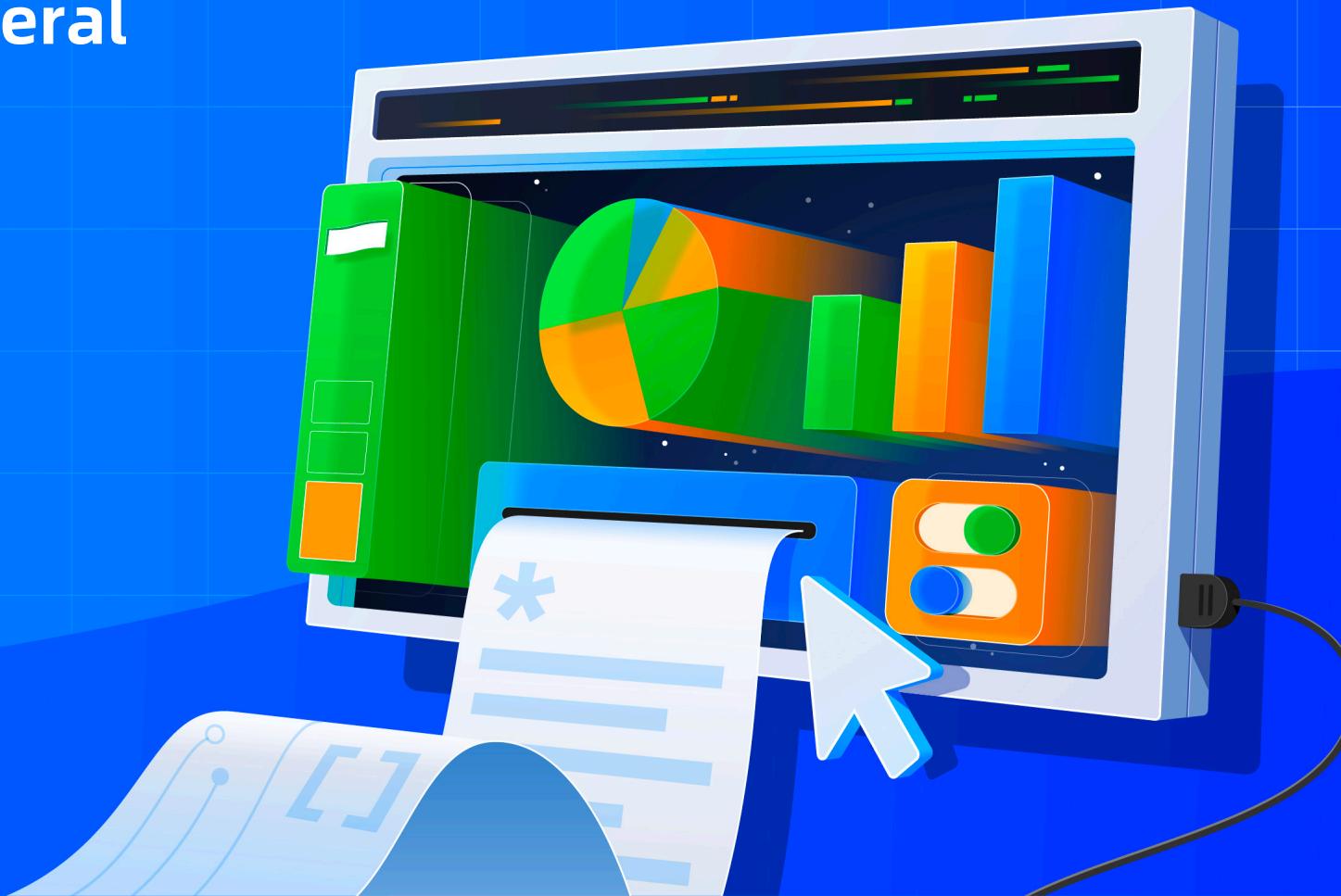


FROM INTRODUCTION TO PRACTICE

Lesson 1: OceanBase General Introduction

Zhongyan Feng

OceanBase OpenSource Team Leader



Contents

- 
- 01 OceanBase Brief Introduction**
 - 02 Architecture and Conception**
 - 03 Key Features**
 - 04 Best Practices**

01

OceanBase Brief Introduction

OceanBase – The only Database which break the world record in TPC-C and TPC-H

TPC Transaction Processing Performance Council

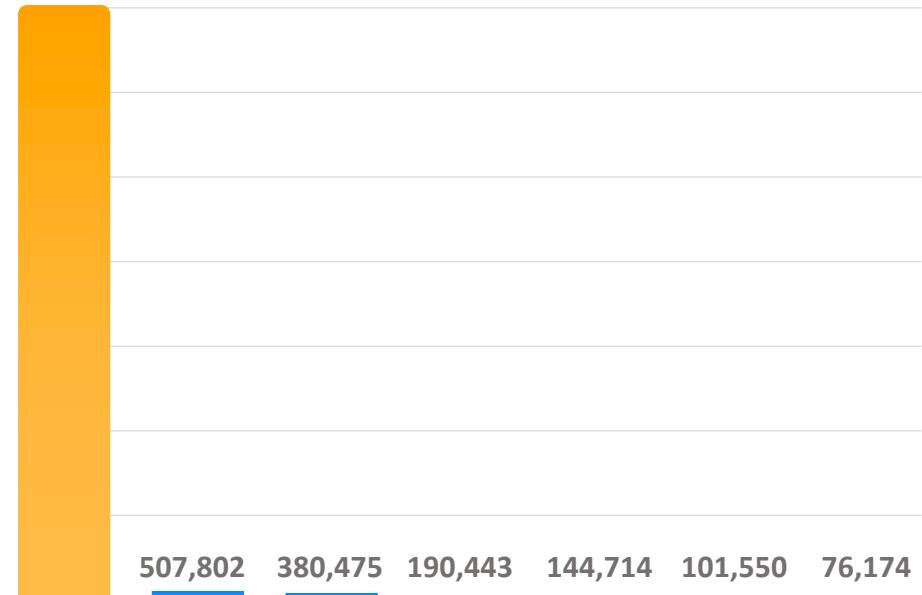
- TPC-C No.2 official world record

TPC-C - Advanced Filtered and Sorted Results				
Sponsor	System	Performance (tpmC)	Price/TpmC	System Availability
ANT FINANCIAL	Alibaba Cloud Elastic Compute Service Cluster	707,351,007	3.98 CNY	6/8/2020
ANT FINANCIAL	Alibaba Cloud Elastic Compute Service Cluster	60,880,800	6.25 CNY	10/2/2019
Oracle	SPARC SuperCluster with 13-4 Servers	30,249,688	1.01 USD	6/1/2011
IBM	IBM Power 780 Server Model 9179-MHB	10,366,254	1.38 USD	10/13/2010
Oracle	SPARC T5-8 Server	8,552,523	0.55 USD	9/25/2013
Oracle	Sun SPARC Enterprise T5440 Server Cluster	7,646,486	2.36 USD	3/19/2010
IBM	IBM Power 595 Server Model 9119-FHA	6,085,166	2.81 USD	12/10/2008
Bull	Bull Escala PL6460R	6,085,166	2.81 USD	12/15/2008
Oracle	Sun Server X2-8	5,055,888	0.89 USD	7/10/2012

- 20 times faster than Oracle

- TPC-H No.3 official world record (30TB dataset)
- On May 19, 2021, **15.26** million QphH@30000GB

OceanBase
707,351,007



* Source: tpc.org

More than 1000 Customer's Choice.

Financial



Government Telecom



Internet



Manufacturer



OceanBase @ Ant Group



"Double Eleven" Performance: 61M QPS
Peaks 544,000 (Alipay Double Eleven) vs Peaks 64,000 (Visa)

61 millions

Responses per second

Peak processing power

> 200

Nodes in one cluster

Cluster Size

> 6

PB data

Data size in one instance

> 320 billions

Billions of rows

One single table size

RPO = 0
RTO < 8

Seconds

Disaster Tolerance

*Statistics above are from real production in Alipay

OceanBase : The Open Source Distributed Database for Mission-critical Workloads at Any Scale

1.0

Firmly embrace the distributed architecture

2010

Product initiation

Taobao

2013

Internal customers

TMALL天猫
Taobao

2014

Support core transactions

支付宝

2.0

Native distributed SQL database

2016

Full service coverage

 余额宝
 芝麻信用
 网商银行

2017

External customers

 四川省农村信用社
 顺德农商银行
 南京银行

2019

Break the TPC-C world record

 中国人民保险
 招商证券
 常熟农商银行
 西安银行

3.0

Unified engine, hybrid deployment

2020

Commercialization

 中国工商银行
 中国移动
 中国石化
 中华保险

2021

Massive adoption

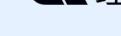
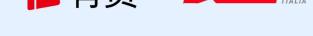
 江西省人力资源和社会保障厅
 交通银行
 国家电网
 携程
 DANA
 中国人寿

4.0

Standalone & distributed integrated architecture

2022 ~

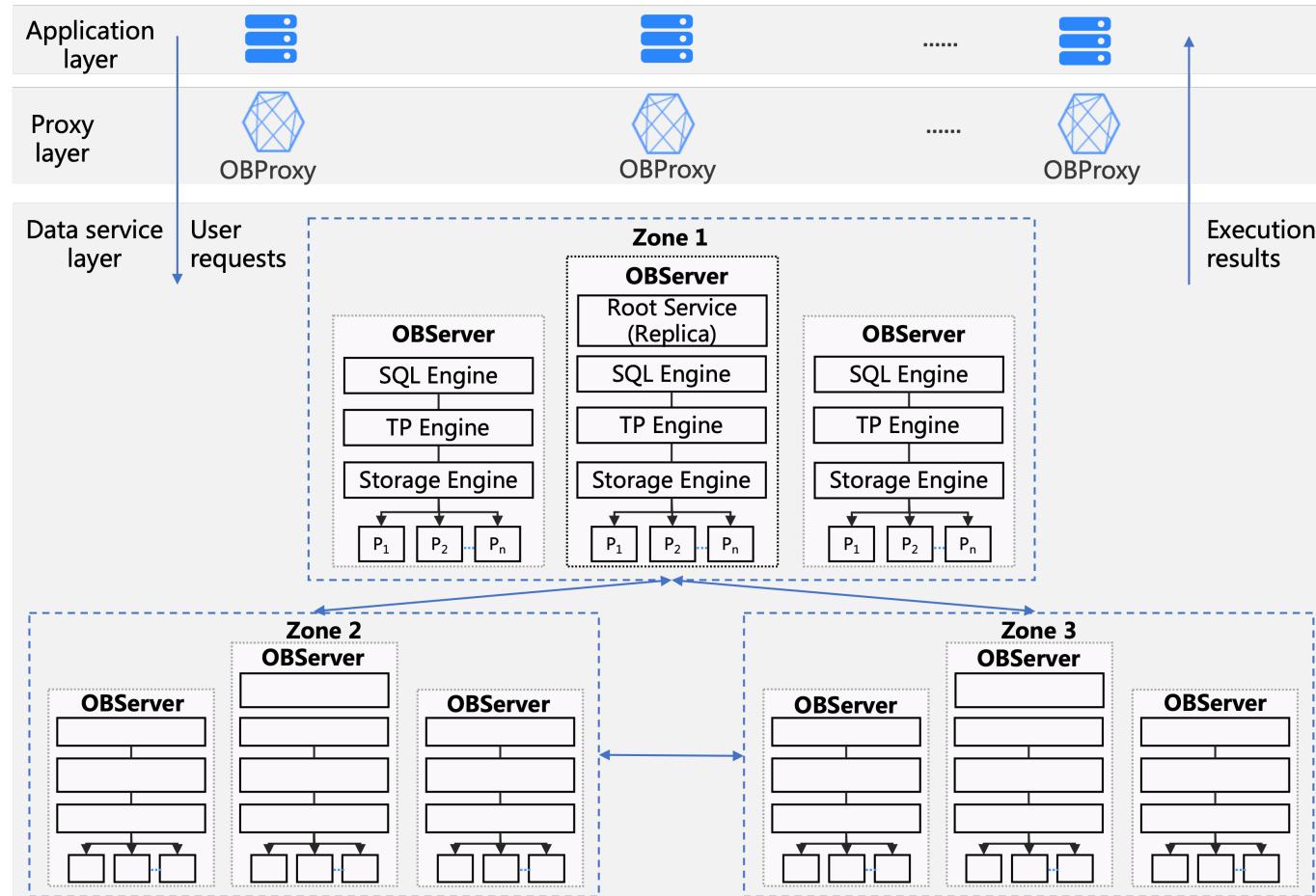
Provide cloud service & go global

 理想
 致欧
 有赞
 TELECOM ITALIA
 GCash
 easypaisa

02

Architecture and Conception

Architecture Overview



- **No single failure point, every node are equal**
- **All OBServers can read and write**
- **Cluster can be deployed on Private Cloud/Public Cloud/Hybrid Cloud**
- **ObProxy is used to route requests**
- **LSM-tree like storaged engine**
- **SQL Engine: local, remote and distributed execution**

Hybrid Deployment



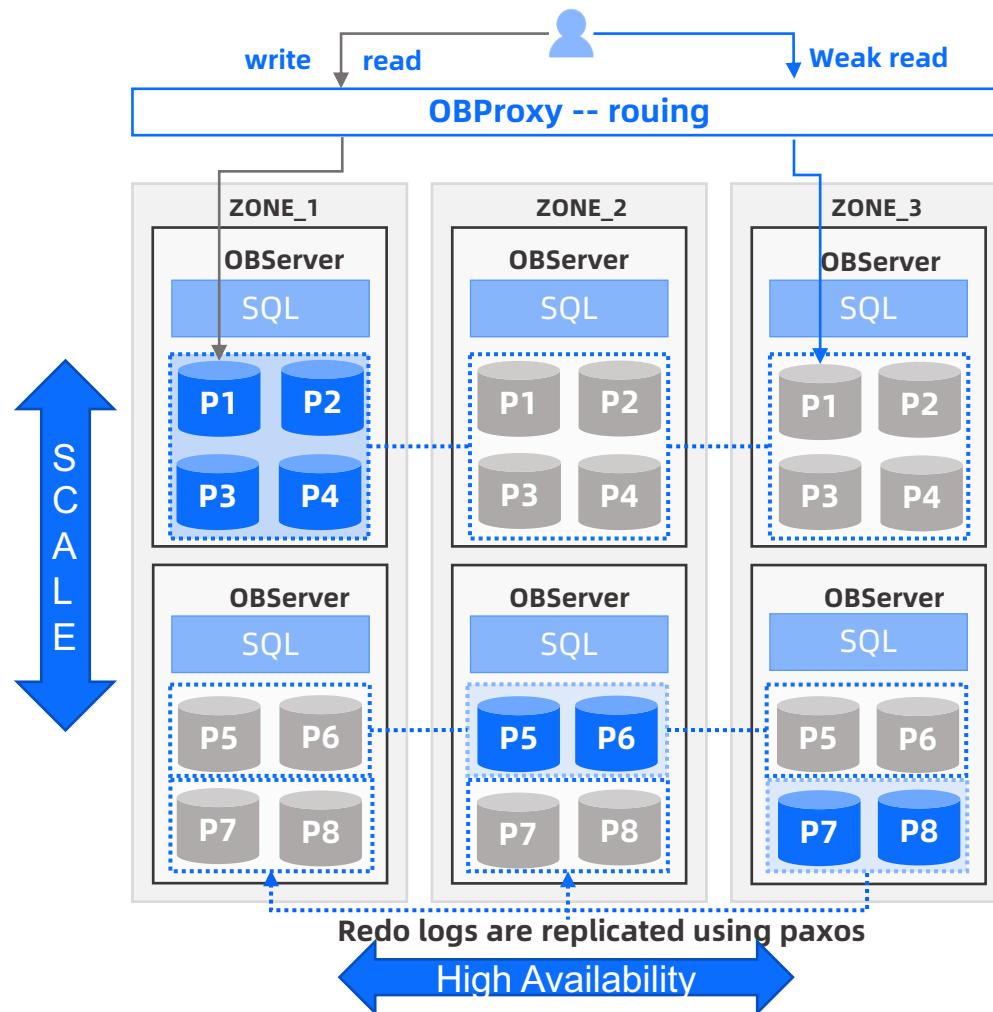
On-Premise

Private Cloud



Public Cloud /Hybrid Cloud

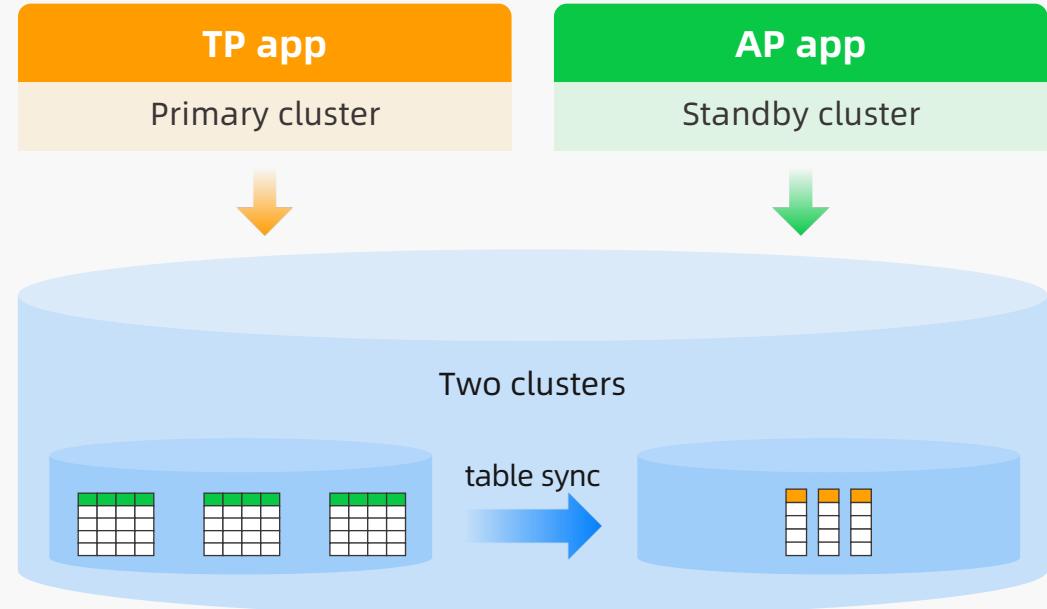
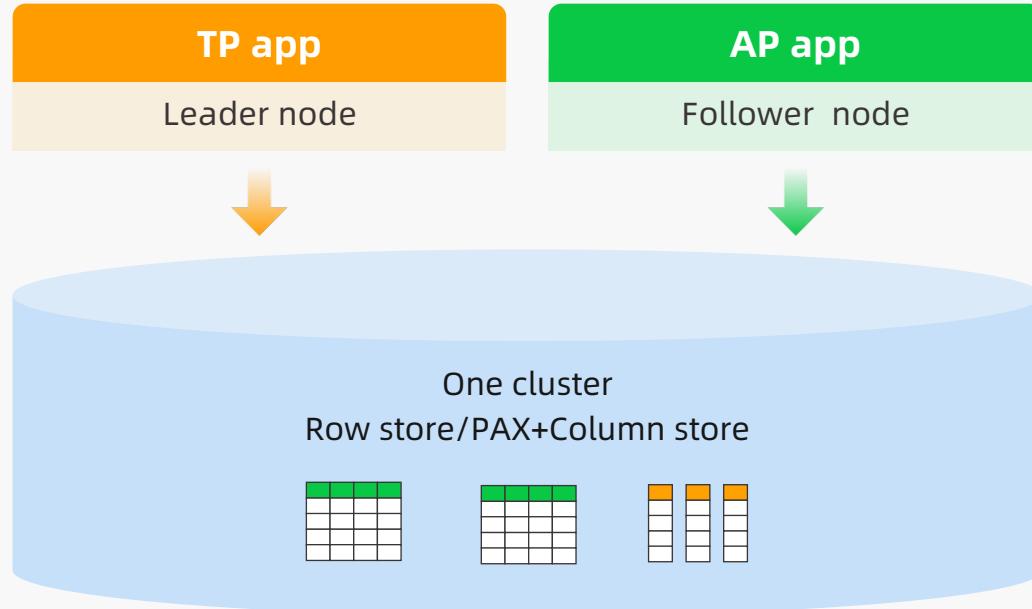
Architecture Overview (2)



- Cluster is composed of zones
- One zone means one replica
- One zone can contain one or more servers
- Data are distributed by Partitions
- The same partition in different zones form Paxos group
- One partition will be elected as leader in one Paxos Group.
- The leader will provide read/write service.
- Redo logs are replicated using Paxos
- Transactions across servers are executed using 2PC
- Transactions on one node, even multiple partitions, local transaction in single node, reduce distributed transaction

One Storage Engine support both Row and Column Store

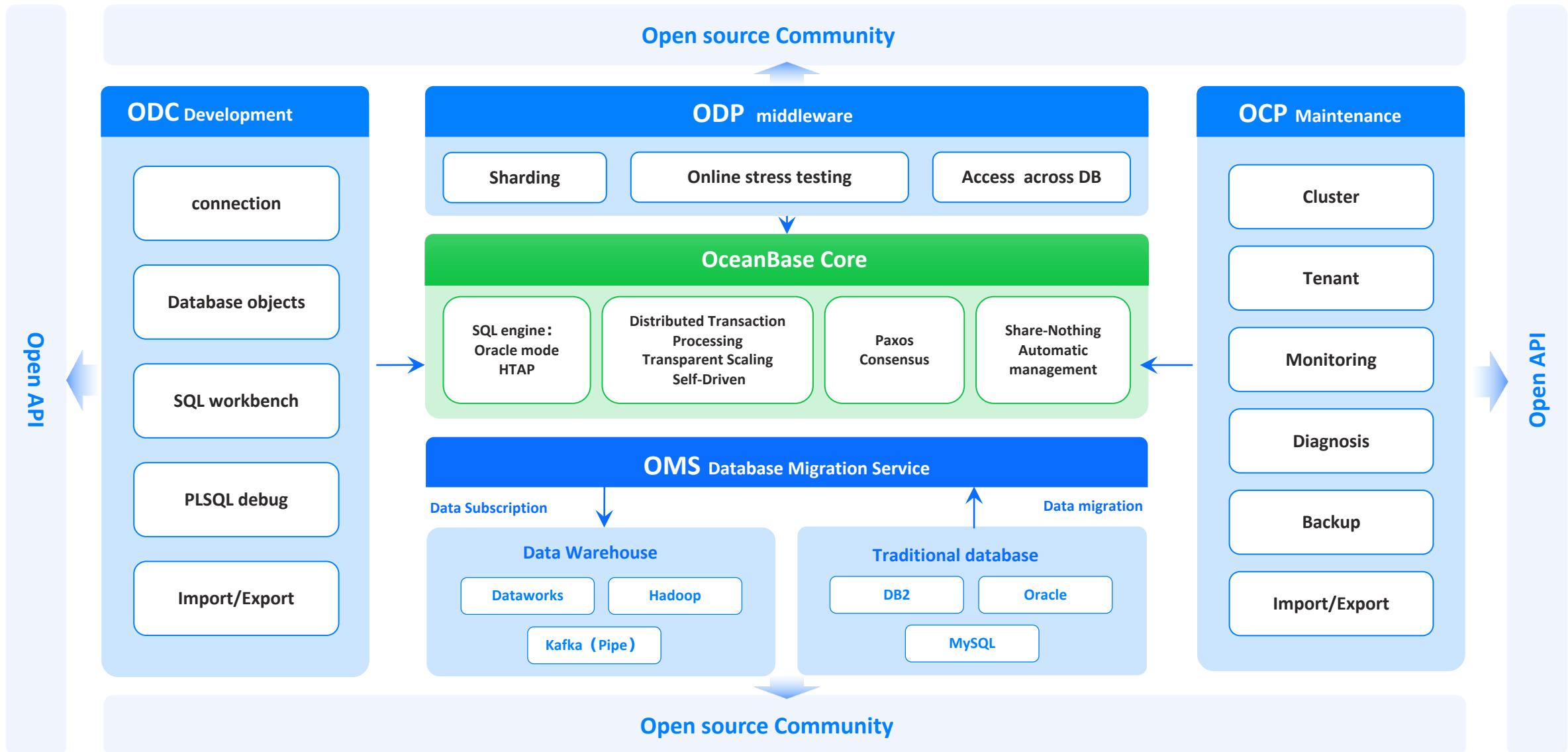
One unified system for both TP & AP



- Low cost, high cost-effectiveness
- Millisecond-level latency and consistency between leader and follower nodes

- High cost, low cost-effectiveness
- Significant latency and data inconsistency between primary and standby clusters

Product Family 1 + 4 Full Stack Solutions



03

Key Features

Features

**High
Availability**

Scalable

**Easy
Operation**

High Performance

Cost Effective

**Compatible
with MySQL**

How to provide 7 x 24 Service?

Key Feature: High availability

High Availability



Flexible DR Tolerance

3 replica in 1
IDC

3 IDCs in 1
city

3 IDCs in 2
cities

Dual IDCs in master/standby
mode

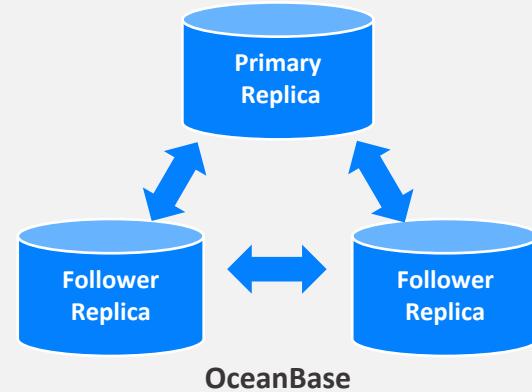
5 IDCs in 3
cities

RPO=0 RTO<8s



Monolithic database

- **Data Loss:** Cannot guarantee consistency in case of failure
- **FO cost:** High maintenance cost for cold standby



- **Zero data loss:** Consensus protocol based on Multi Paxos
- **Unattended HA:** Available when a minority of nodes fails RPO = 0, RTO < 8s
- **Peer Node:** All replicas are active

Standby cluster/tenant

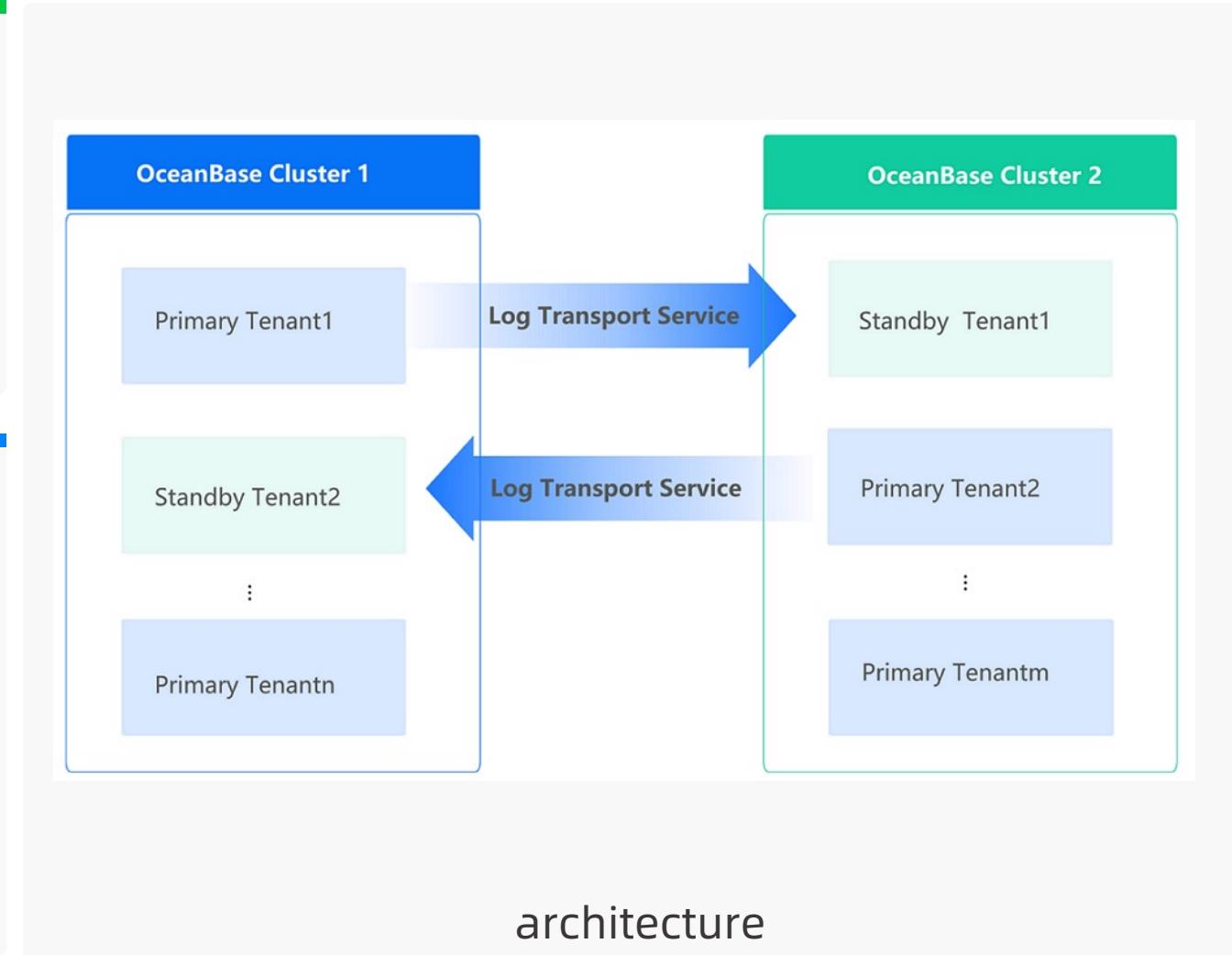
Similar like primary/slave cluster

User scenarios

- Single replica with high available
- Long network latency such as across city

Advantage

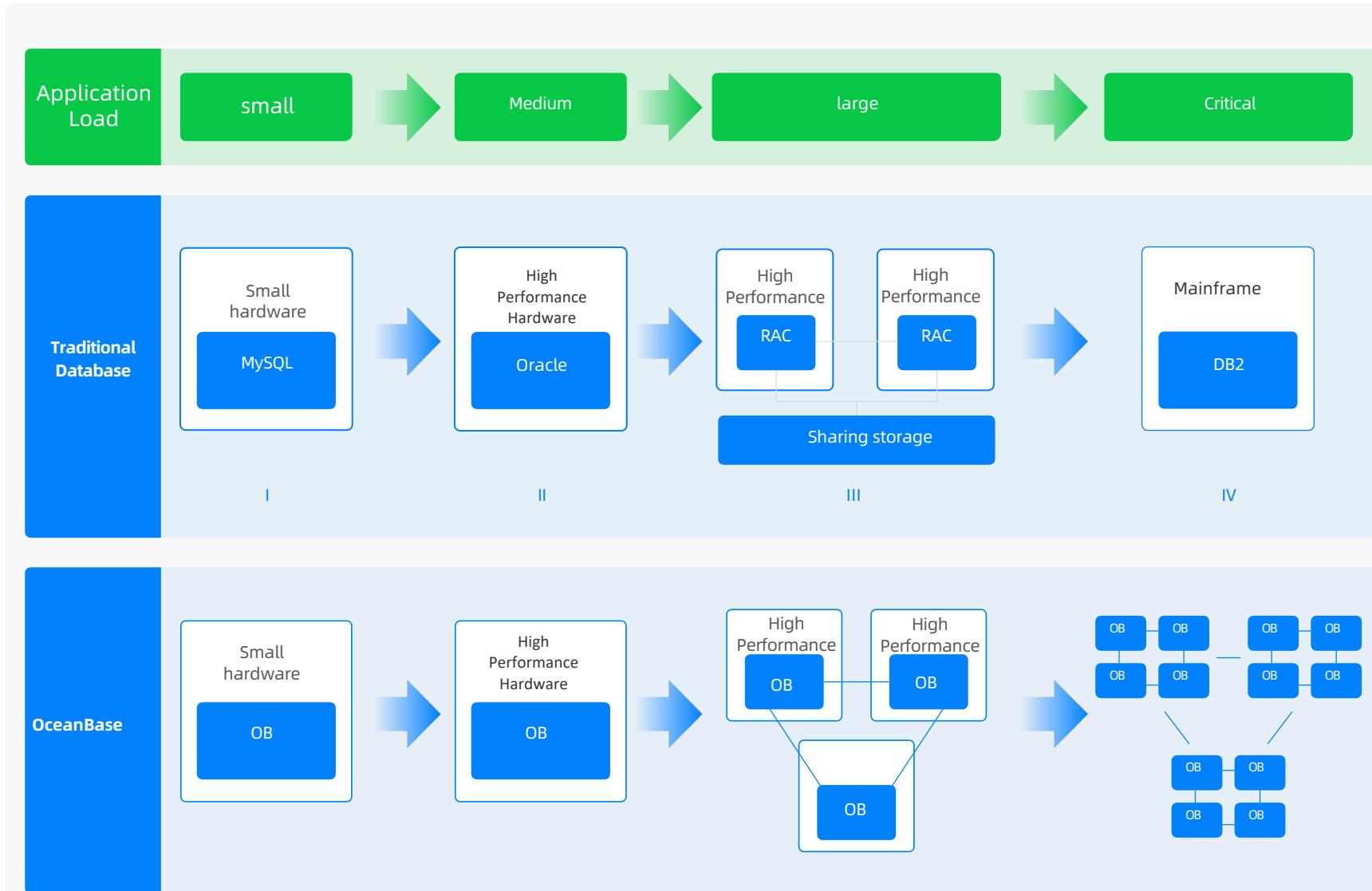
- Simple
 - Implementation by kernel, without any other component
 - The primary cluster hardware can be different from that of the slave cluster.
- More robust
 - Sync DDL automatically
 - Easily switchover/failover, support RTO < 8, RPO=0
- Performance better
 - Sync physical modification log(clog)
 - compress



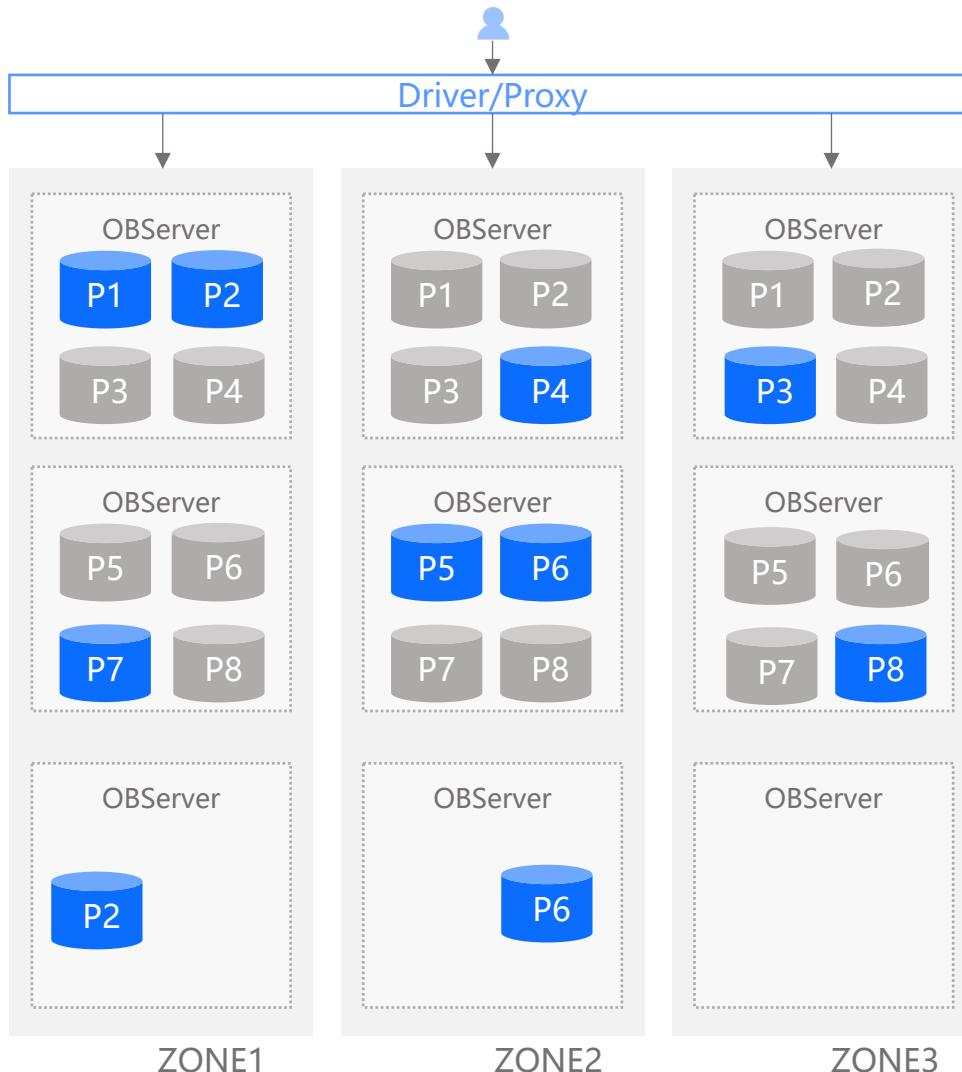
How to scale ?

Key Feature: Scalable

Flexible deployment - different size different architecture

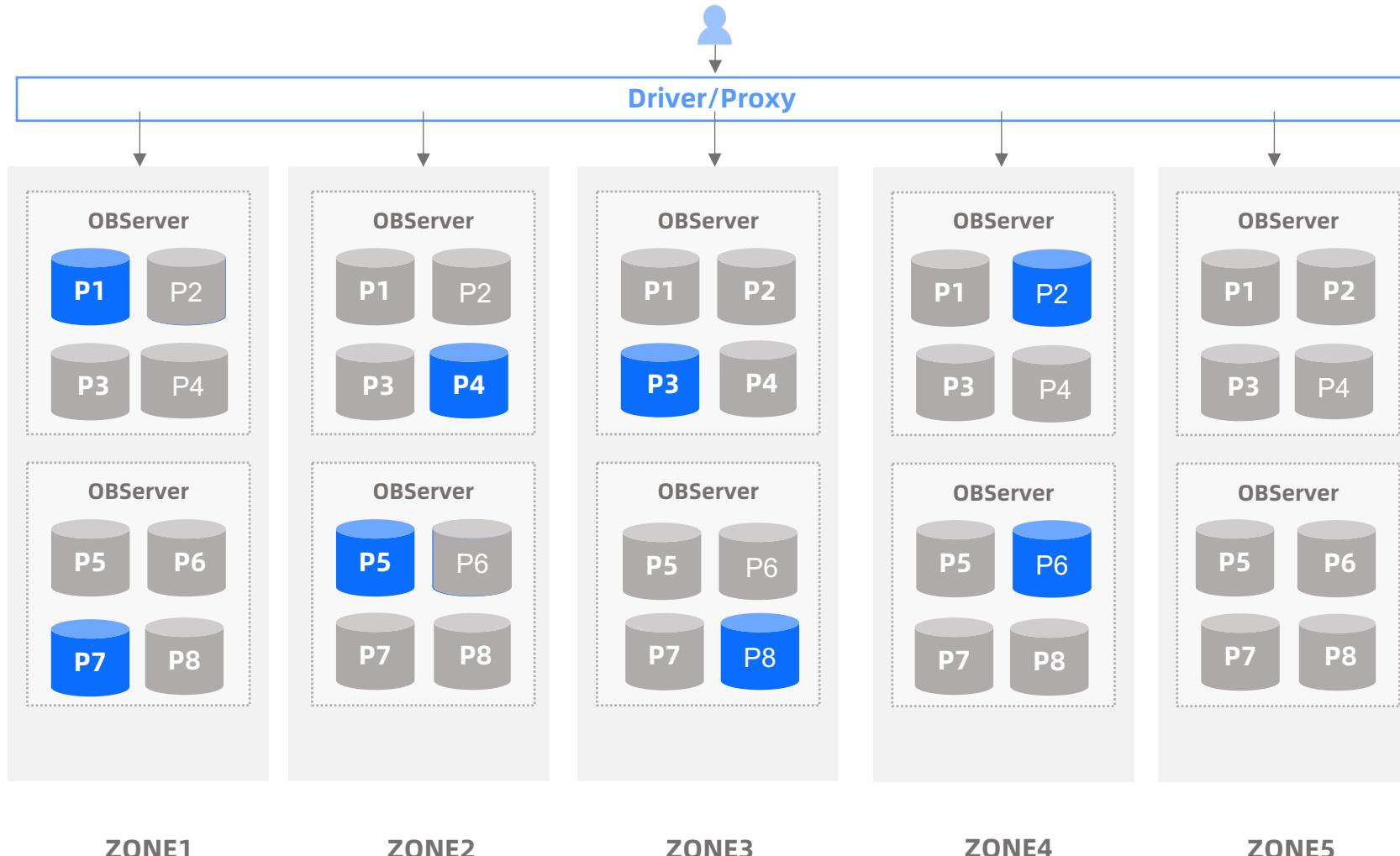


Transparent scalability for applications - vertical scaling



- Transparent for application
 - Add machine to cluster
 - Modify the number of tenant unit
- ✓ Auto/manual rebalance
- ✓ Sync speed under control
 - ✓ 70MB/s per thread, default 500MB/s

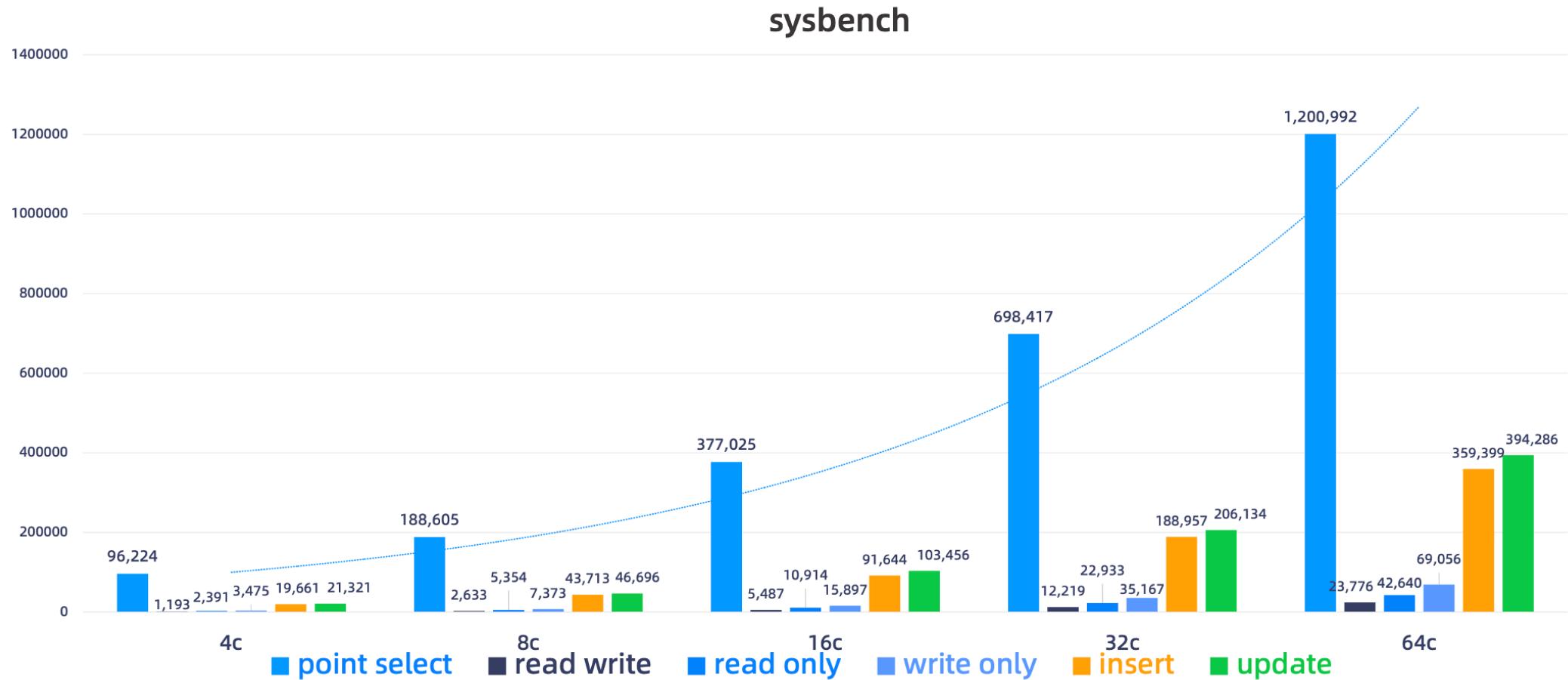
Transparent scalability for applications - horizontal scaling



- Transparent for application
- Add zone to cluster
- ✓ Auto/manual rebalance
 - ✓ Sync speed under control
 - ✓ 70MB/s per thread, default 500MB/s
- Auto select Leader
 - ✓ User can set zone priority

Leader
 Follower

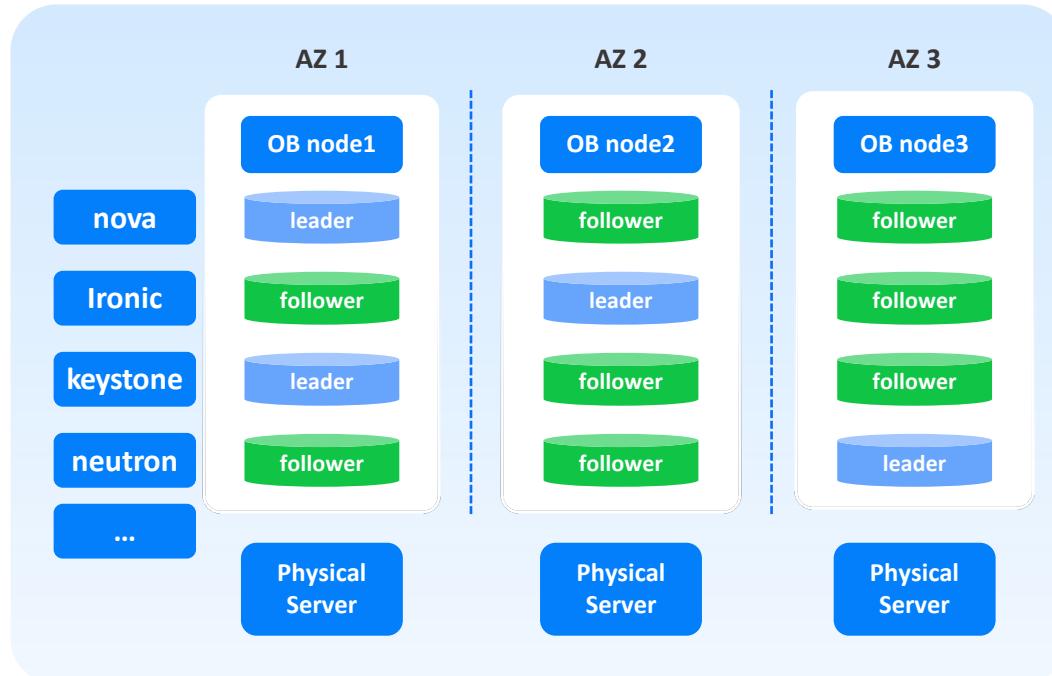
linear scalability



How to simplify maintenance task?

Key Feature: Easy Operation

Multi tenant – One Cluster for all user, but one tenant per module

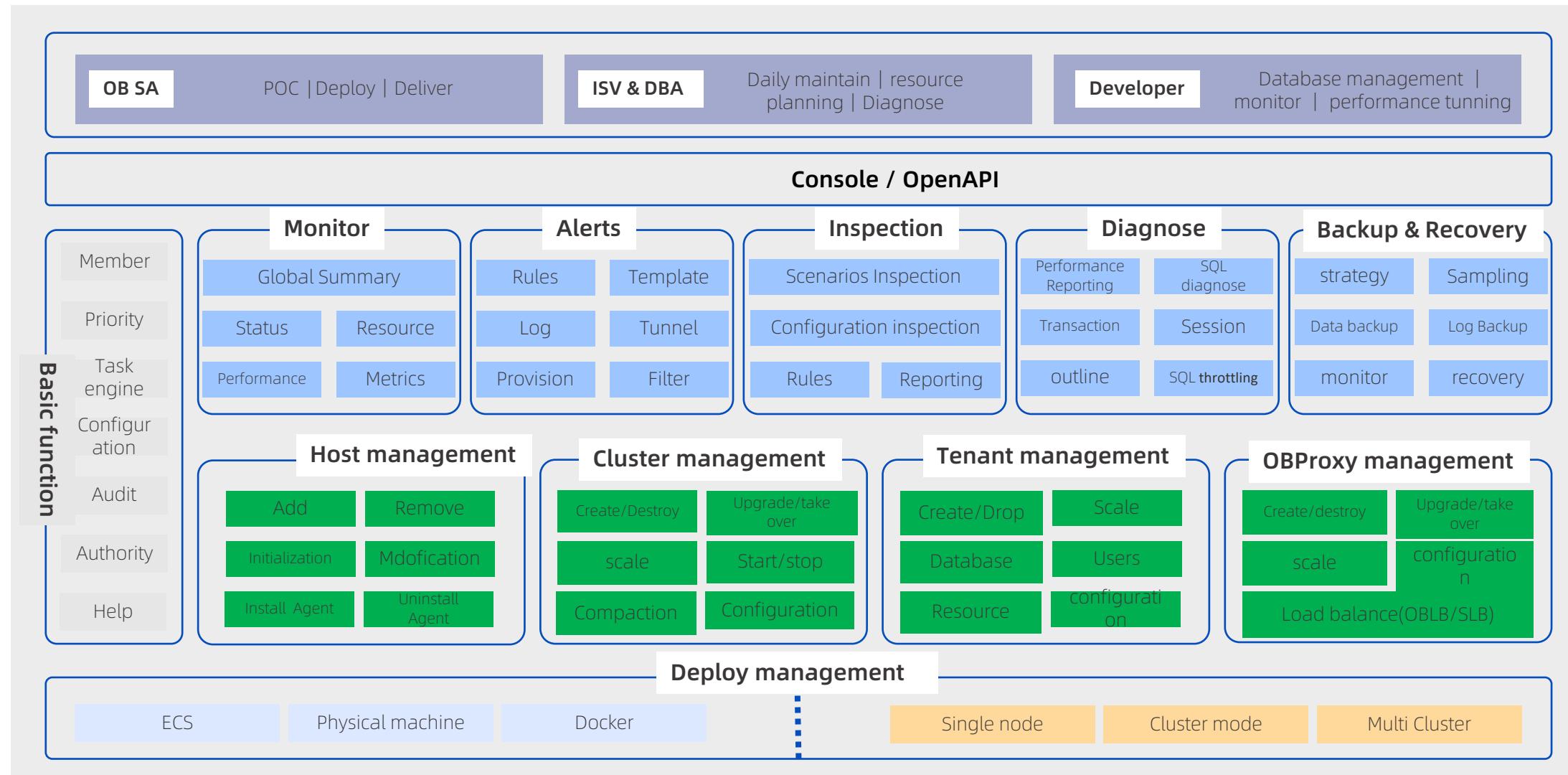


- Multi-tenancy (CPU/Mem/disk resource isolated)
- Scale up for a tenant in second

Zone	sqa.eu95_1	sqa.eu95_2	sqa.eu95_3
区租户名	副本大小		
agds_sit10_2085	7.7G		
asttship_1965	64.9G		
faasset_1966	45.0G		
fincore0_2431	1.3G		
finflux1_1964	142.1G		
gdtrans_1967	14.2G		
gzcomm1_2079	19.5G		
mapitool_1969	41.4G		
mix1_1968	12.4G		
mix2_1970	314.4G		
sys	0.0G		
trans_mix4_1973	2.7G		
trans_mix5_1974	5.3G		
内存分配	74.0-76.0G	58.0-60.0G	65.0-67.0G
总内存	80.4-321.8G	80.4-321.8G	80.5-321.9G
CPU分配	14-29.5/30-120	16.5-37.5/30-120	14-32.5/30-120
磁盘使用%	7.8%	11.1%	4.6%
磁盘空闲	3311.9G	3193.3G	3427.4G
总磁盘大小	3593.5G	3593.5G	3593.5G
unit数目	3/6	1/7	2/7
leader数目	21081	4472	12993
partition数量	34923	39837	37161
ob版本	1.4.79	1.4.79	1.4.79

One of application of Alipay
Tenant strategy is flexible

OceanBase Console Platform - all operation in one



test:3
ocp_meta
Tenant ID: 1010 (Running)

SQL Diagnosis

Features

View SQL Collection History View Outlines View Request Analysis

Duration: Last 30 Minutes Jan 8, 2024, 16:51:16 Jan 8, 2024, 17:21:16 OBServer: All OBservers Internal SQL?

Keyword: Please remove the literal quantity such as string and numerical value.

Conditions: + Add Reset Search Hide

TopSQL SlowSQL Suspected SQL ParallelSQL Export TopSQL Custom Column Column Management

SQL Text	Database	Tenant Name	Total Executions	Total Response Time (ms)	Response Time (ms)	CPU Percent	Actions
+ select outline_name as ...	oceanbase	ocp_meta	14	104.07	7.43	24.29	Throttling
+ select database_id, ten...	oceanbase	ocp_meta	14	71.87	5.13	15.77	Throttling
+ SHOW PARAMETERS W...	oceanbase	ocp_meta	2	44.01	22	12.17	Throttling
+ SELECT event_id, svr_i...	oceanbase	ocp_meta	1	37.19	37.19	13.95	Throttling
+ SELECT COUNT(*) FRO...	oceanbase	ocp_meta	1	18.94	18.94	6.32	Throttling
+ select database_id, ten...	oceanbase	ocp_meta	1	16.97	16.97	6.36	Throttling
+ set autocommit=1, sql_...	oceanbase	ocp_meta	3	11.53	3.84	4.24	Throttling
+ SELECT ? FROM DUAL	oceanbase	ocp_meta	58	11.3	0.19	2.4	Throttling
+ SELECT @@max_allowable...	oceanbase	ocp_meta	3	11.25	3.75	4.01	Throttling
+ SHOW GLOBAL VARIAB...	oceanbase	ocp_meta	2	10.89	5.44	4.03	Throttling
+ set ob_query_timeout =	oceanbase	ocp_meta	36	9.5	0.26	2.7	Throttling

03

Best Practices

User Scenario

Data Center (stored historical data)

Cost Effective/Simple management

Multi-Tenant

Cost Effective/suitable for SOA & SaaS

High availability

Anti disaster by deploy across IDC/Cities/Cloud

High Performance

Replace MySQL Sharding architecture

HTAP & OLAP

Realtime data warehouse

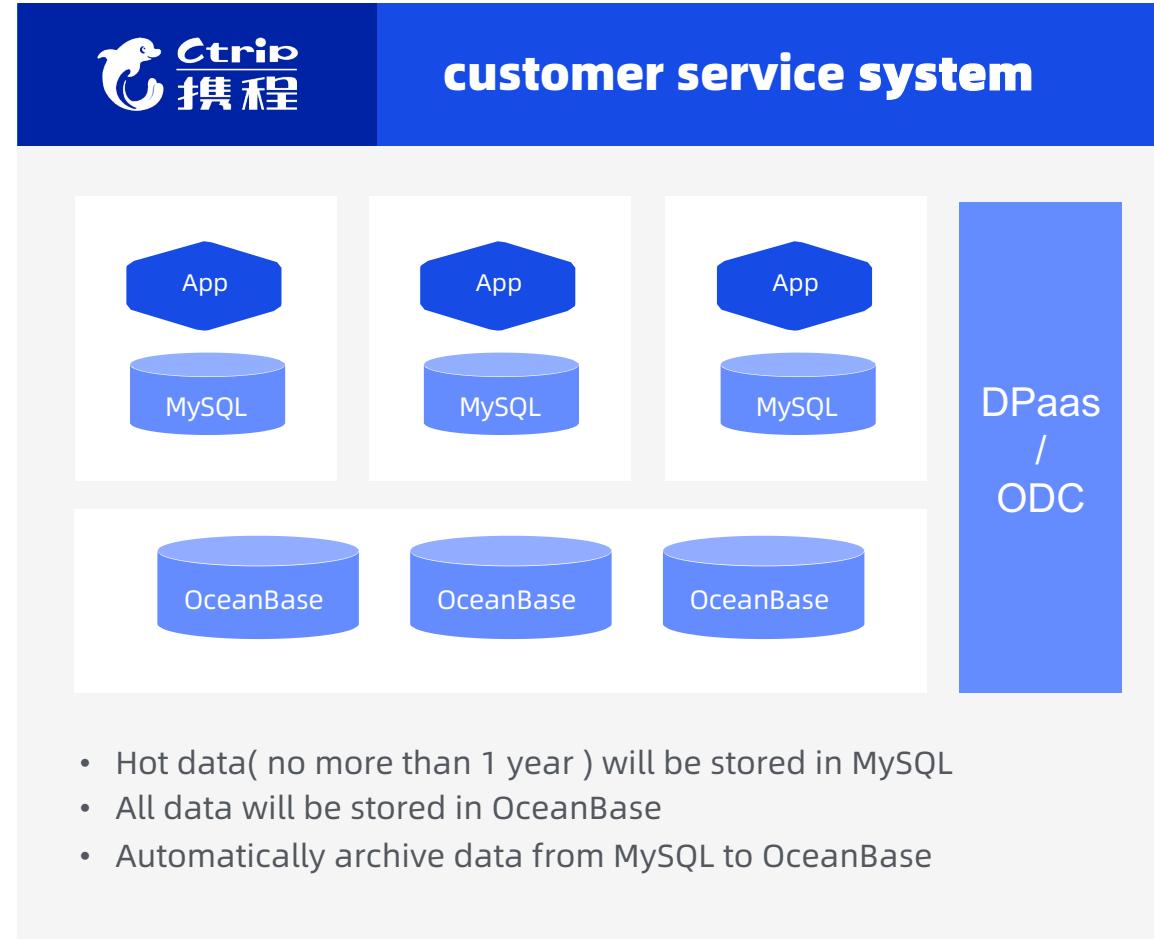
OBKV

Reuse all features of OB, cost effective

Typical user scenarios (1): data center (stored historical data)

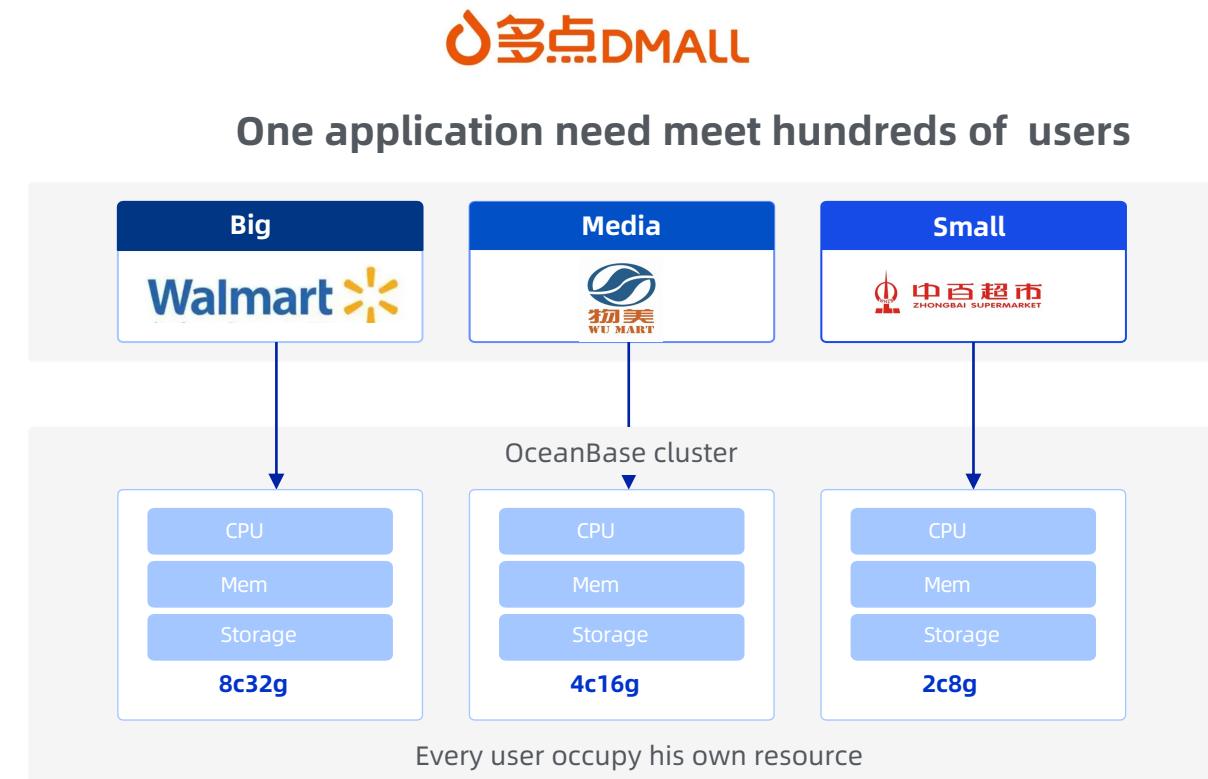
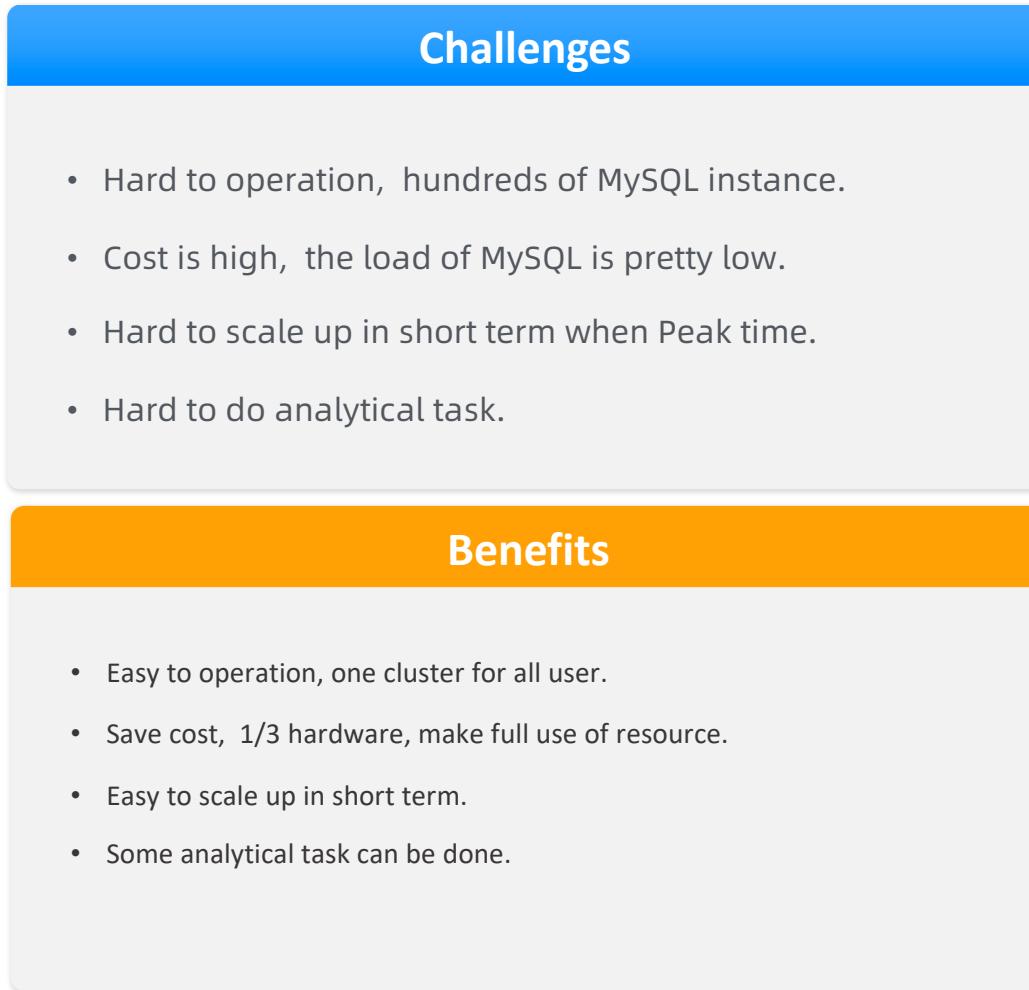
- Cost effective

Challenges
<ul style="list-style-type: none">• Cost is high - the data is huge.• Using sharding middleware<ul style="list-style-type: none">• hard to scale• hard to do complex query.
Benefits
<ul style="list-style-type: none">• Save Cost, the disk usage of OceanBase is only 15% of MySQL.• Easy to scale, no sharding middleware .• Improve writing performance 3x.



Typical user scenarios (2): Multiple Tenants

- Make full use of resource

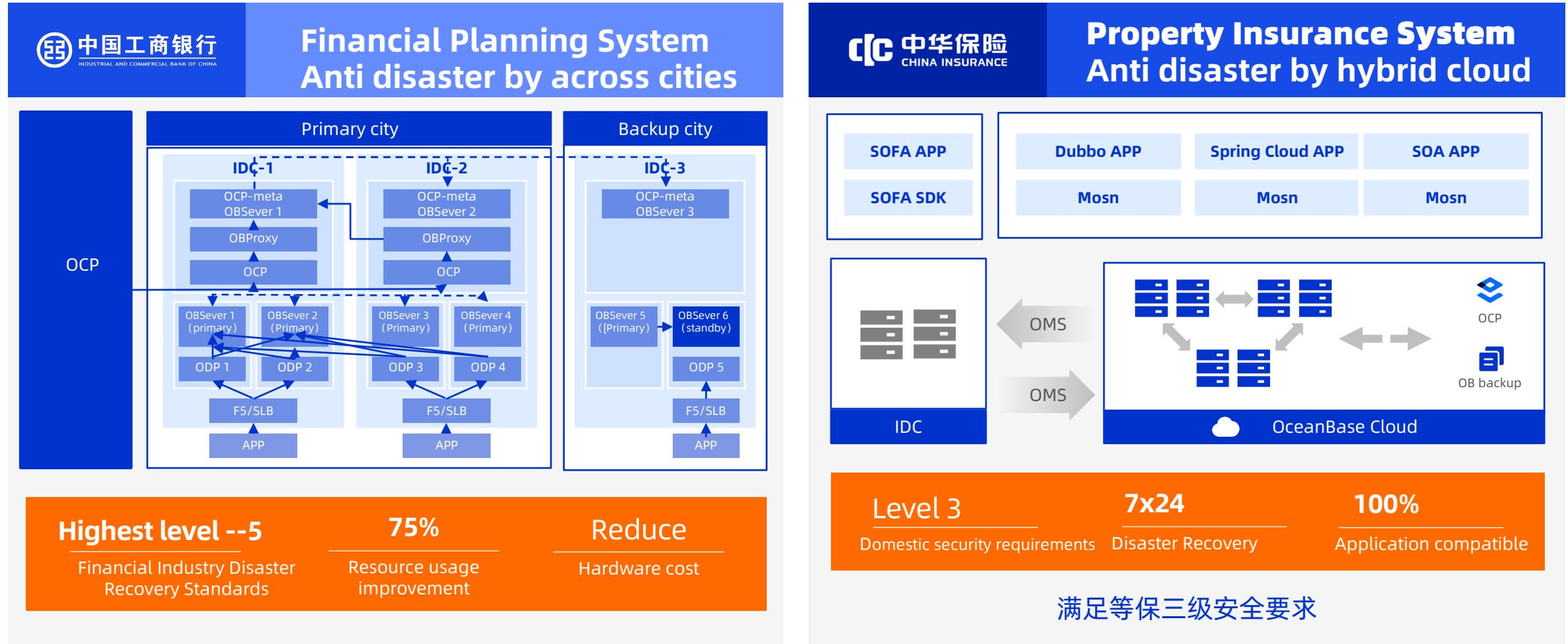


Retail SaaS

- Resource & data isolation
- Resource elastic scale in seconds

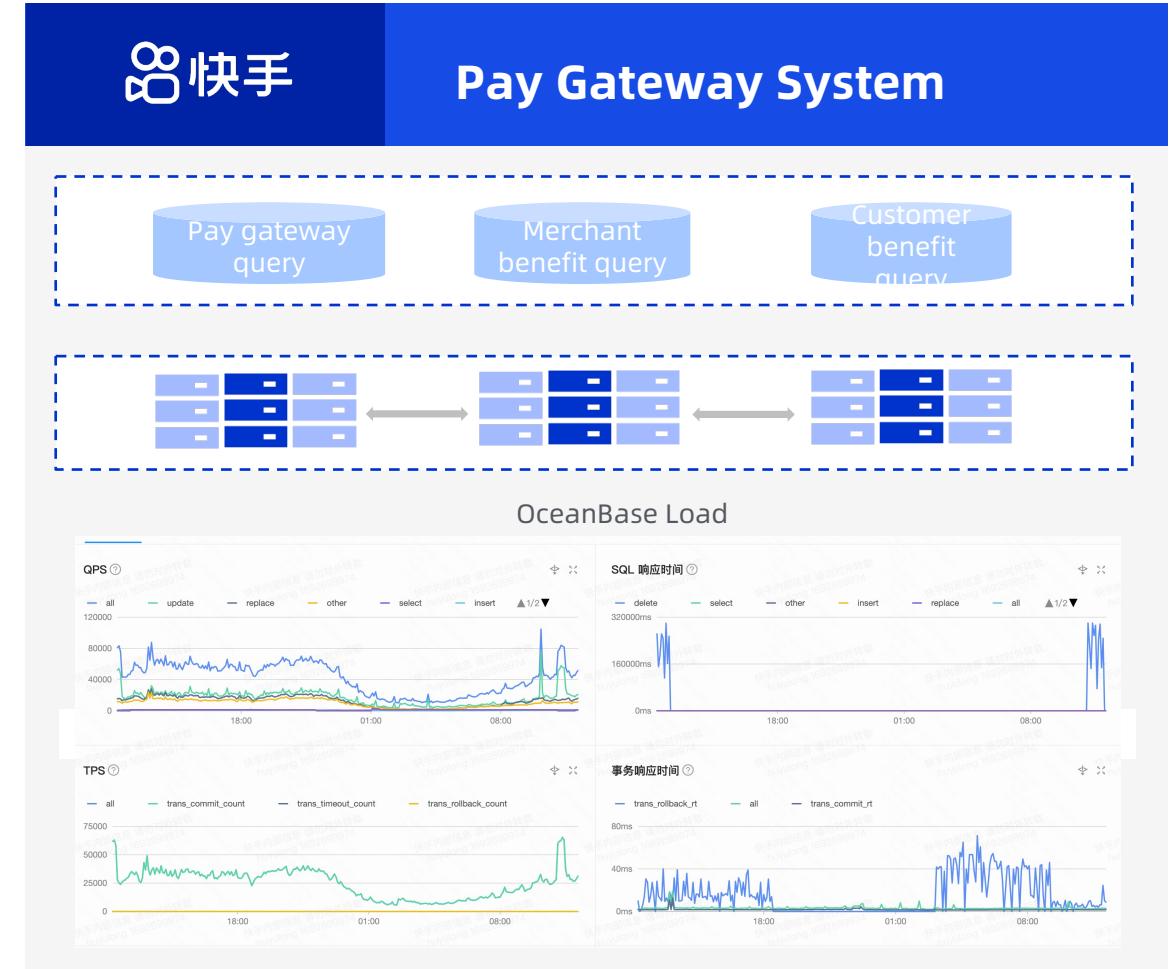
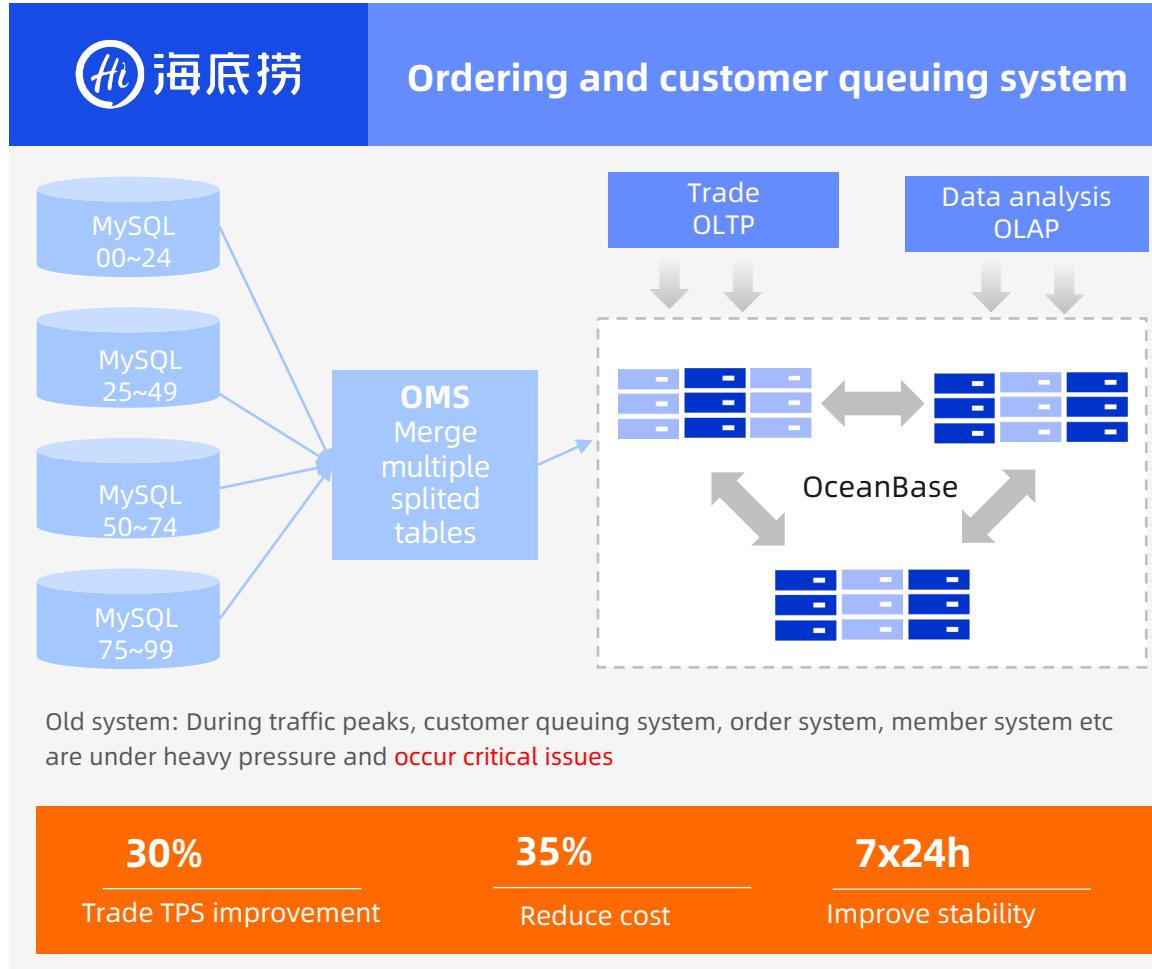
Typical user scenarios (3): High Availability

- 7 X 24 - Anti disaster



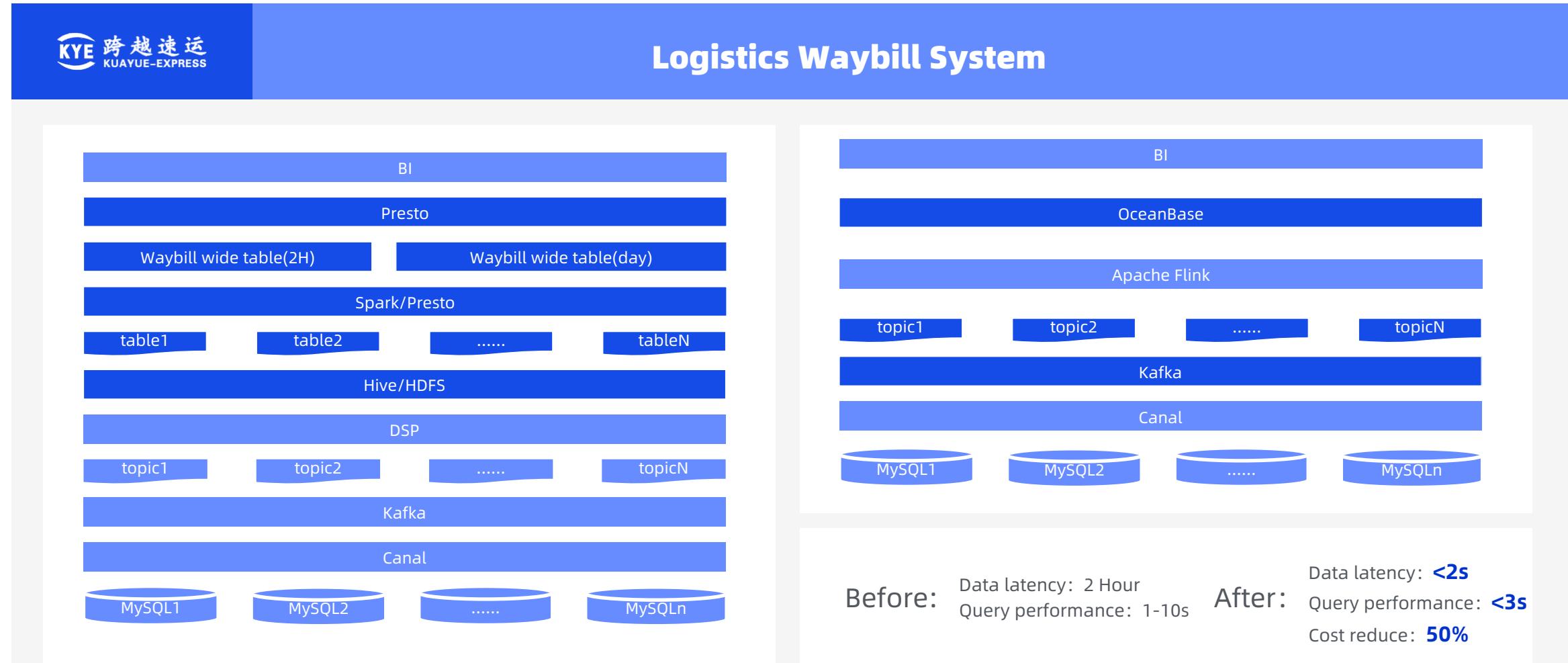
Typical user scenarios (4): High Performance

Replace MySQL Sharding architecture



Typical user scenarios (5): HTAP & OLAP

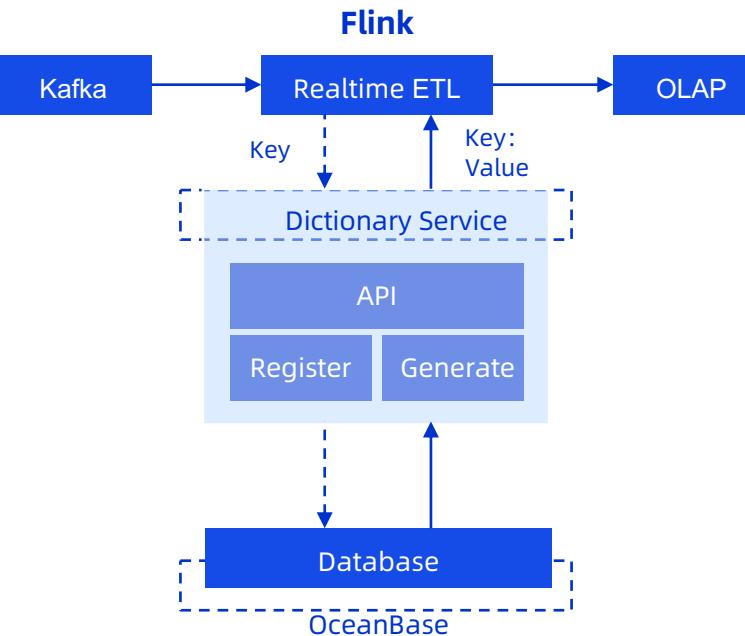
- Realtime data warehouse



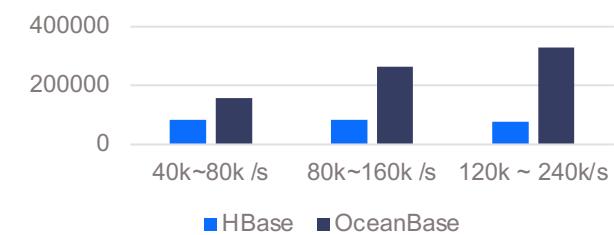
Typical user scenarios (6): OBKV



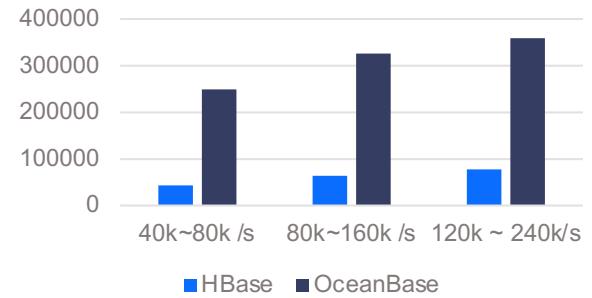
Realtime Dictionary Service



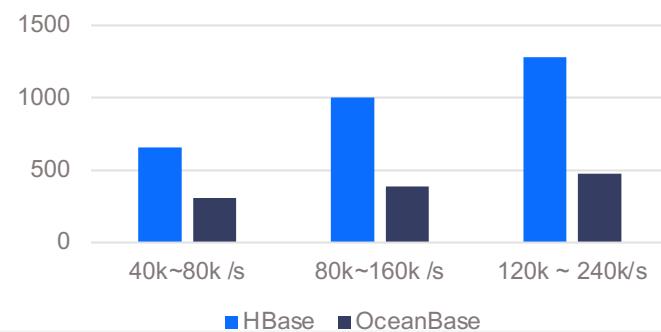
Query throughput
(row/s)



Write throughput(row/s)



Latency(ms)



Performance improvement

- Query: 4.3
- Batch insert: 4.6
- ETL : 2.7

Thank You!



OceanBase Official website: <https://oceanbase.github.io/>



Forum: <https://github.com/oceanbase/oceanbase/discussions>

