

数据来源：数据库产品上市商用时间



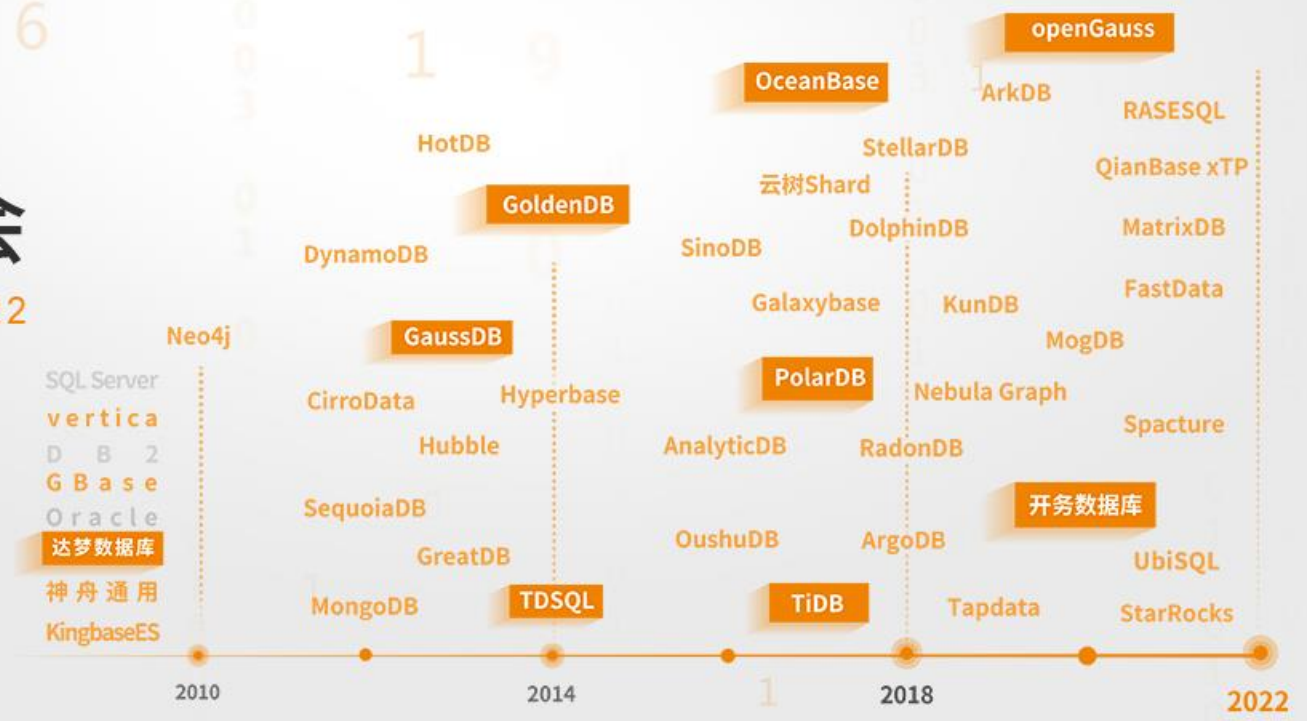
第十三届中国数据库技术大会

DATABASE TECHNOLOGY CONFERENCE CHINA 2022

数据智能 价值创新



线上直播 | 2022/12/14-16



融合普惠的云数仓 ——华为云GaussDB(DWS) 3.0

王传廷 华为云计算技术有限公司 技术专家

数据仓库趋势：下一个十年，智慧数仓提供开放、融合、云化、实时、全场景分析

描述型数仓

报表应用 (T+1)：固定查询

探索型数仓

灵活查询 (T+0.x)：分析师应用

运营型数仓

实时分析 (T+0)：实时运营/IOC

智慧型数仓

数智融合：BI + AI，数流、智流融合

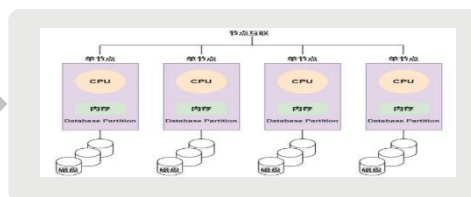
单机架构：TB



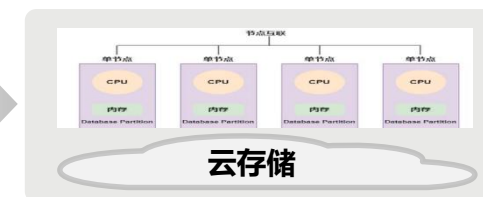
集中式架构：100TB



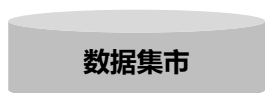
分布式架构：10PB



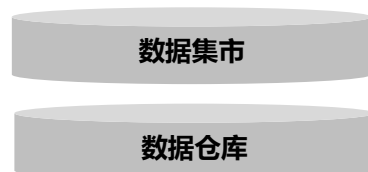
云原生架构：存算分离，EB级



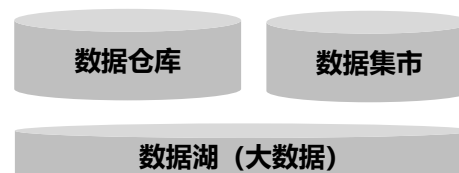
数据集市



数据仓库



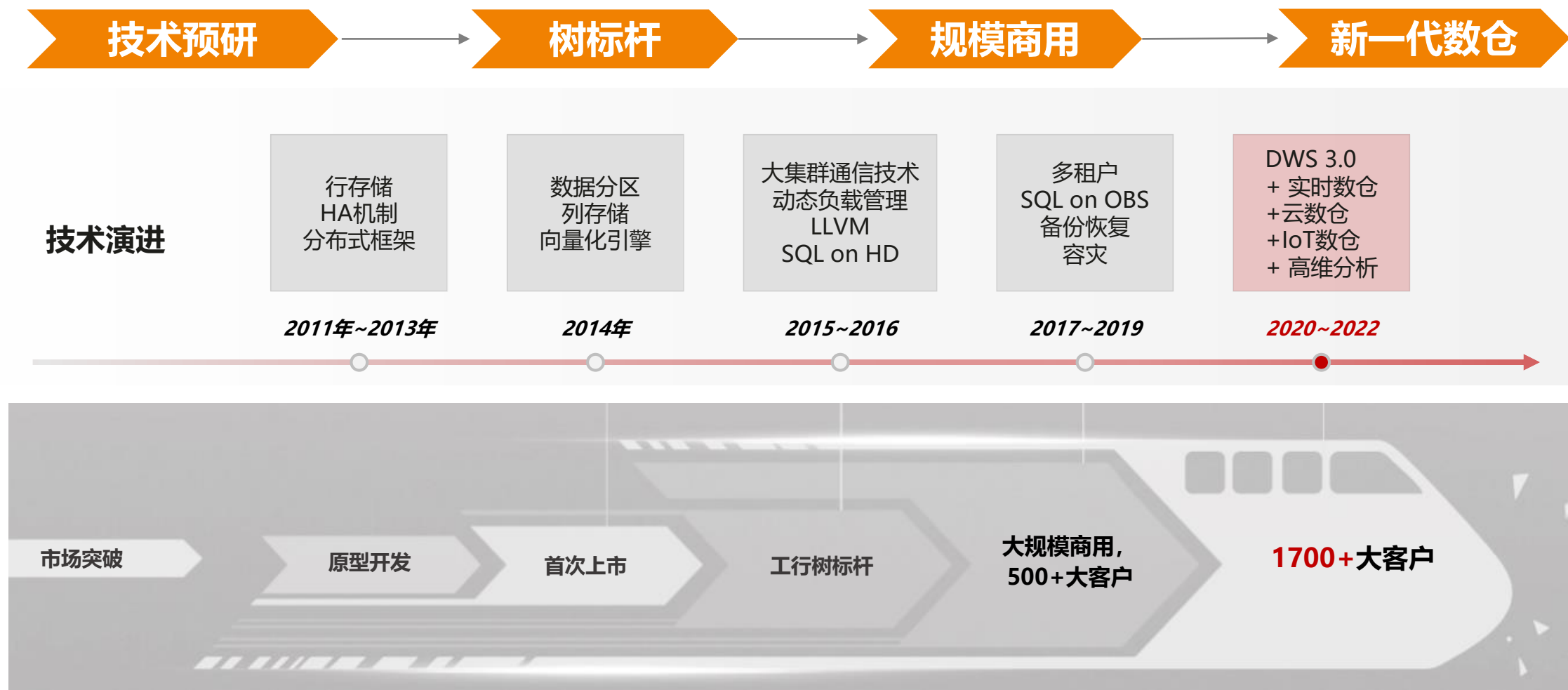
湖仓一体



数智平台



GaussDB(DWS)演进历程：十年技术沉淀，厚积薄发



GaussDB(DWS): 新一代全场景云数据仓库, 简化IT架构使能人人用数



全场景一站式

- 一套内核、一套架构支撑公有云、混合云、On-Premise部署, 用户体验一致
- 支持虚拟机(ECS)、裸金属(BMS)、物理机(HCS)模式
- 支持批量分析、实时分析、交互式查询、HTAP, **实时数据和历史数据关联分析**

融合分析

- 湖仓一体**: +HD Connector, +ORC, +OBS, 与大数据互联互通
- 数智融合**: +AI Connector, 结构化和非结构化数据关联分析
- 高维分析**: 时空分析 (时序/+GIS), 特征分析 (Text Search), 关系分析 (+GES), 打通关系型数据和非关系型数据分析边界

核心技术

高性能

万亿数据分析秒级响应

高扩展:

2048节点, >100PB

高可用

强一致性; 集群内 RPO=0, RTO<30s;
在线扩容业务零中断

高安全:

CC EAL2 + ALC_FLR.2

Note: 标准数仓: T+1, OLAP | 实时数仓: T+0.x, HTAP | IoT数仓: T+0, 时序计算

云数仓技术特点：分层弹性、横向融合、数智融合



分层弹性

- ✓ 存算管的三层分离
- ✓ 计算存储独立伸缩 -> Serverless化
- ✓ 数据共享



横向融合

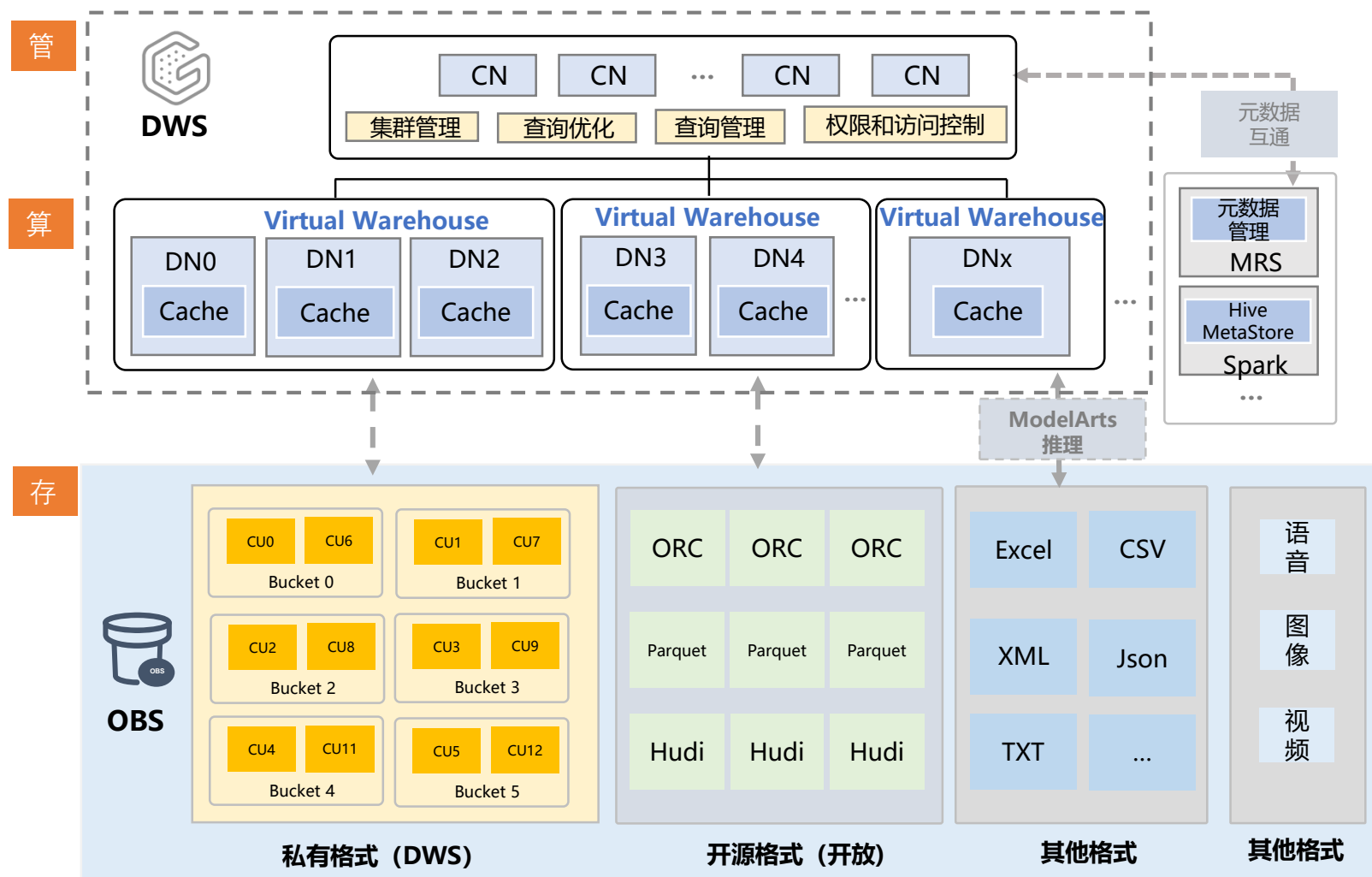
- ✓ 多编程语言支持
- ✓ 湖仓一体、数智融合，
打造端到端一站式的
服务



数智融合

- ✓ 面向系统内部，智能化降低用户维护成本和使用门槛；精细化的资源调度，提升性价比
- ✓ 面向用户，提供良好的AI编程能力

GaussDB(DWS)3.0的Serverless云原生架构，极致弹性，湖仓一体，数智融合



Serverless的云原生架构

- ✓ 存算管分离，分层独立弹性
- ✓ 吞吐线性提升
- ✓ 资源隔离

极致弹性

- ✓ 多种形式弹性
- ✓ 数据共享
- ✓ 业务负载隔离，承载能力线性扩展

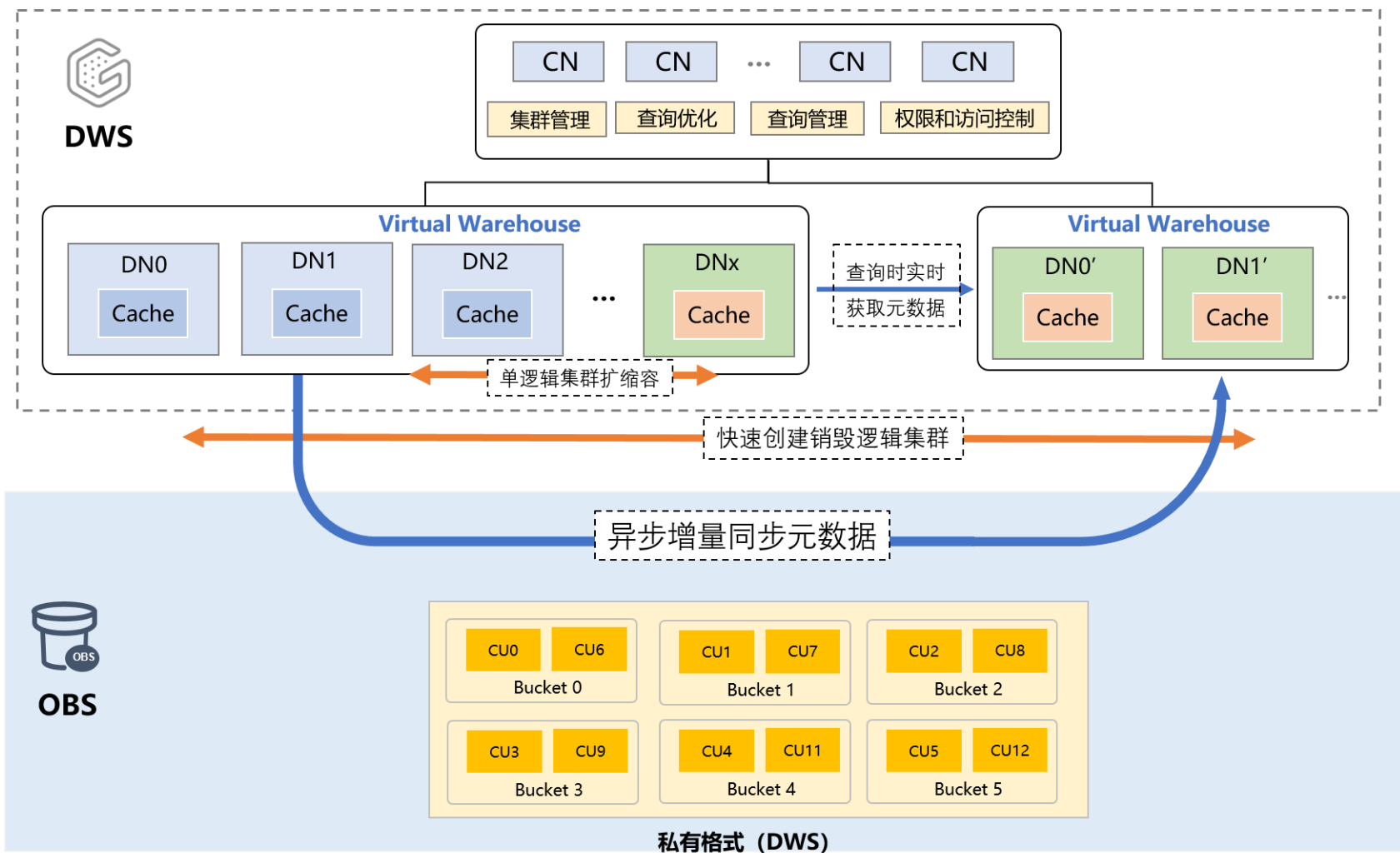
湖仓一体

- ✓ 在数据湖上体验数仓性能和管控度
- ✓ 纵向加速

数智融合

- ✓ 数据生产线与AI生产线的无缝对接

弹性优势：极致弹性、数据共享，赋能高灵活度、高性价比的使用体验



三层解耦

- ✓ 存算管三层分离，独立伸缩

灵活弹性

- ✓ 分钟级单逻辑集群扩缩容
- ✓ 分钟级快速创建销毁逻辑集群
- ✓ 快速扩缩容，无数据重分布、拷贝

一数多用

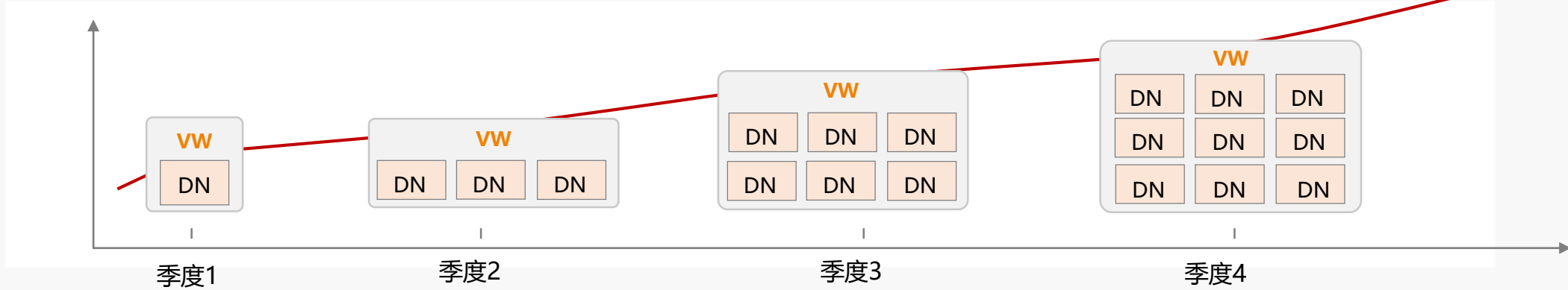
- ✓ 多逻辑集群间共享数据，无需拷贝
- ✓ 提供实时和近实时两种数据共享方式

按需配置

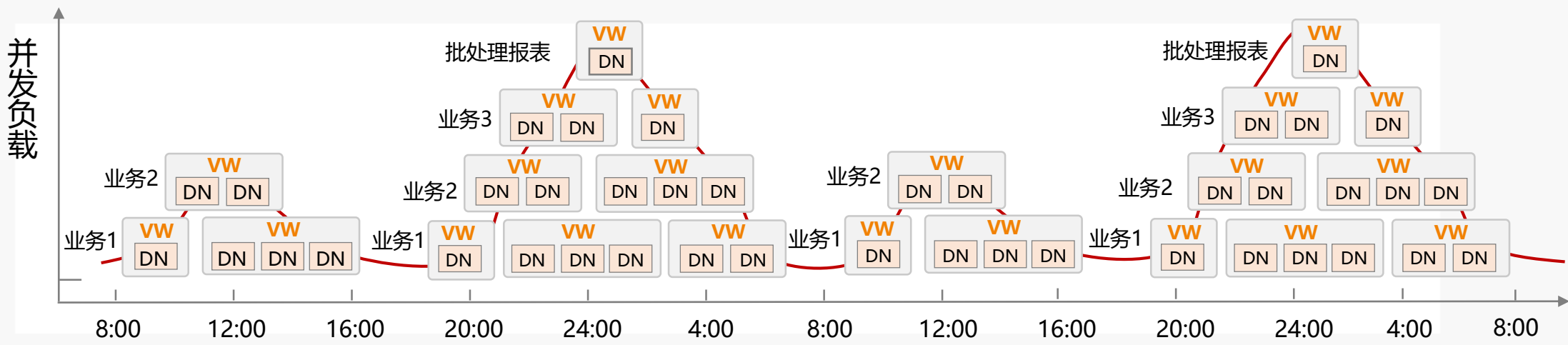
- ✓ 逻辑集群隔离不同业务
- ✓ 业务承载量/并发量的线性扩展
- ✓ 读写分离、不同负载隔离

按需弹性实践：高性价比地适应灵活多变的业务需求

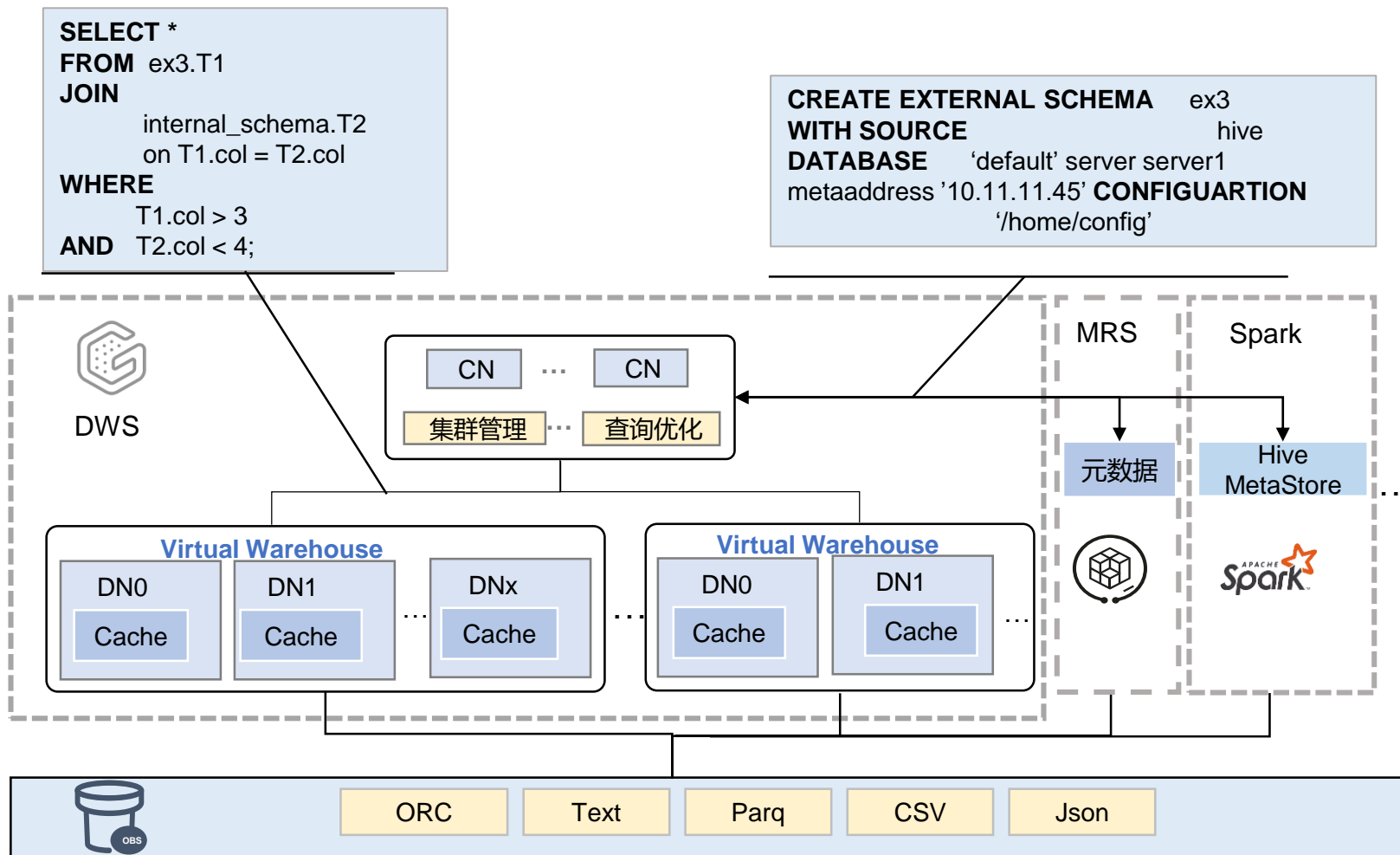
长期时间推演



短期单日弹性



湖仓一体：与大数据互联互通，体验横向融合分析



无缝访问数据湖

- ✓ 对接Hive Metastore元数据管理，直接访问数据湖的数据表定义
- ✓ 支持主要数据格式：ORC, Parquet, Hudi, Carbon

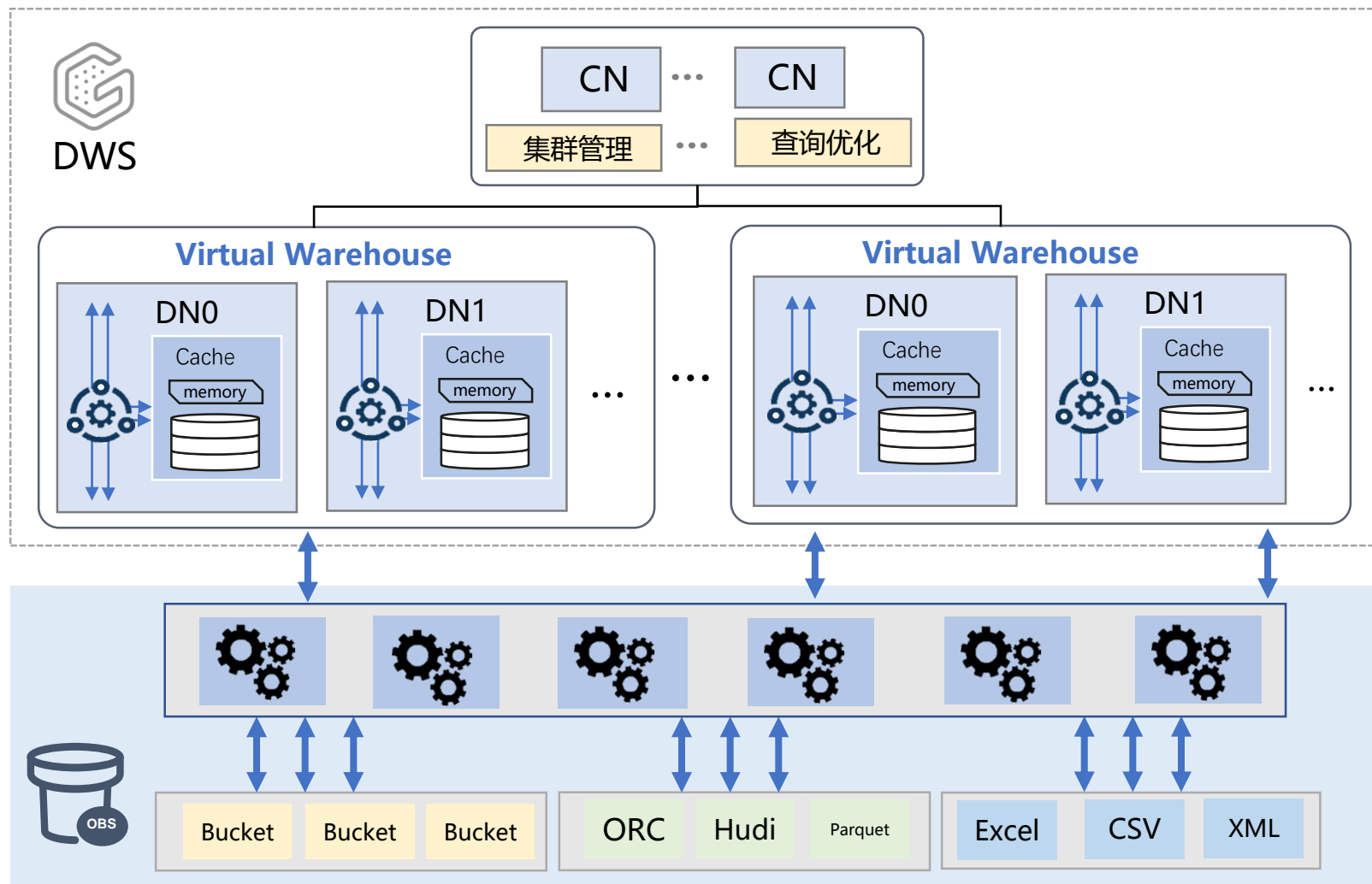
融合查询

- ✓ 混合查询数据湖和仓内的任意数据
- ✓ 一步到位，查询输出到仓内/数据湖

极致查询性能

- ✓ 使用数仓高质量的查询计划和高效的执行引擎
- ✓ 使用数仓的负载管理手段，精准控制

纵向加速：灵活可配的性能优化选择，保持优异性能



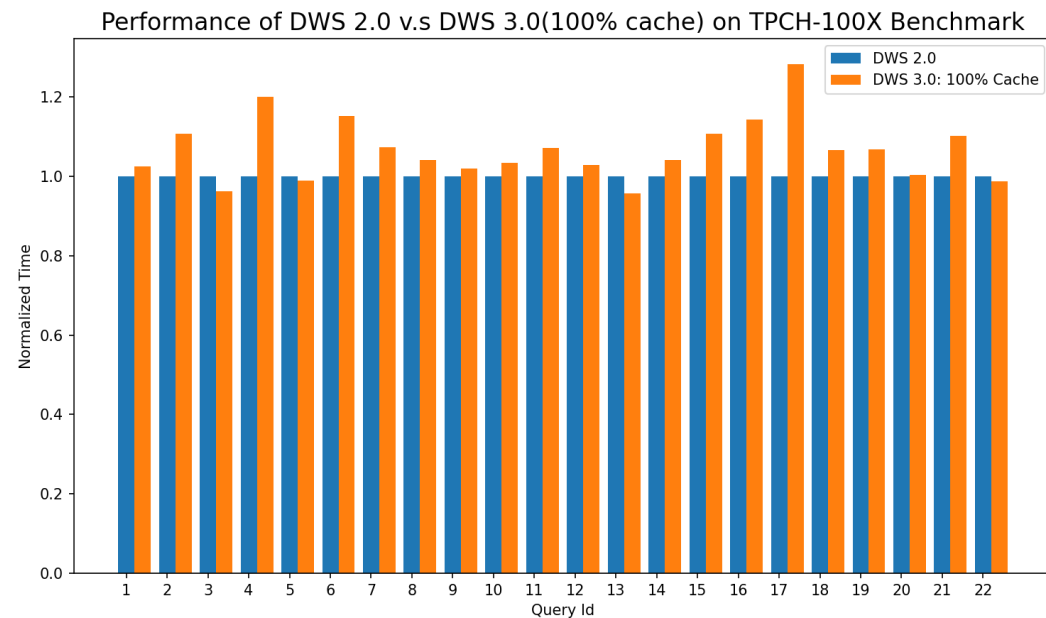
一体化性能优化

冷热分区高效缓存

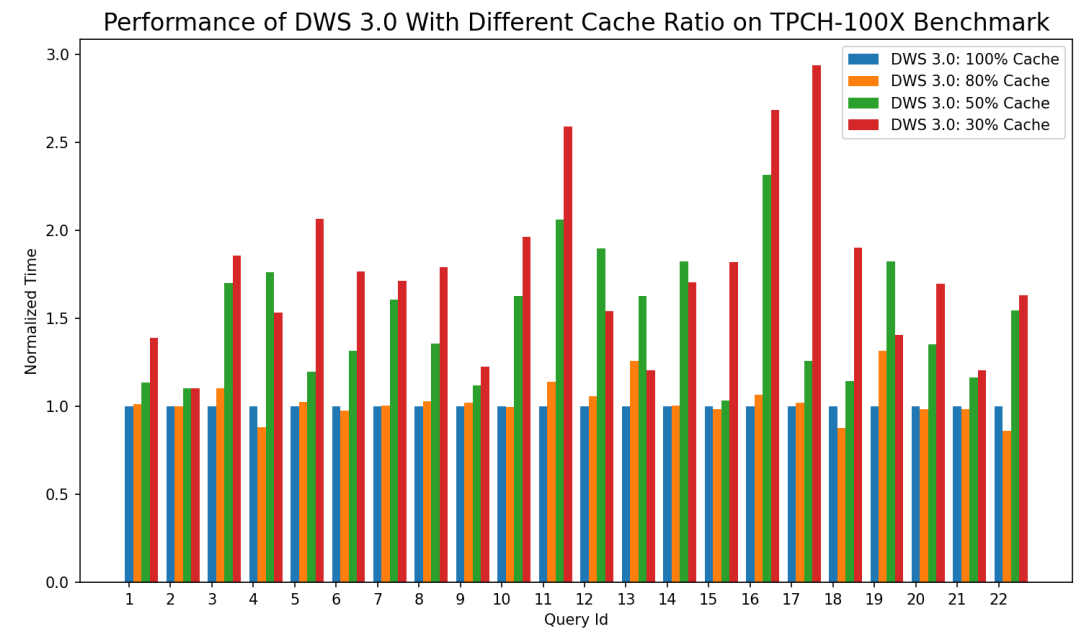
近数据计算

大带宽云存储

灵活可配的性能优化选择，保持优异性能

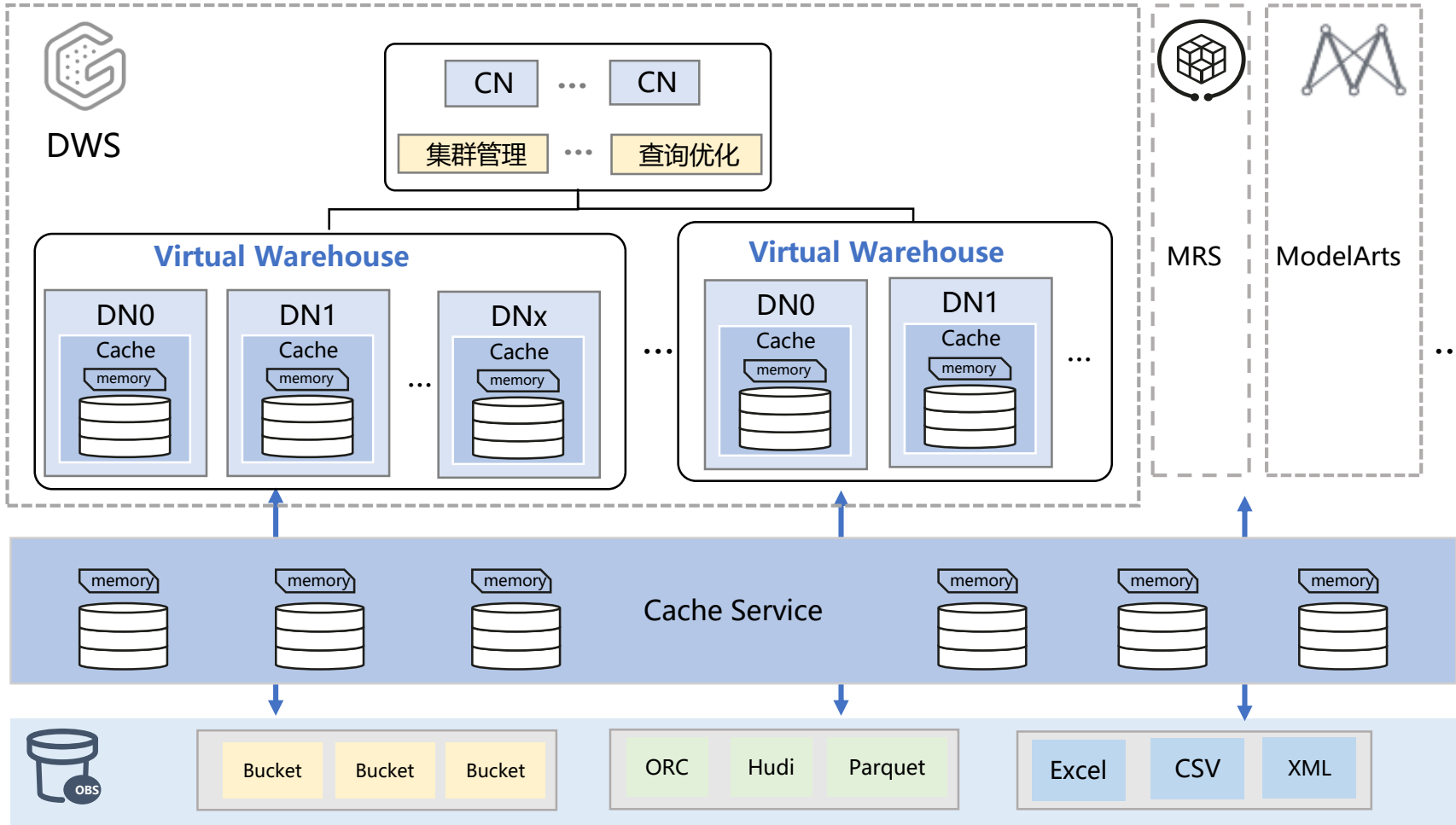


DWS 3.0在100%磁盘缓存情况下与DWS 2.0本地盘性能基本打平



DWS 3.0的磁盘缓存有效保证查询性能，查询性能随缓存配比递增

灵活可配的缓存，提供性价的按需权衡



多级缓存

- ✓ 内存、磁盘的分级缓存
- ✓ 基于MPP架构，充分利用并行算力和高扩展性

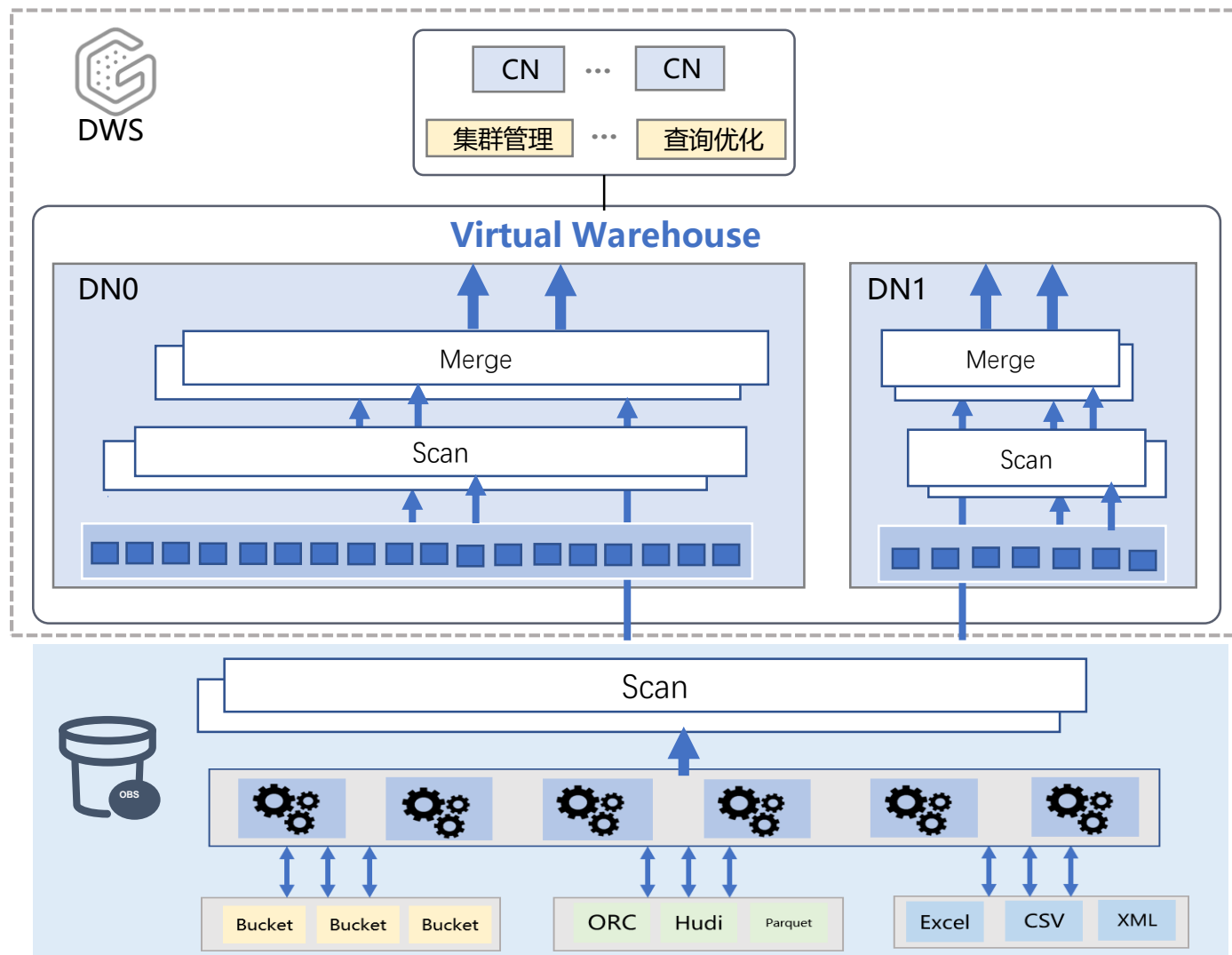
灵活的缓存策略

- ✓ 缓存大小可配置
- ✓ 表级分区配置
- ✓ 数据访问冷热程度可配置

跨计算引擎缓存

- ✓ 单VW缓存
- ✓ 跨VW缓存

运用云存储近数据计算能力优化网络读取



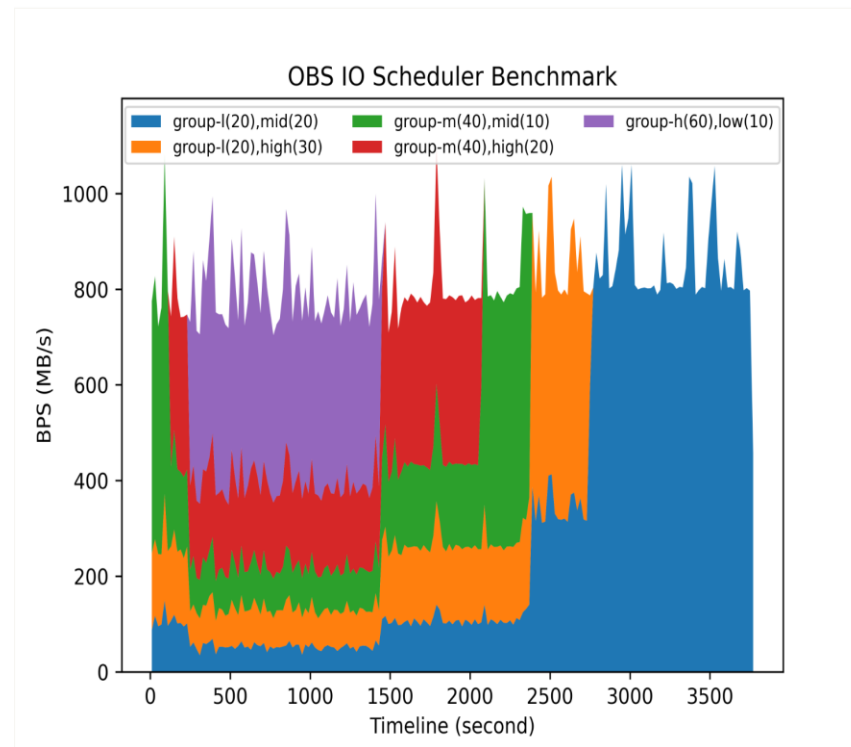
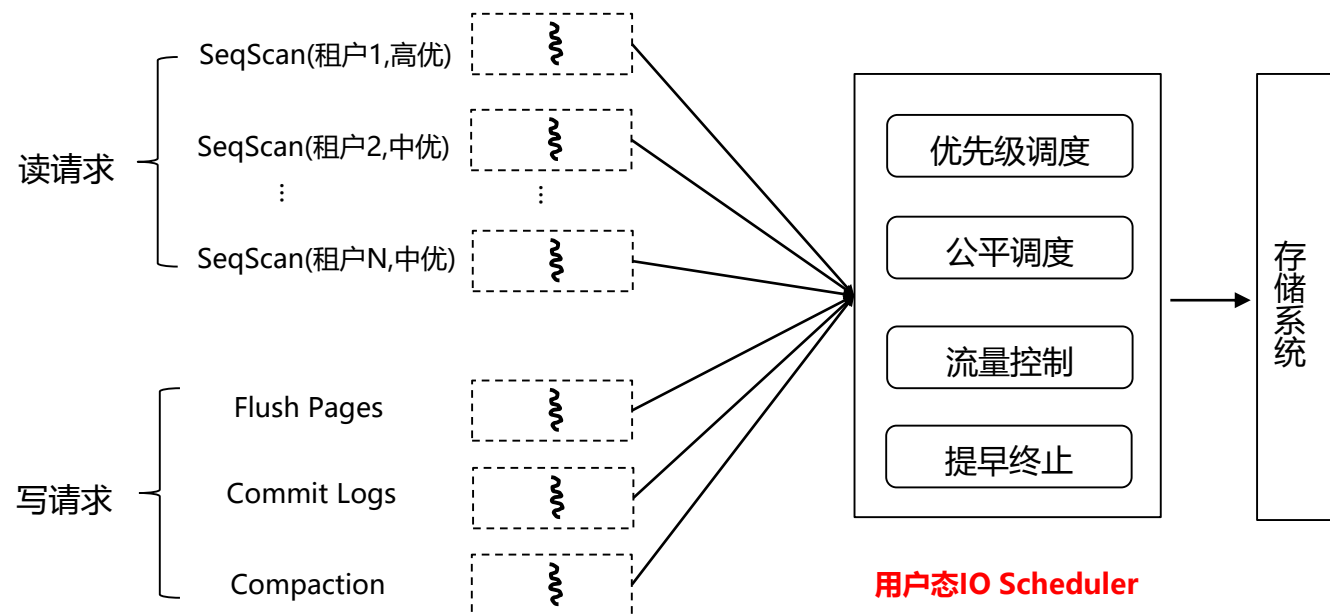
无缝配合缓存:

- ✓ 热数据优先缓存，使用本地的算子下推能力
- ✓ 冷数据优先下推，使用云存储的近数据计算资源池

近数据计算:

- ✓ 将计算下推到云存储，显著降低数据读取量

充分利用和精准控制云存储能力，深度优化存算分离架构



更低时延

- ✓ 充分利用云存储的带宽优势，弥补其相较传统MPP的高延迟劣势

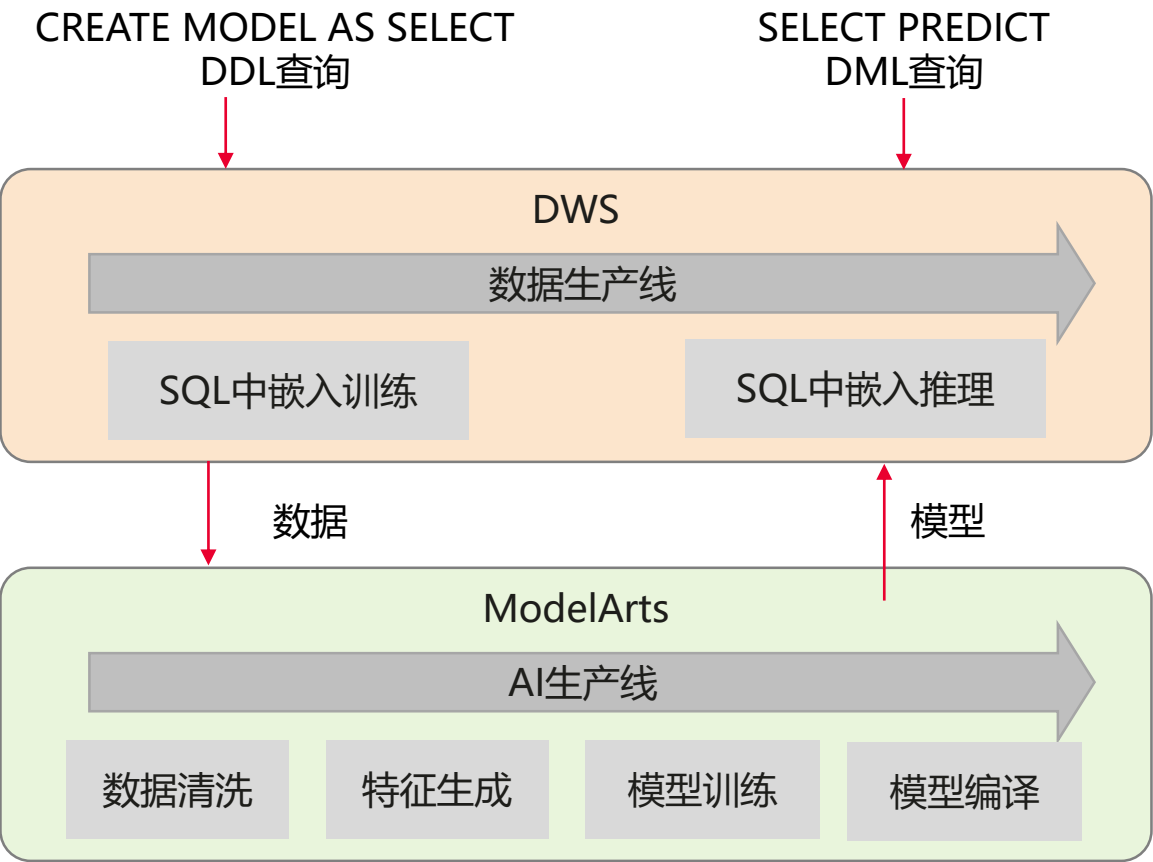
更优资源调度

- ✓ 单查询充分利用资源，为并发查询提供稳定、可预测的性能保证

更灵活配置

- ✓ 多级资源池灵活配置

数据生产线与AI生产线的高效配合



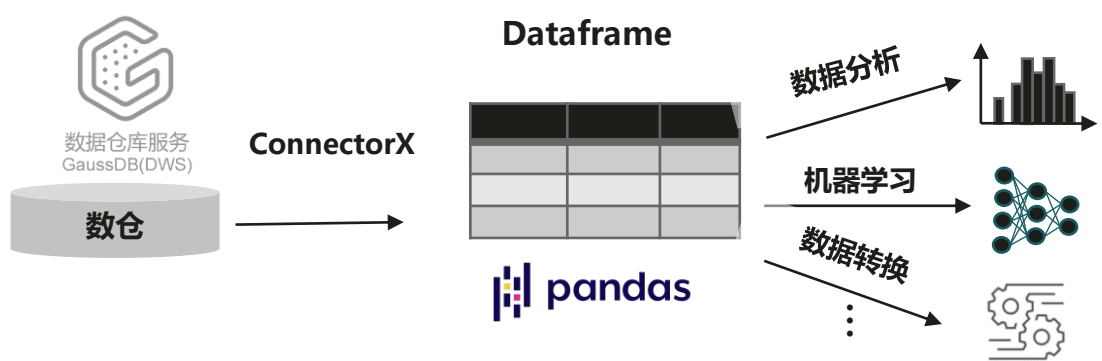
数据生产线→AI生产线：无缝数据通路

- ✓ **面向批量生产**：通过OBS共享开放格式数据
- ✓ **面向快速开发**：通过ConnectorX等以查询取数的方式嵌入Python开发生态，重点是Pandas

AI生产线→数据生产线：AI for Data

- ✓ 提供SQL语法，在数据分析过程中提供**驱动AI训练、应用AI推理**的能力
- ✓ **将推理能力引入分析**：直接调用部署的推理服务端点，灵活性好；将模型二进制部署为UDF，性能好

数智融合，高效沟通数据与AI两生产线

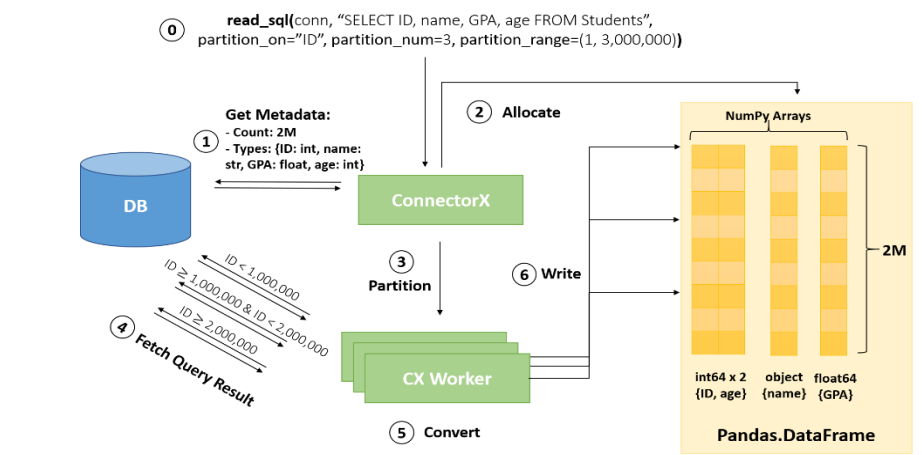


```
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.pipeline import make_pipeline
from sklearn.impute import SimpleImputer
from sklearn.feature_selection import SelectKBest
from sklearn.datasets import fetch_openml
from sklearn.linear_model import LogisticRegression
import connectorx as cx

# Read data from DB
df = cx.read_sql("postgresql://username:password@server:port/database", "SELECT * FROM titanic")

# Data Preprocessing
X = df.drop(["survived"], axis=1)
y = df["survived"]
numeric_features = ["age", "fare"]
numeric_transformer = make_pipeline(SimpleImputer(strategy="median"), StandardScaler())
categorical_features = ["embarked", "pclass"]
preprocessor = ColumnTransformer(
    [
        ("num", numeric_transformer, numeric_features),
        ("cat",
         OneHotEncoder(handle_unknown="ignore", sparse=False),
         categorical_features),
    ],
    verbose_feature_names_out=False,
)

# Building ML models
log_reg = make_pipeline(preprocessor, SelectKBest(k=7), LogisticRegression())
log_reg.fit(X, y)
```

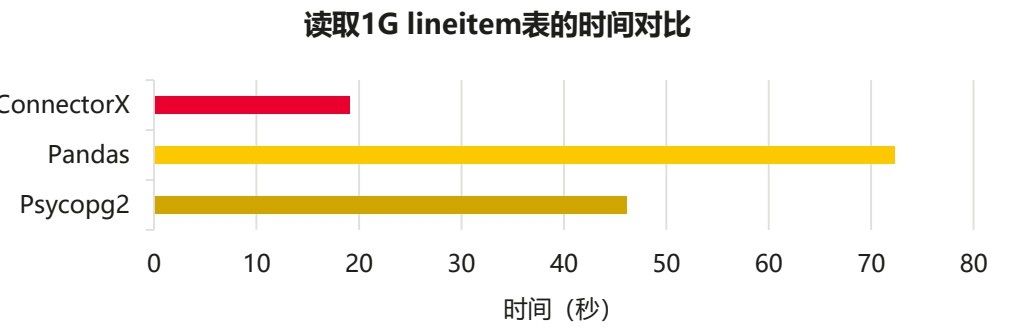


流式处理 + 优化技术：加速数据流转

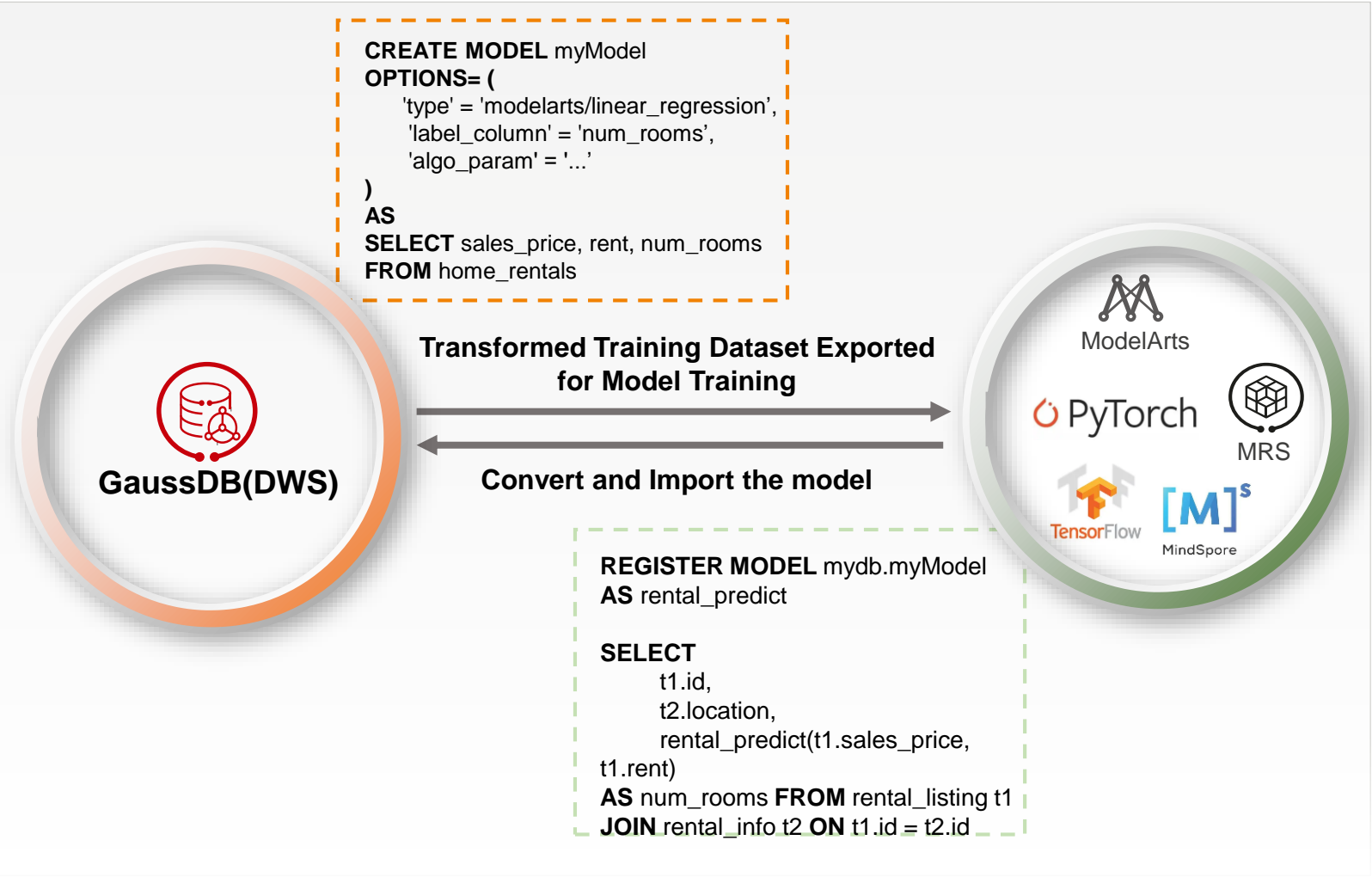
并行执行

空间预分配

Python字符串分配优化



数智融合，赋能用户BYOM体验



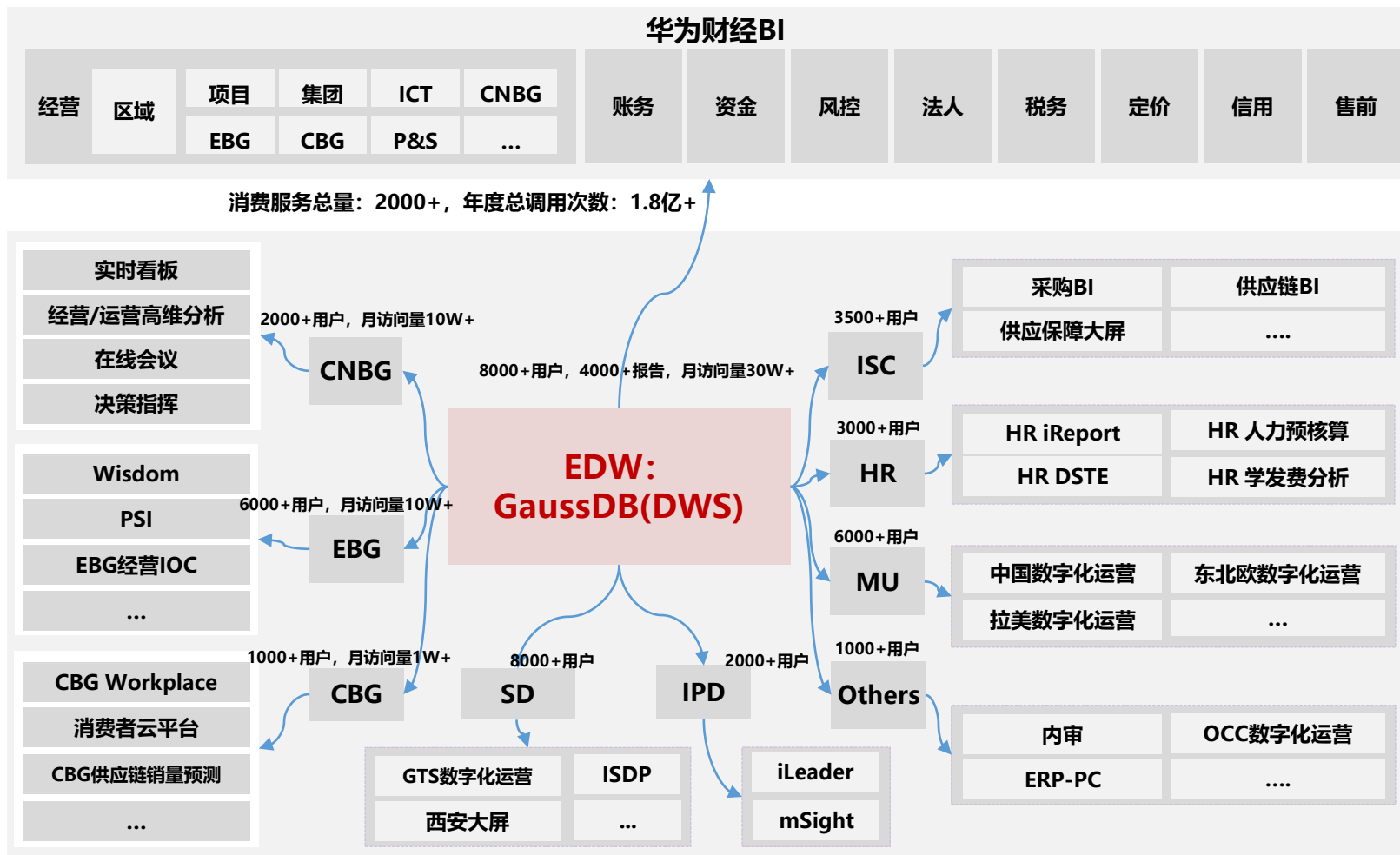
BYOM体验

Data for AI

原生SQL语句

模型推理

企业级数据底座上云的最佳实践：华为集团财经EDW



业务挑战

数据耦合度高
18个业务领域，相关性强，数据共用

多源异构场景
交易系统400+
Oracle实例1400+

复杂计算资源消耗巨大
5年以上历史数据，海量计算模型

业务规则变化
2万多条动态规则

业务效果

- 月度结账周期从15天缩短到**3天**
- 年度报告周期从月级缩短到**10天**
- 长查询时间缩短**60%以上**
- 短查询时间缩短**30%以上**

THANKS

SQL Server
vertica
D B 2
G B a s e
O r a c l e
达梦数据库
神舟通用
KingbaseES

2010

2014

2018

openGauss
OceanBase
ArkDB
RASESQL
HotDB
StellarDB
QianBase xTP
GoldenDB
云树Shard
MatrixDB
DynamoDB
SinoDB
DolphinDB
FastData
Galaxybase
KunDB
GDB
GaussDB
PolarDB
KunDB
Spacture
SequoiaDB
OushuDB
ArgoDB
开务数据库
GreatDB
MongoDB
TDSQL
TiDB
Tapdata
StarRocks
UbiSQL