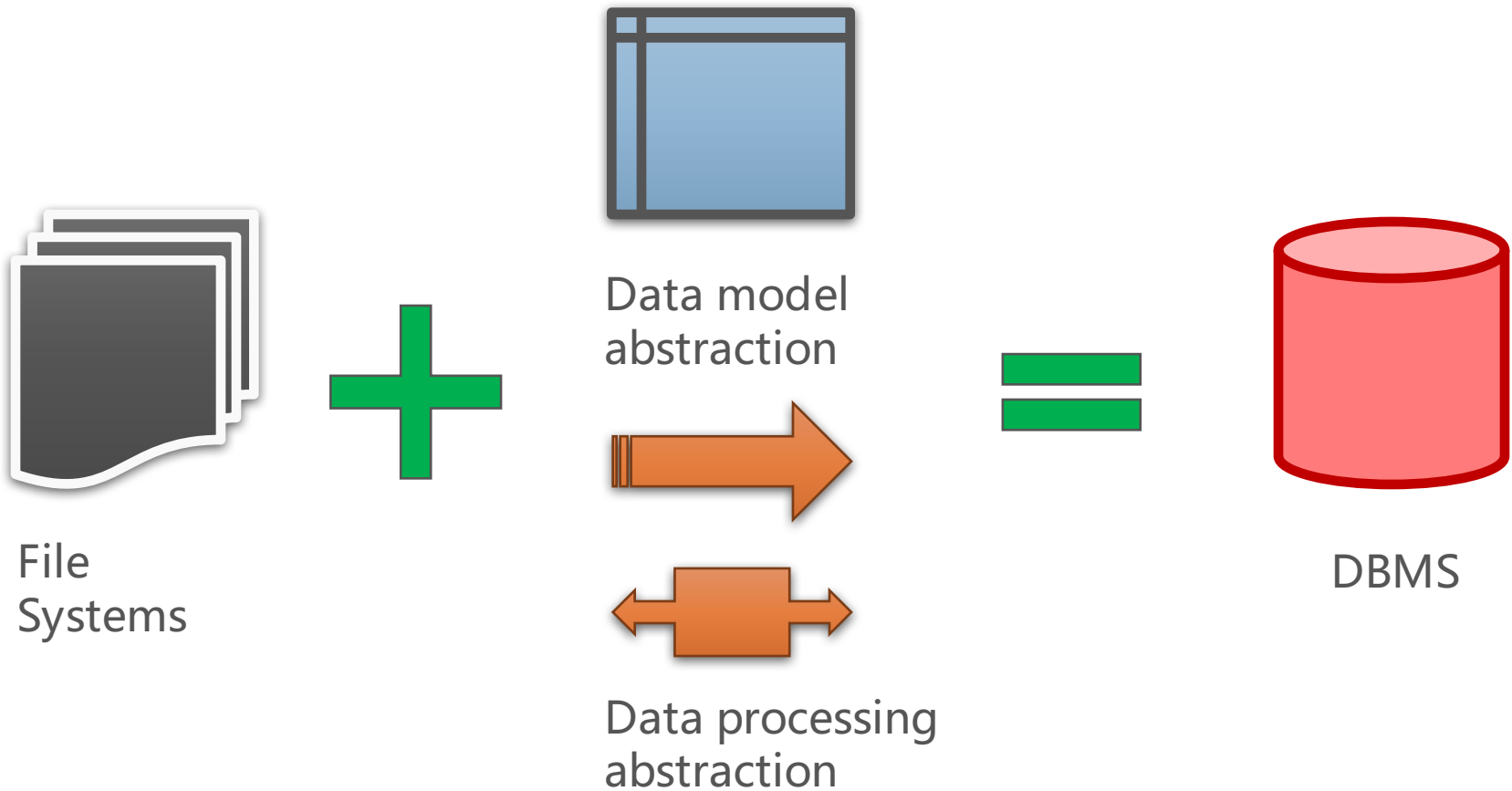
## 2015 - One Size Fits None!!

- Data Warehouse market
- OLTP market
- NoSQL market
- Complex Analytics
- Streaming market
- Graph analytics market



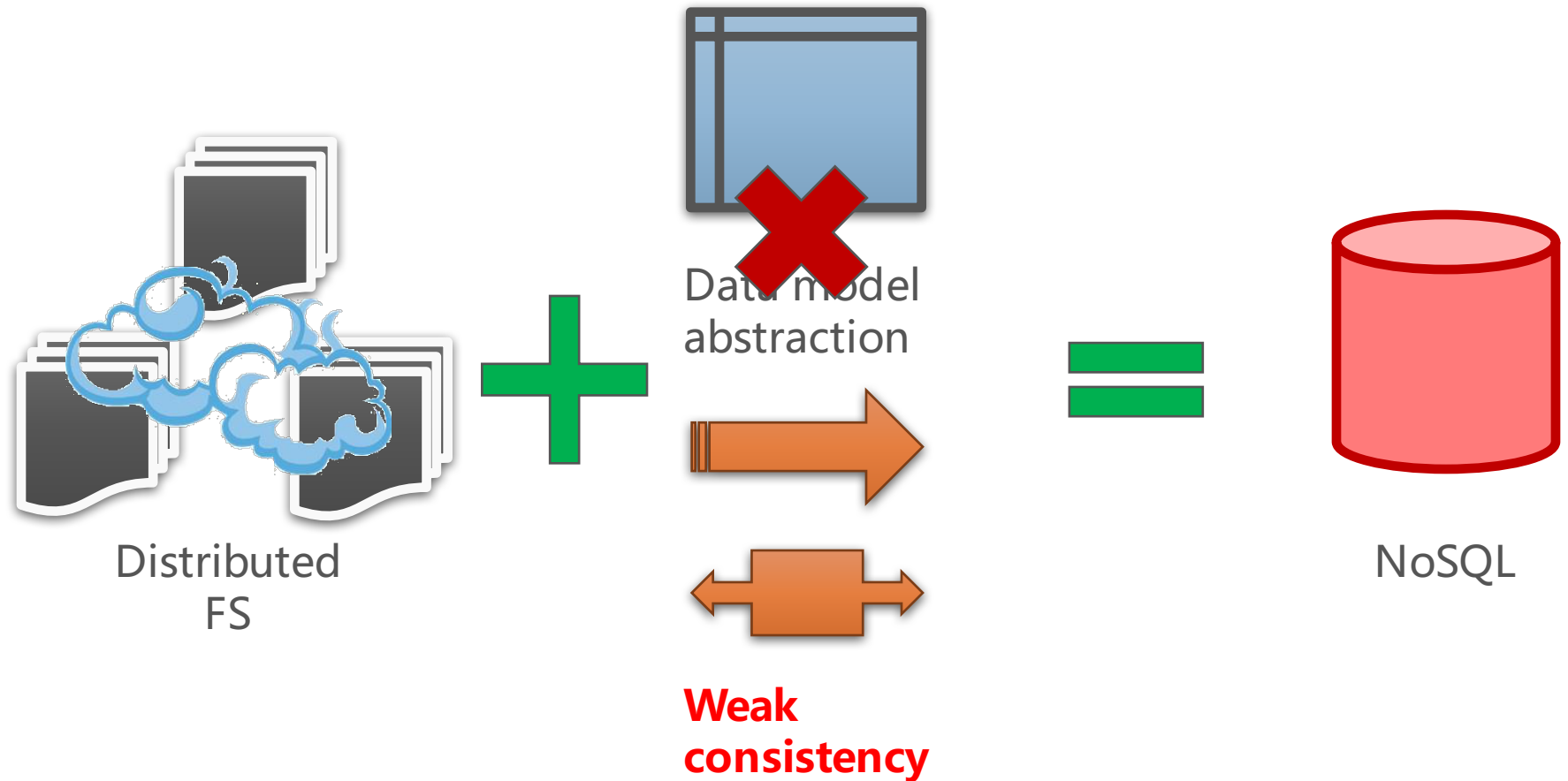
# DBMS

7



# NoSQL

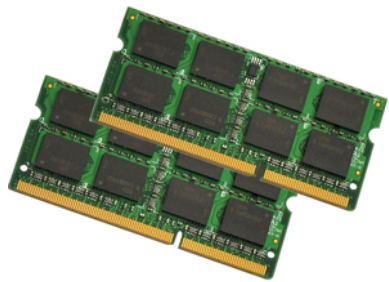
8



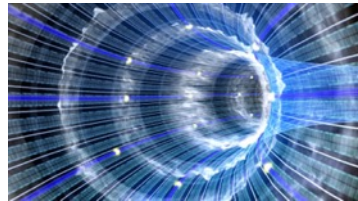


# NewSQL

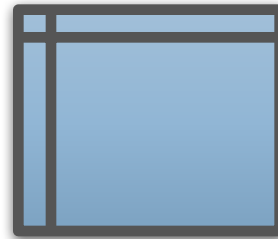
9



In-memory  
computing



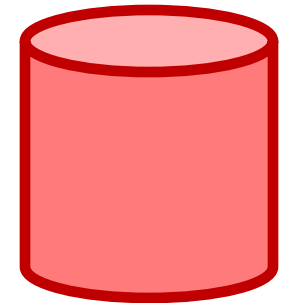
Fast networking



Relational  
model



Transaction  
processing



NewSQL

# OldSQL vs. NoSQL vs. NewSQL

10

	OldSQL	NoSQL	NewSQL
Data model	Relational	---	Relational
Interface	SQL	Variance	SQL
Consistency/Concurrency control	Strong	Weak	Strong
Fault tolerance	Strong	Fine	Strong
Performance	Poor	Good	Very good
Scalability	Poor	Good	Fine

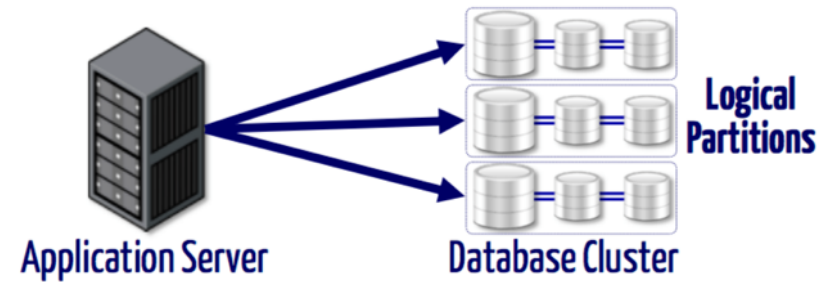


# How to scale a database system?

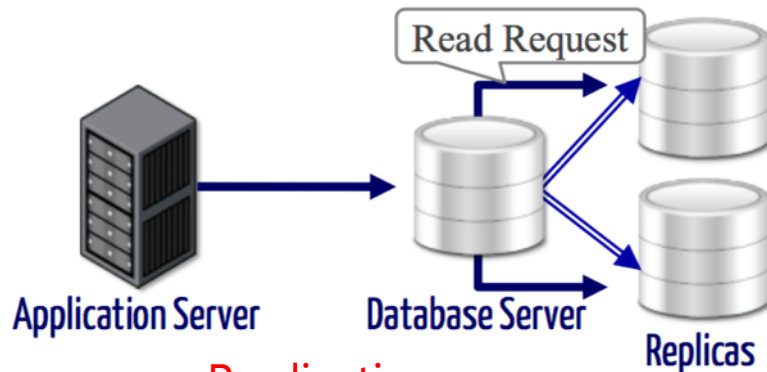
11



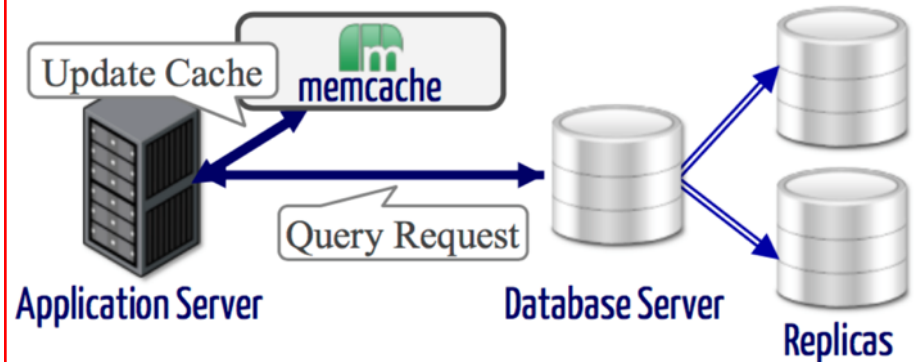
Scaling-up



Sharding/partitioning



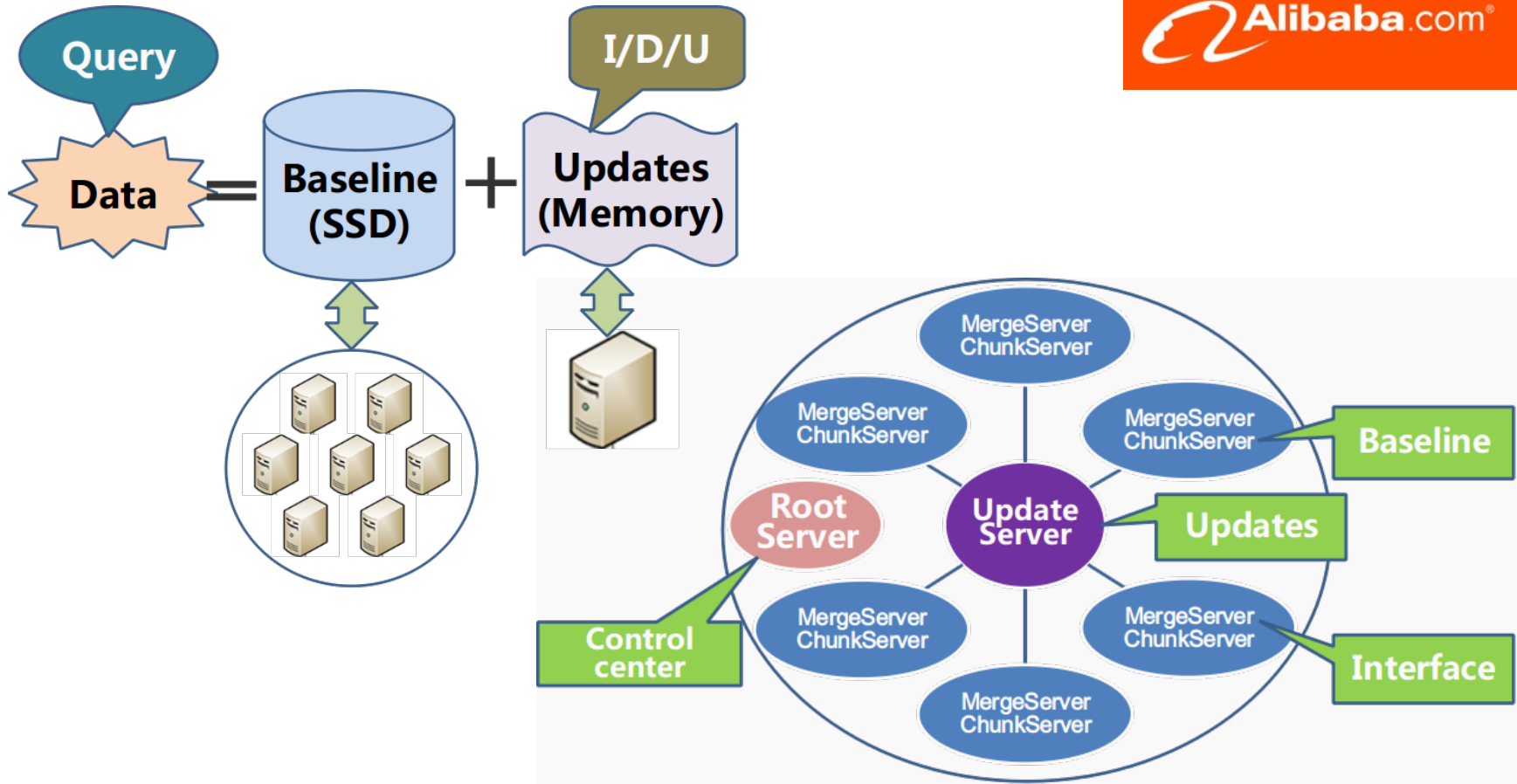
Replication



Caching

# The open-source OceanBase (0.4)

12



# OceanBase 0.4 is not enough

13

- ❑ Simple transactions only
- ❑ Weak availability support
- ❑ Single-point transaction
- ❑ More optimization needed for query/storage
- ❑ Interface adaptability

Not enough for mission critical apps in banks with strong-consistency, high-availability, and high-throughput complex transaction processing requirements

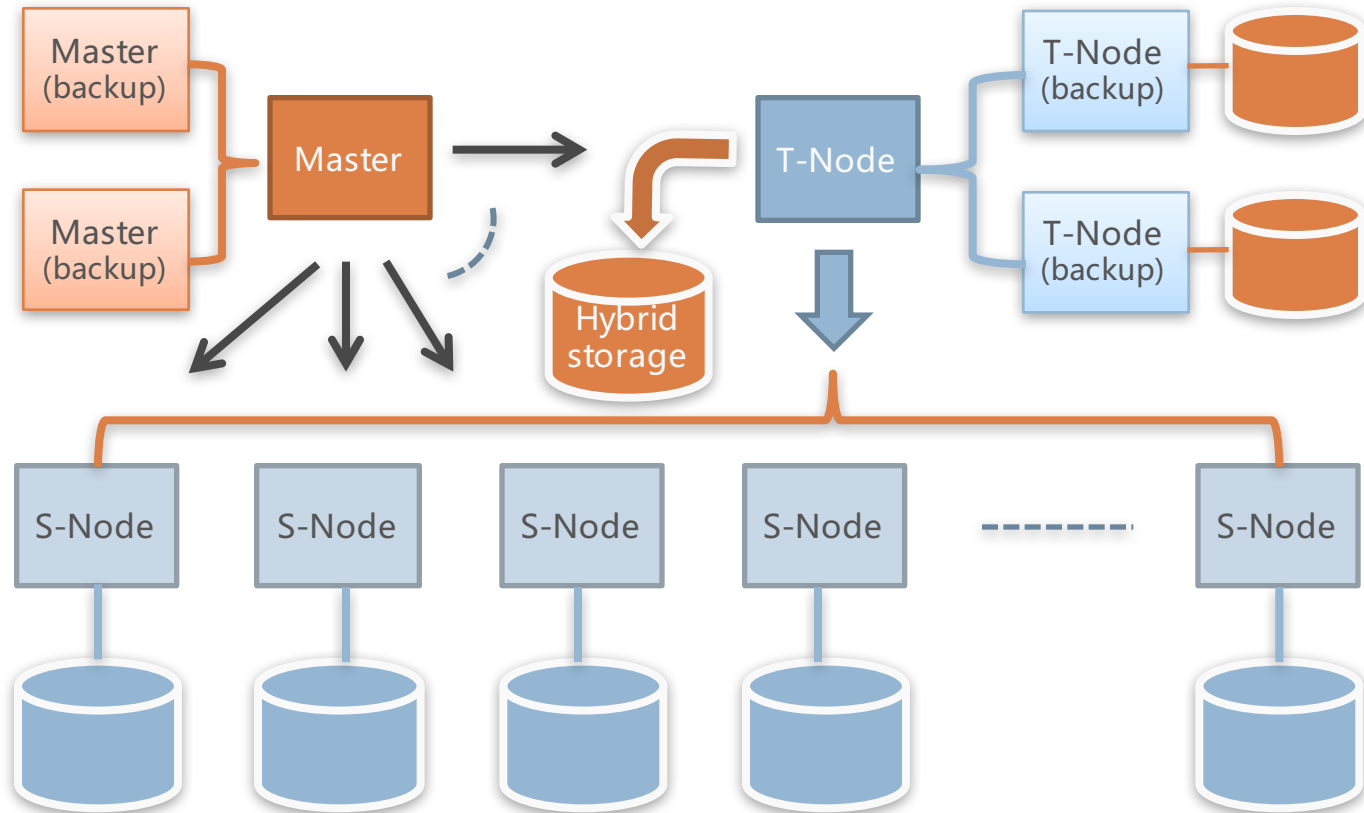
# Features we need

14

- Complex transactions
- High-performance:
  - ⦿ High-throughput
  - ⦿ Low latency
- High-availability
- Scalability
- Elasticity

# Overview (CEDAR core)

15



# Design choices

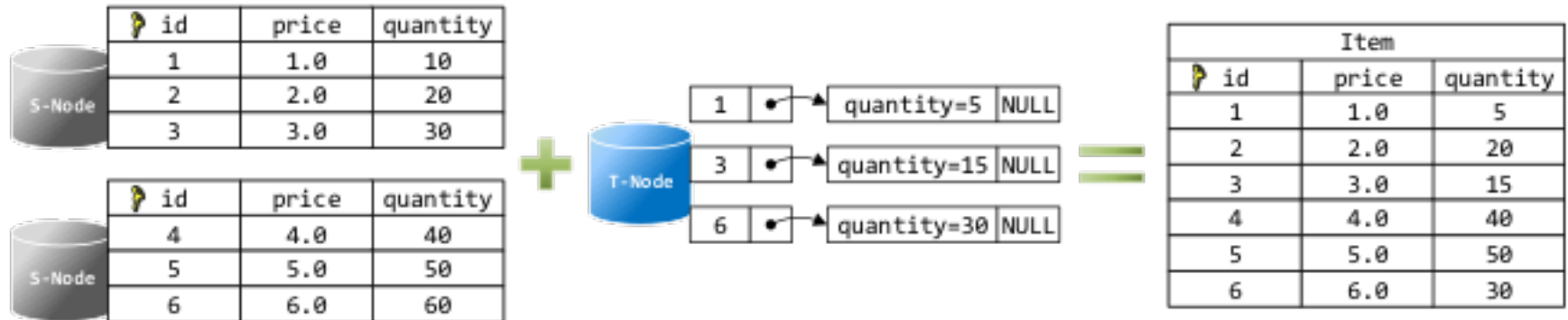
16

- Separating read and write operations on different nodes
- Scalable reading on multiple nodes
- Writing to memory in one node only
  - ⊙ No expensive distributed concurrent control or synchronization
- “Deep” optimization for transactions
  - ⊙ To optimize data transmissions, query execution plans, and executions
- High-availability is guaranteed by log synchronization



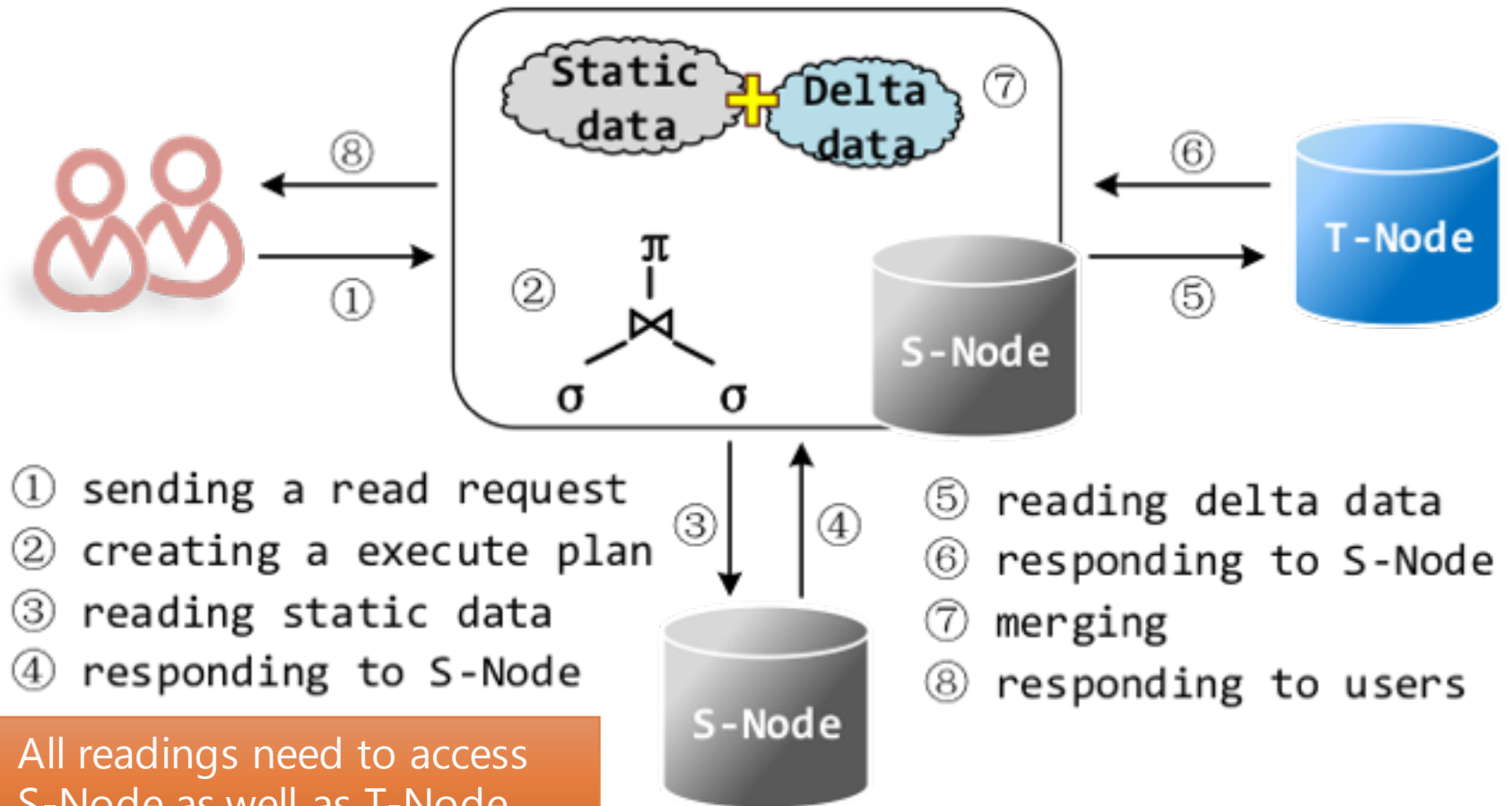
# Status = Baseline + Delta

17



# Reading

18

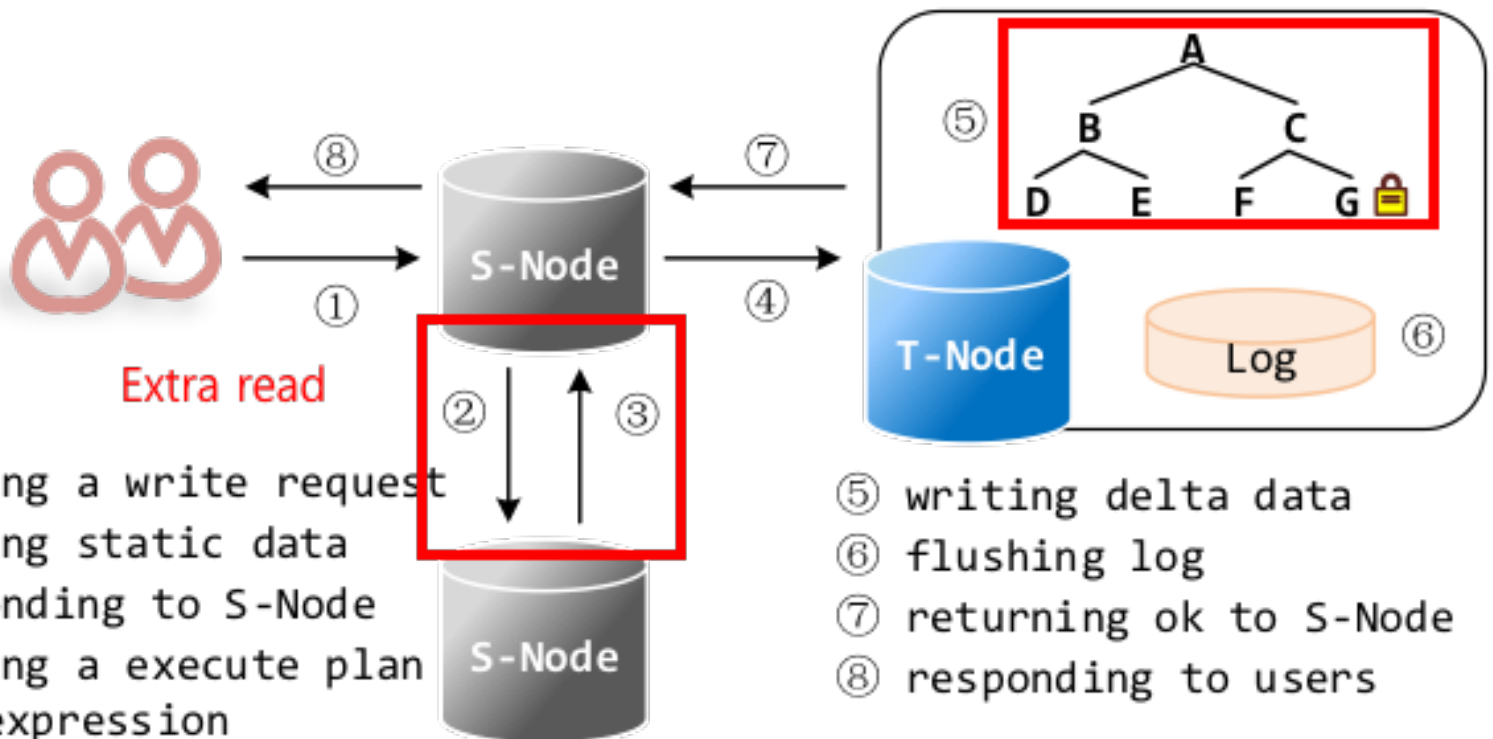


- All readings need to access S-Node as well as T-Node

# Writings

19

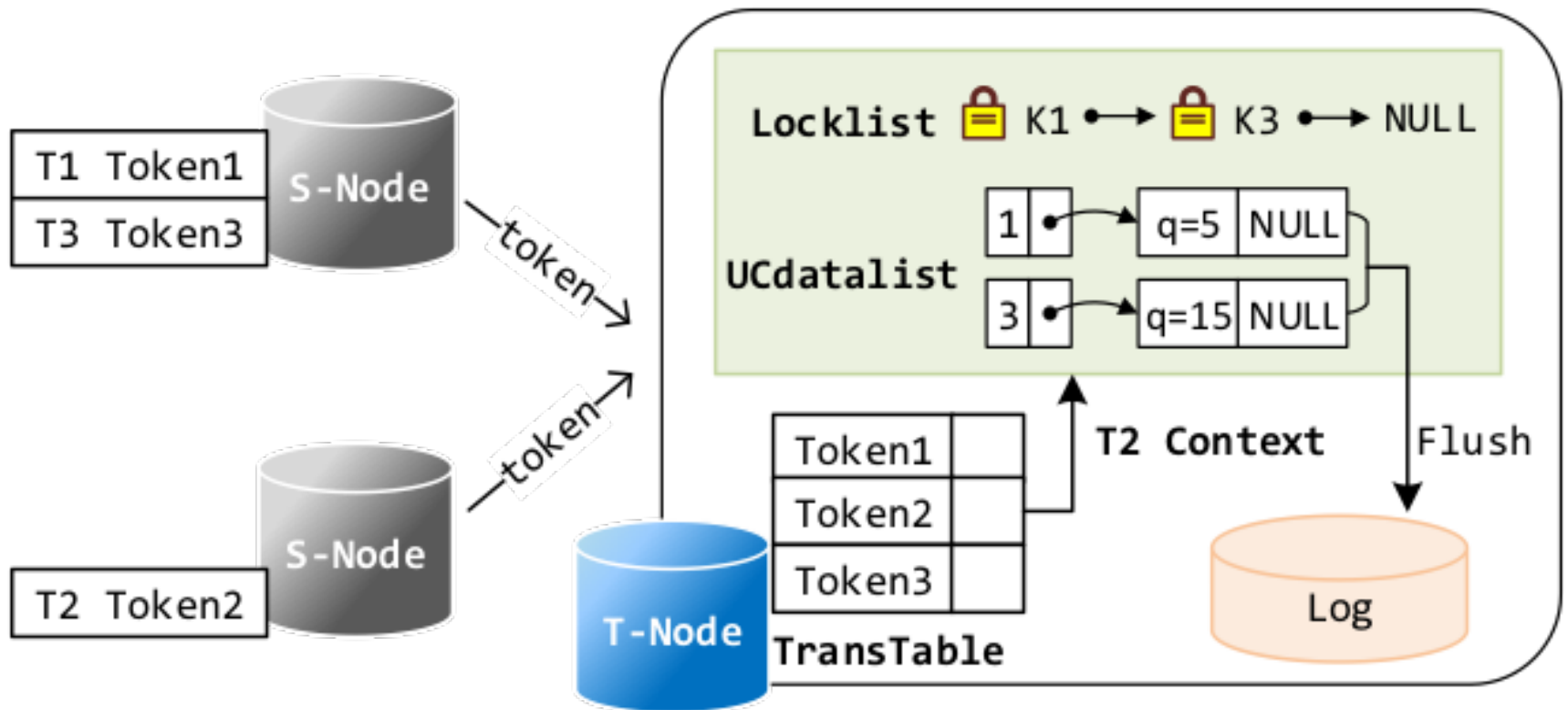
In-memory



- All writings need to access S-Node as well as T-Node

# Transaction management

20



- Transactions are only processed on the T-Node

# Pros and cons

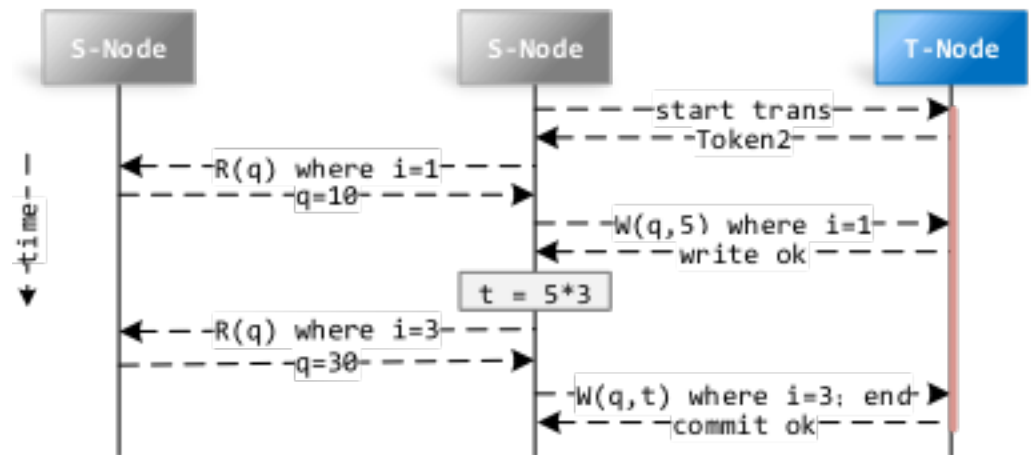
21

## □ Pros

- ⊙ Massive storage
- ⊙ Scalable read
- ⊙ Efficient transaction management

## □ Cons

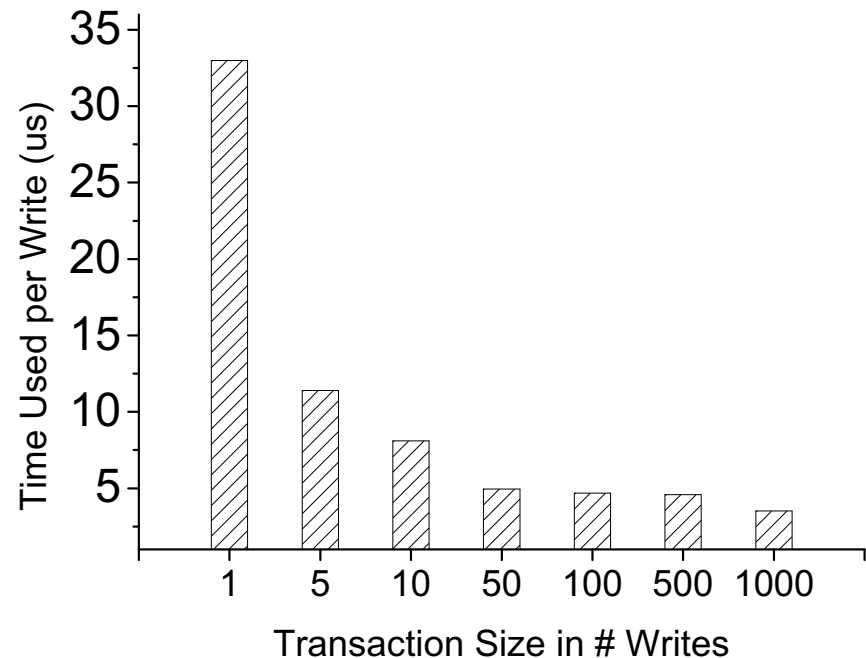
- ⊙ Expensive data transmission



# Performance is affected by

22

- The lengths of locks affect
  - ⦿ Degrees of parallelization
  - ⦿ Latency
- Capability of S-Nodes
  - ⦿ Throughput of readings
- Capability of the T-Node
  - ⦿ Throughput of writings

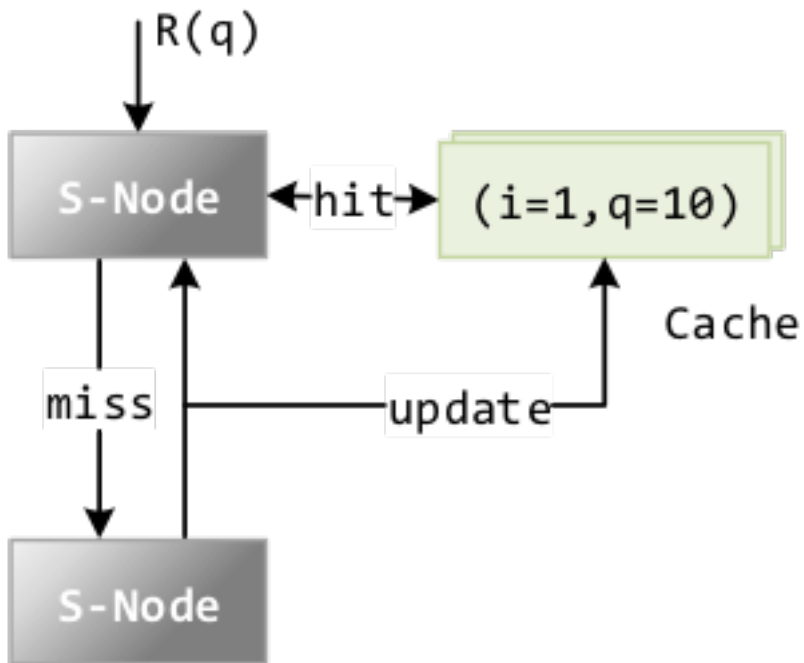


- Most cost is for short/simple transactions
- They are easier to be scheduled

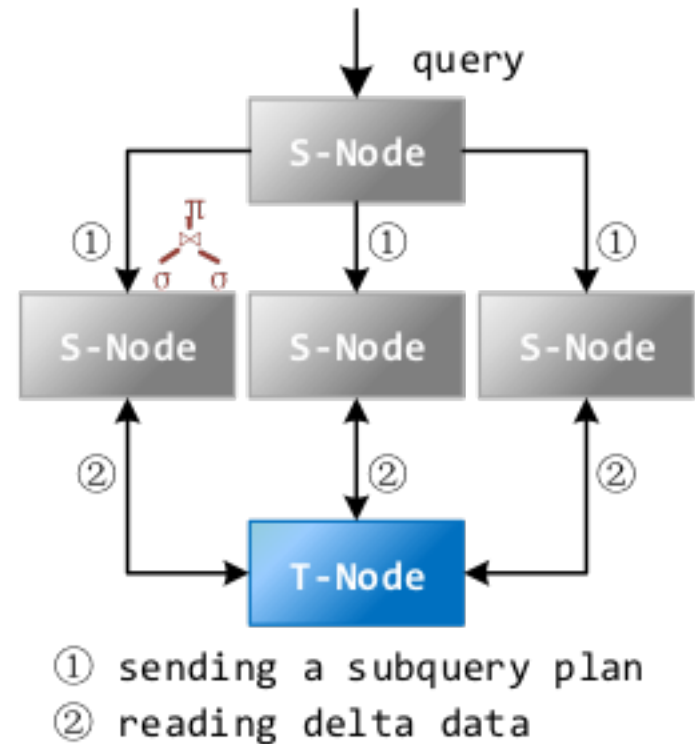
# S-Node optimizations

23

## Static data caching



## Parallel readings



# T-Node


24

## □ Storage node?

- ⊙ All computation resources are for TP
- ⊙ High communication cost

## □ Computation node?

- ⊙ Low communication cost
- ⊙ How much work should it do?

- 
- Balancing T-Node's computation for queries and transaction processing



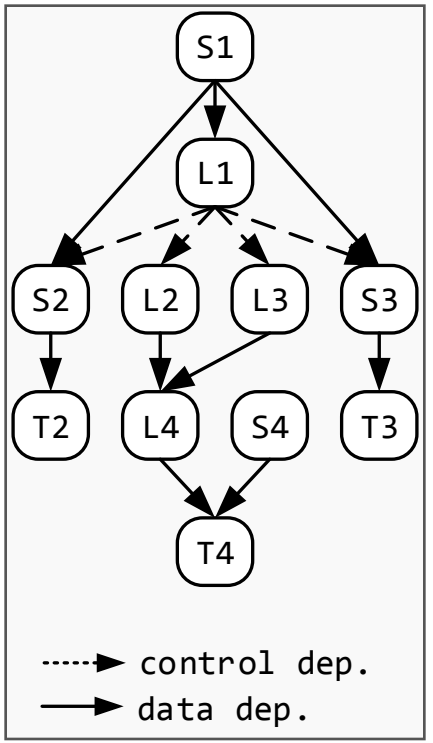
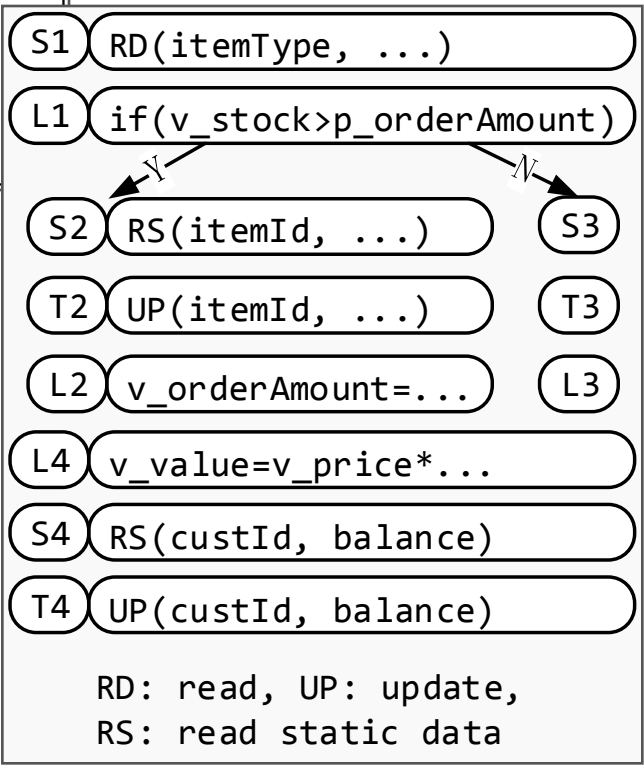
# Transaction compilation

```

Procedure Order(p_itemType int, p_custId int,
               p_orderAmount int)
declare v_price, v_value double;
declare v_itemId, v_stock, v_orderAmount int;
select itemId, price, stock
into v_itemId, v_price, v_stock
from item where itemType = p_itemType
order by price desc limit 1;
if( v_stock > p_orderAmount )
update item set stock = stock - p_orderAmount
where itemId = v_itemId;
v_orderAmount = p_orderAmount;
else
update item set stock = stock - v_stock
where itemId = v_itemId;
v_orderAmount = v_stock;
end
v_value = v_price * v_orderAmount;
update customer set balance -= v_value
where custId = p_custId;
    
```

Transactions

Execution plan

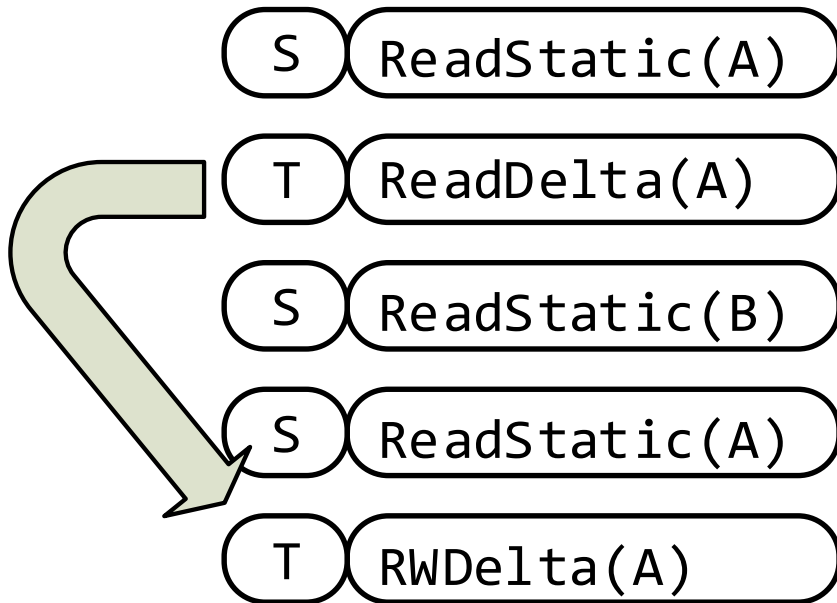


Dependency graph

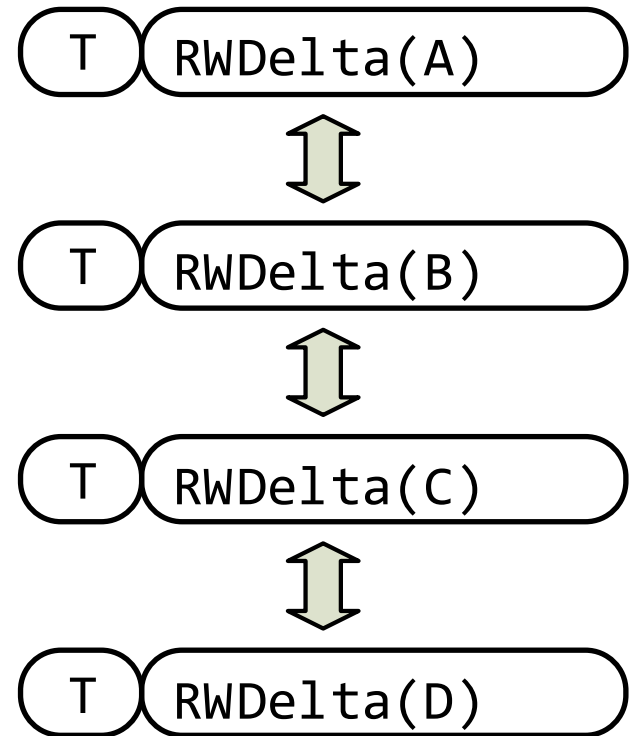
# T-Node optimization

26

## Re-order T-Node ops

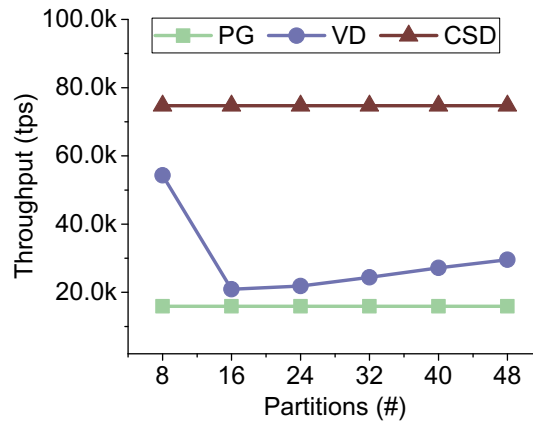
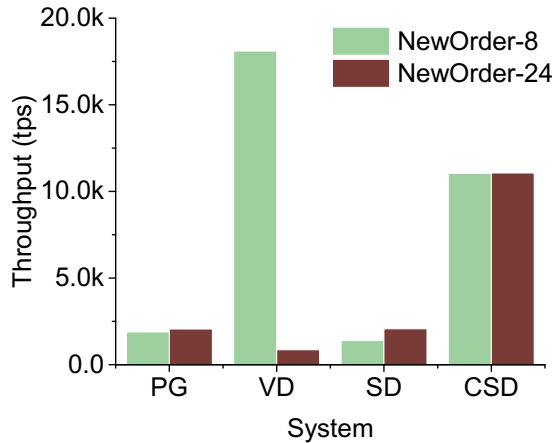


## Postpone conflict ops

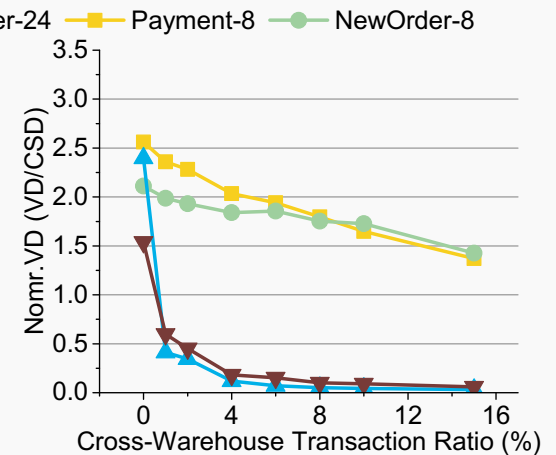
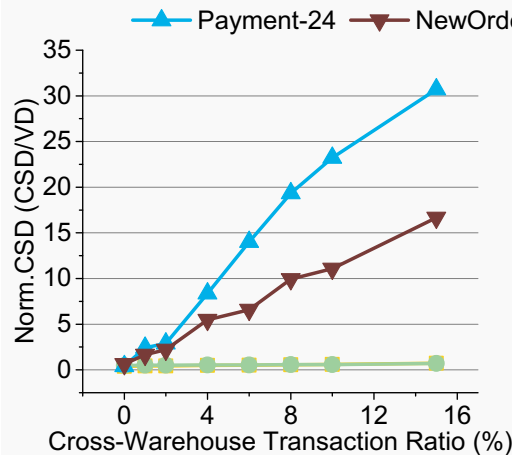
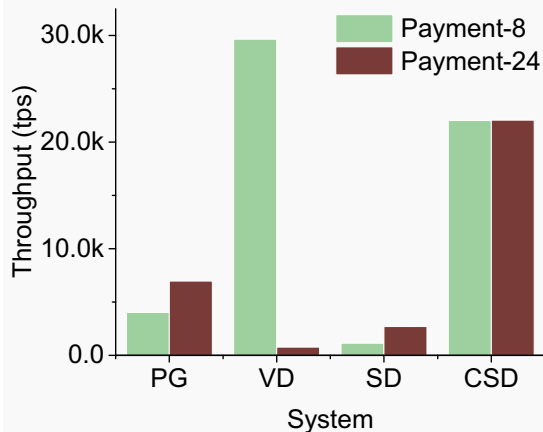


# TPC-C, Smallbank, TATP Benchmarks

27



- Better than PG
- Better than VoltDB under complex workload
- With significant advantage when data cannot be naturally partitioned



# Indexing

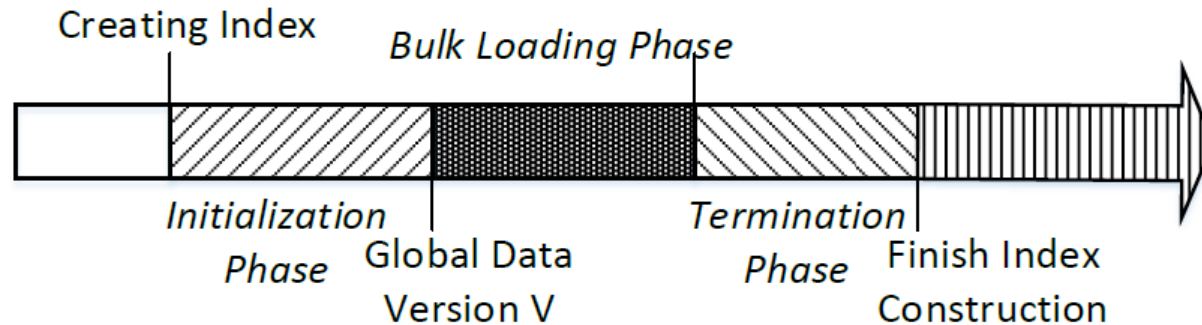
28



- Index is organized as a table.
- No distributed transactions.
- Taking advantage of the load balancing, availability of system.

# Indexing overview

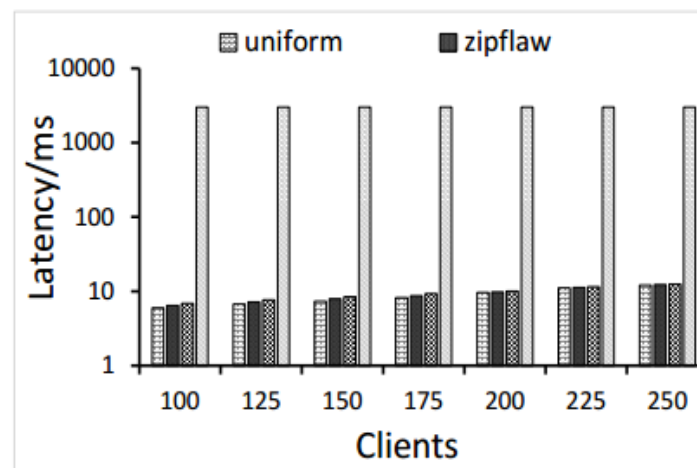
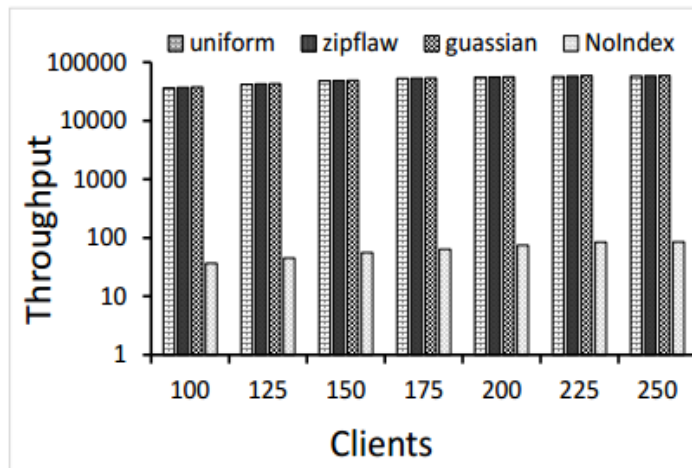
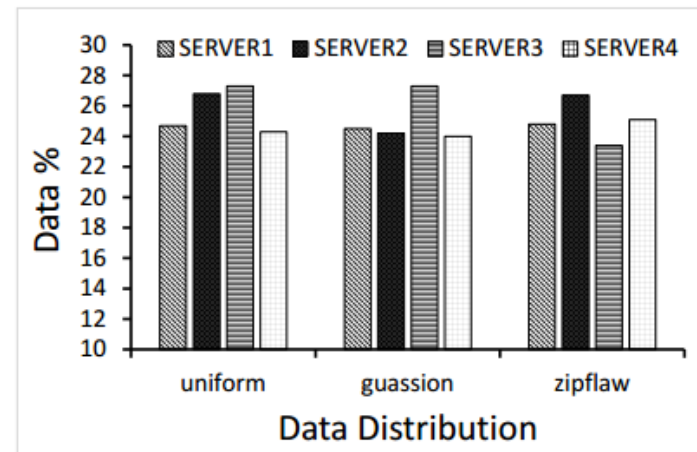
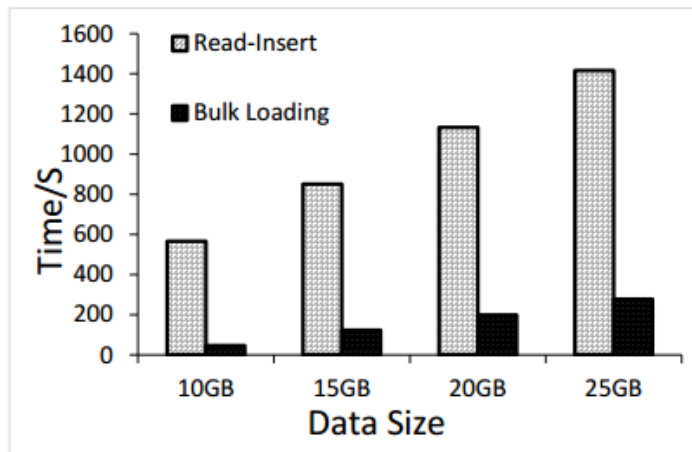
29



- Initialization: Preparing for the start of index construction.
- Bulk Loading
  - Local processing: collecting statistical information
  - Global processing: achieving load balancing based on an equi-depth histogram
- Termination: Scheduling the task for replication of the index for high availability.

# Experimental results

30



# Other works on the CEDAR

31

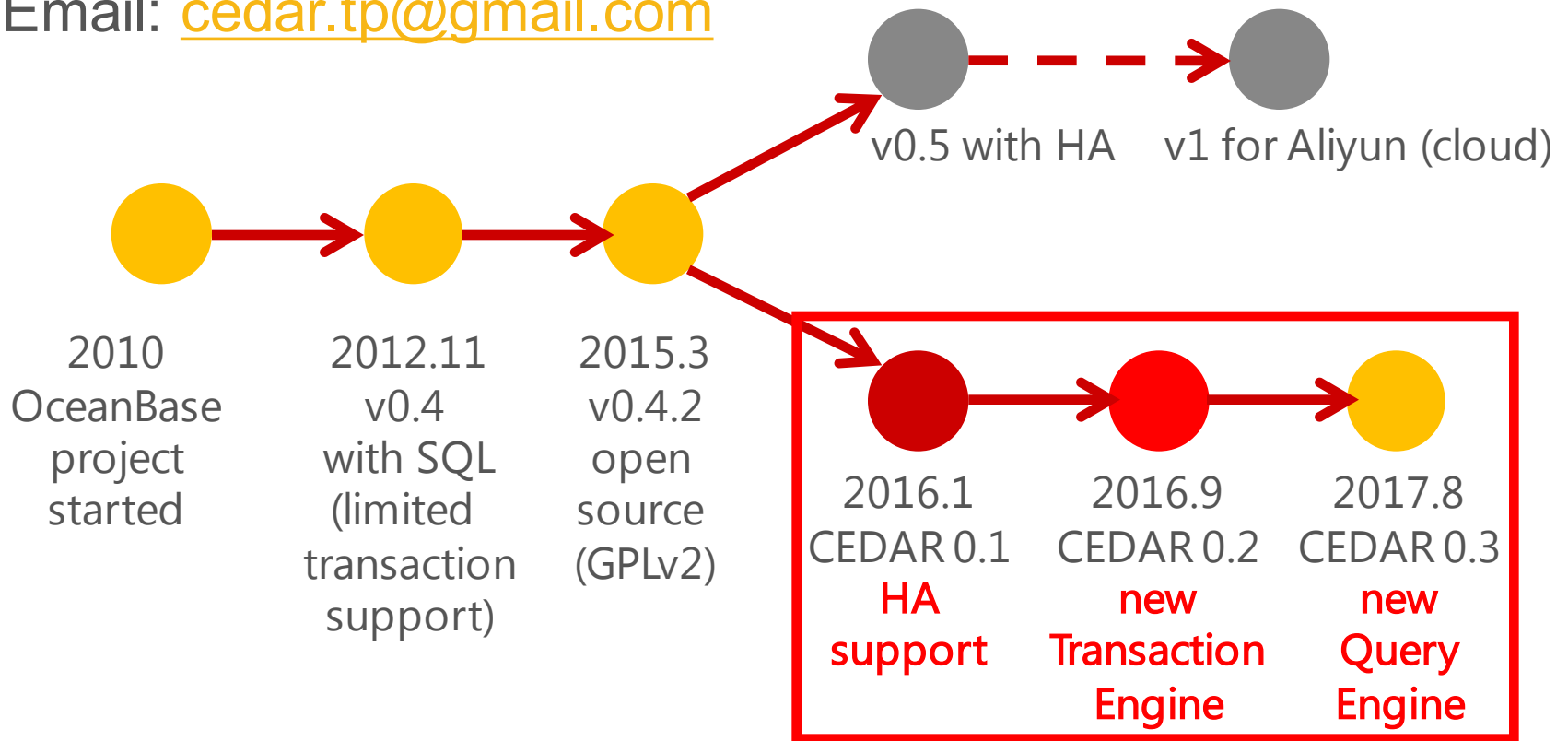
- Scalable range optimistic concurrency control
- Global snapshot isolation with Paxos replication
- Scalable commit log recording/synchronization
- Data transmission optimization
- Distributed statistics monitoring
- Rule-based and cost-based query optimization

# Release

32

Homepage: <https://github.com/daseECNU/CEDAR>

Email: [cedar.tp@gmail.com](mailto:cedar.tp@gmail.com)





# Homepage:

<https://github.com/daseECNU/CEDAR>

33

CEDAR是华东师范大学数据科学与工程学院（简称“DaSE”）基于 OceanBase 0.4.2 研发的可扩展的关系数据库。2016年2月1日，CEDAR项目组完成了CEDAR 0.1 版本的开发与测试，2016年9月26日，CEDAR 0.2 版本发布。

## 版本特性

CEDAR在OceanBase 0.4.2 的基础上新增了如下11个功能模块：

CEDAR 0.1 版本新增的功能有：

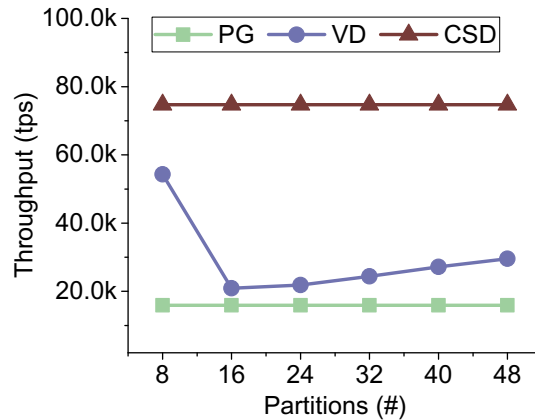
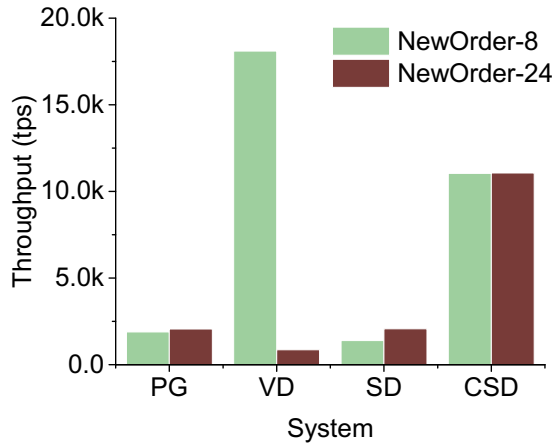
- 高可用的三集群架构（集群间选主、集群角色自动切换、日志强同步及恢复等机制）
- 多线程网络IO处理框架Libonev
- 游标
- 存储过程
- 二级索引
- 非主键多行更新
- 半连接

CEDAR 0.2 版本新增的功能有：

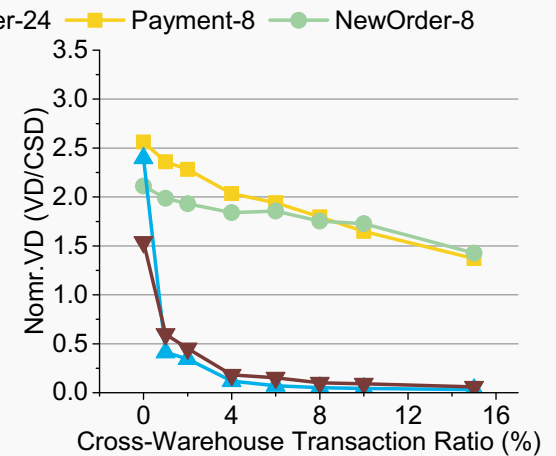
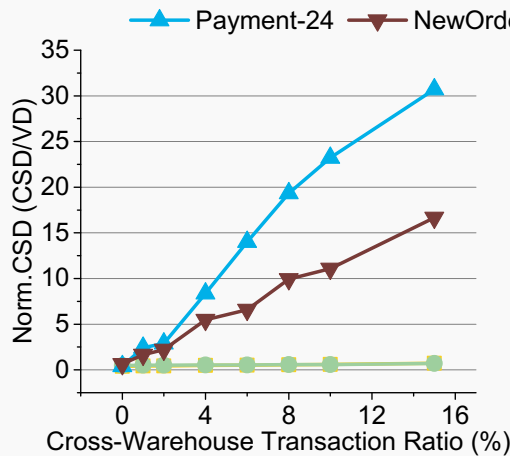
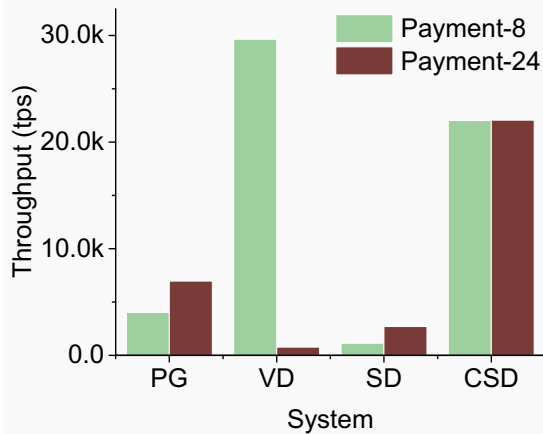
- SNAPSHOT ISOLATION 隔离级别
- 表锁
- 基于布隆过滤器的连接
- 日志同步优化

# TPC-C, Smallbank, TATP Benchmarks

34



- Better than PG
- Better than VoltDB under complex workload
- With significant advantage when data cannot be naturally partitioned



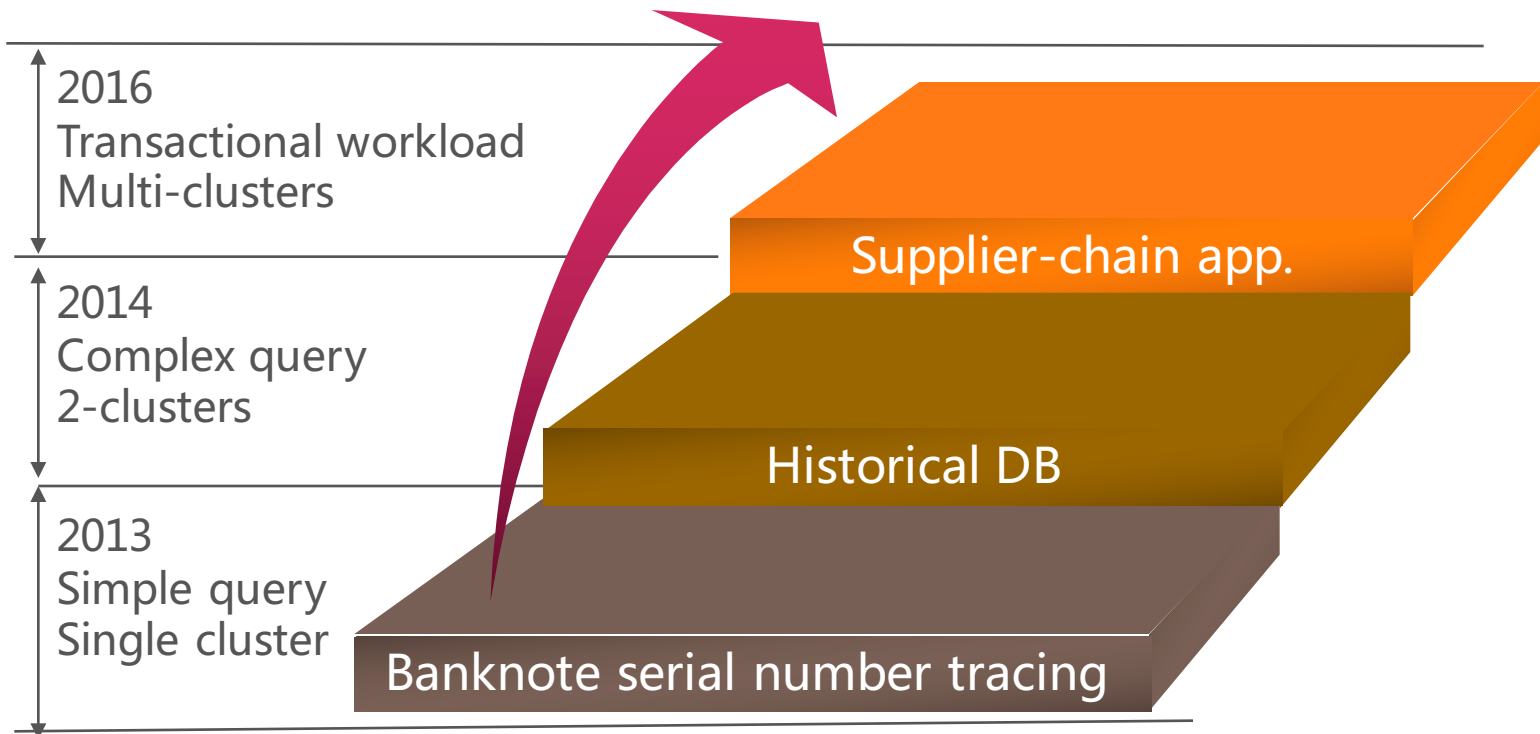
# Application

35

- One of the largest banks in China
  - ⊙ To use local/open-source DBMS for transactional applications (to replace DB2)
  - ⊙ Three stages: historical DB => Hybrid DB => transactional DB
  - ⊙ 2013 – present (in stage 2)
  - ⊙ Code name: CBase within the bank

# Stage 1 => Stage 2

36



# Banknote serial number tracing

37

- All records of banknote serial numbers are stored for 90 days in the CBase
- More than 10TB data totally
- For counterfeit money detection etc.



# Status

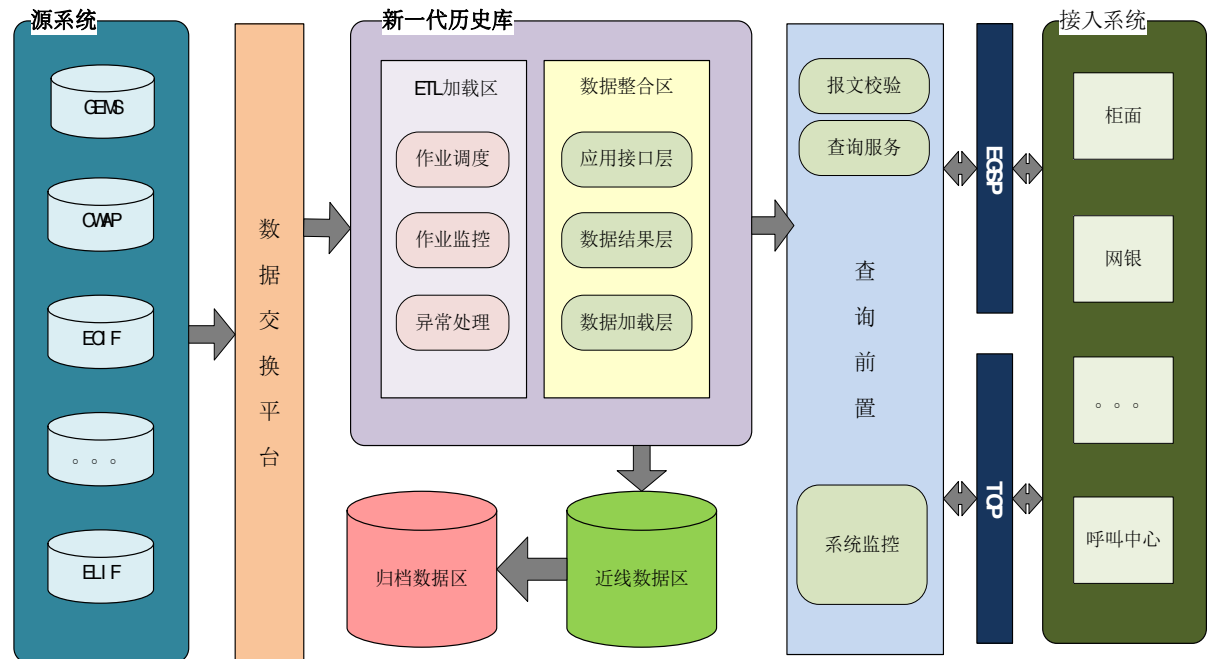
38

- Went online in 2013-12
- Single cluster with 9 servers
- Features
  - ⊙ Batch load with about 100GB data/day
  - ⊙ Latency in several milliseconds over TBs of data
  - ⊙ Linear scalability with respect to number of servers

# Historical DB

39

- Move historical data to historical DB so that the online system is lightweight
- 283 kinds of query services are supported by CBase



# Status

40

- Went online in 2014-09
- 2-clusters, with 21 servers each
- Features
  - ⊙ About 400GB data/day
  - ⊙ Latency in several milliseconds for most queries
  - ⊙ Active/active high-availability
- Has become the single-point sign-in gateway of nearly all business logic

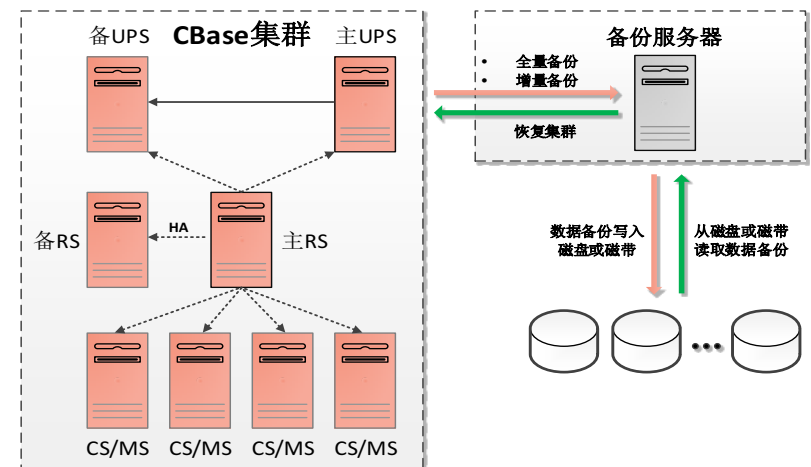
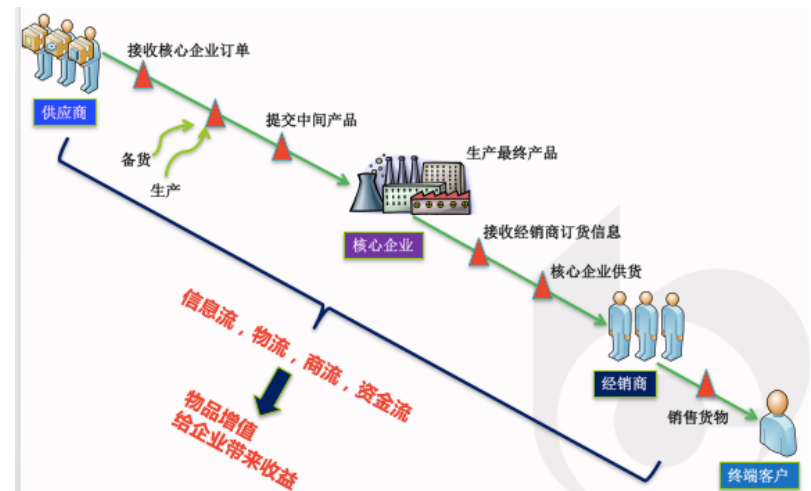


# Supplier-chain applications

41

## □ Hybrid workload

- ⊙ Complex transactions with tens or even hundreds of SQL operations
- ⊙ Complex analytical queries with many joins of large tables
- ⊙ High-availability requirements



# Status

42

- Went online in 2017-03
- Multi-clusters, each with 12 servers
- Features
  - ⊙ Scale-out like most NoSQL systems
  - ⊙ 5000tps for complex workload
  - ⊙ 5ms latency for key-search
  - ⊙ Complex queries are answered within 3s
  - ⊙ Multi-clusters are synchronized with Paxos-like protocol to provide high-availability

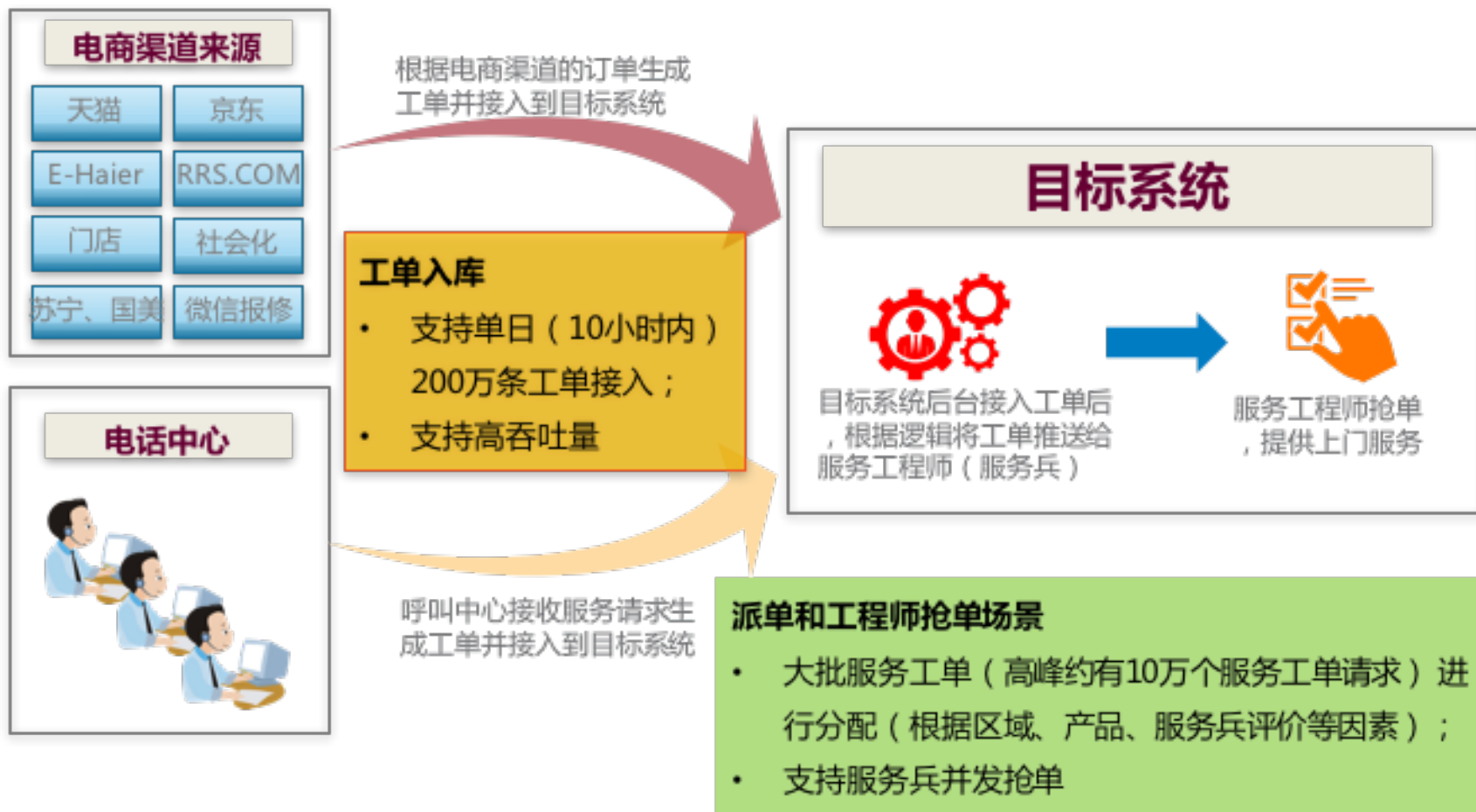
# More applications this year

43

- Network Alliance (网联) e-payment clearinghouse
  - ⊙ 2017-03 to 2017-04
- Loaning
  - ⊙ 2017-09 to 2017-12

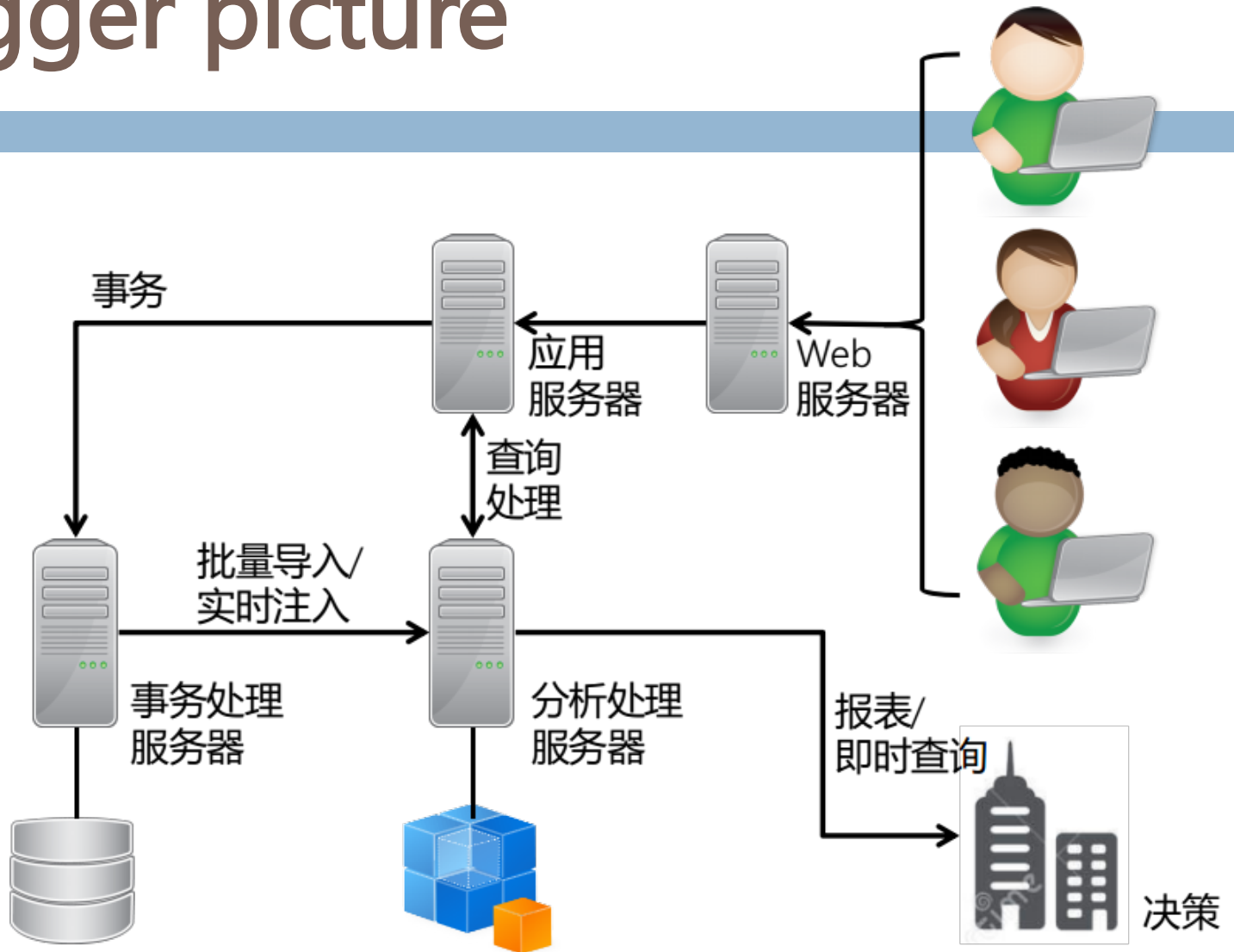
# POC: O2O task assignment/taking

44



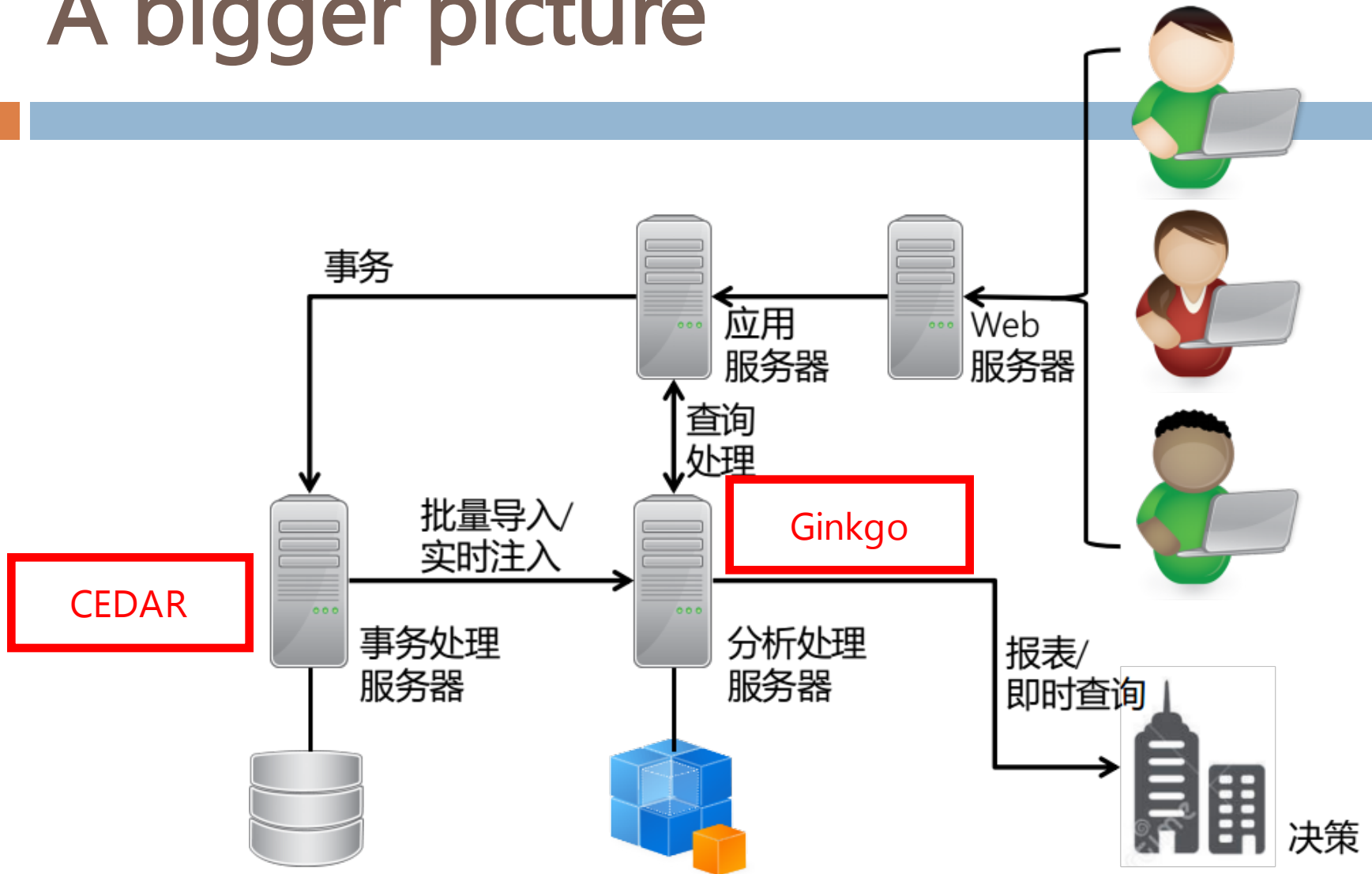
# A bigger picture

45



# A bigger picture

46



# Summary

47

## CEDAR: 雪松

- **C**: **C**luster-oriented
- **E**: for **E**nterprise applications
- **D**: scalable **D**bms
- **AR**: non-traditional **AR**chitecture



# Summary

48

- **SQL** support: ODBC/API interfaces
- **Transaction** support: ad-hoc transactions and store-procedures
- **Efficient** query optimization/execution: various distributed join implementation + indexing schemes
- Deployment with **High-Availability** support
- **Highly Scalable:**
  - ⊙ Read/write separation
  - ⊙ Hot/cold separation
- **Management/maintenance-friendly:** Import/export toolkit and monitoring/diagnose toolkit



# Summary

49

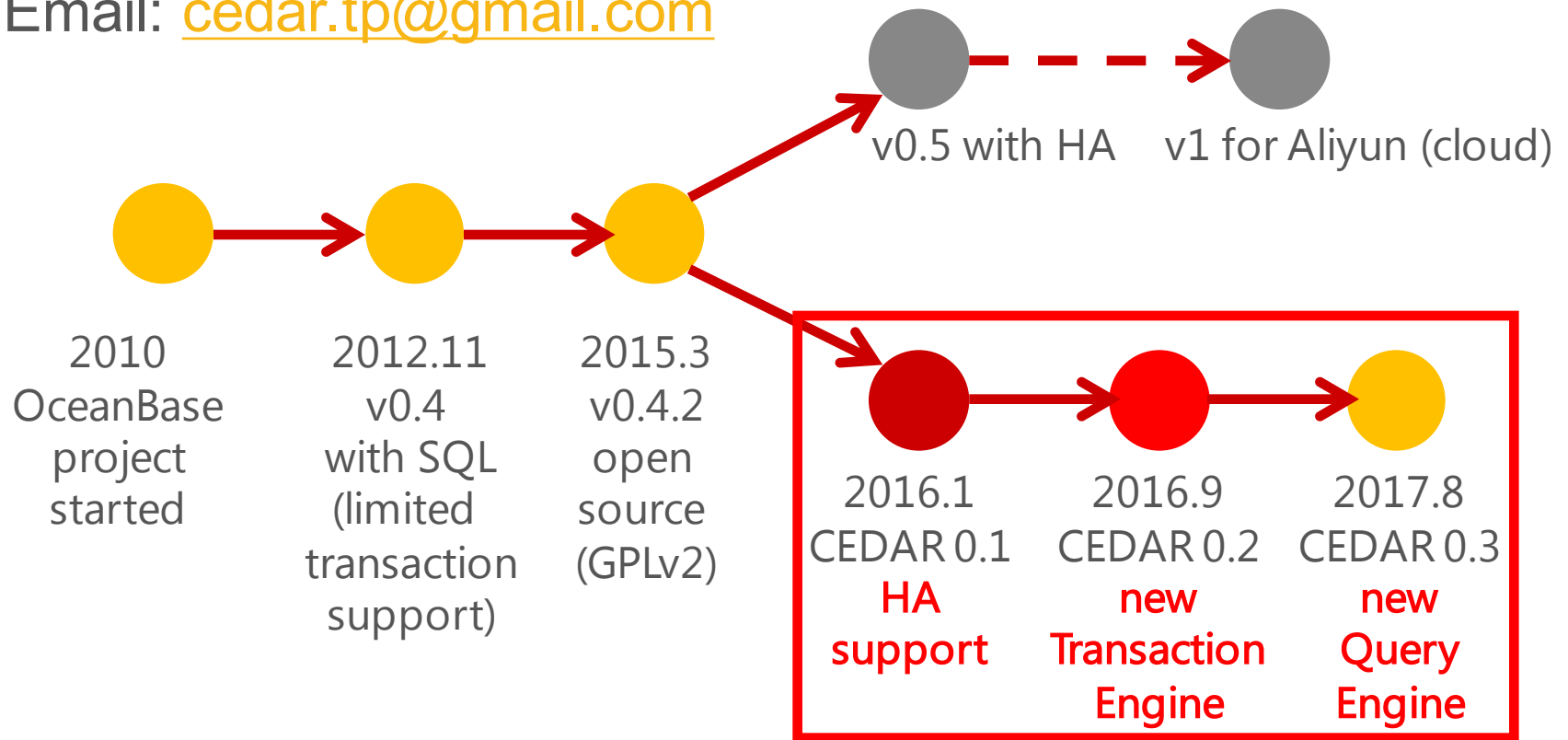
- **Full-fledged DBMS**
  - ⦿ with **SQL** and **Transaction Processing** support
- **Scable** architecture
  - ⦿ cluster-oriented: commodity PC server with large memory, SSD drive, and high-speed network
- **Mission-critical**-app.-oriented
  - ⦿ apps in enterprises, banks, communications, etc.
- GPLv2

# Thanks!

50

Homepage: <https://github.com/daseECNU/CEDAR>

Email: [cedar.tp@gmail.com](mailto:cedar.tp@gmail.com)



# Acknowledgement



華東師範大學 中國人民大學 印孚瑟斯

数据科学联合实验室

ECNU-RUC-Infosys Data Science Joint Lab

**Thanks!**