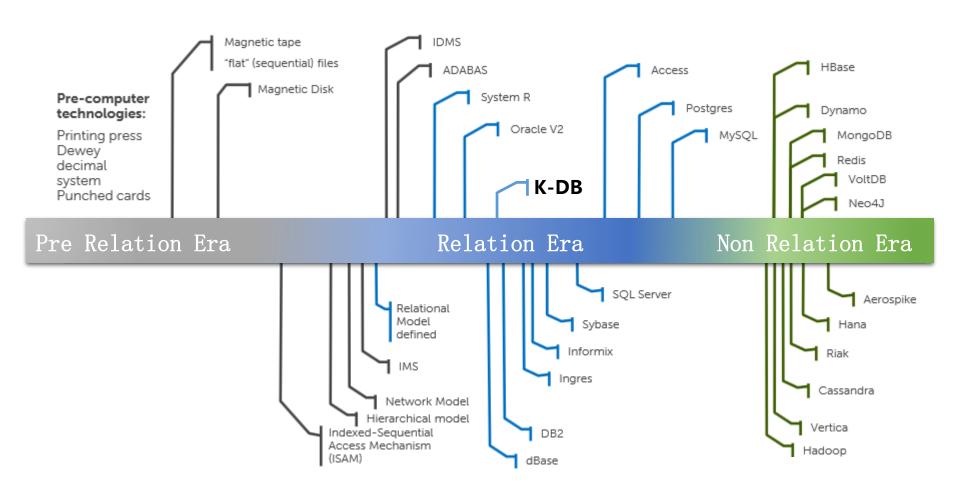
传统数据库运维及面临的挑战

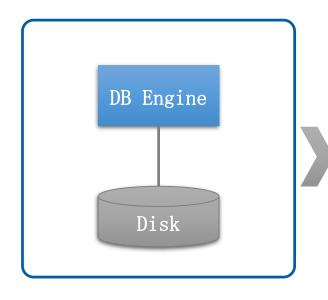
浪潮 架构师

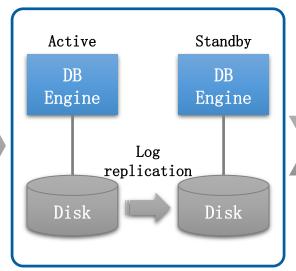
金学东

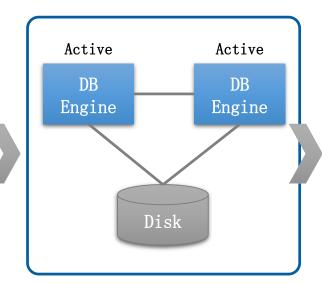
数据库技术发展

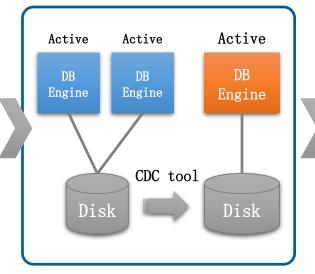


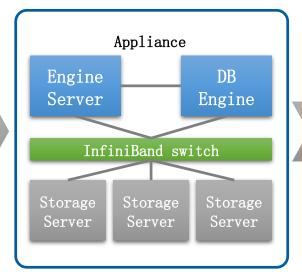
我们的故事

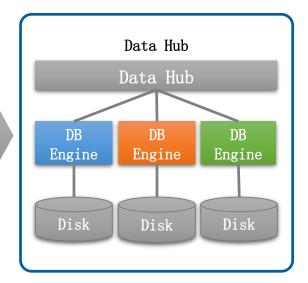










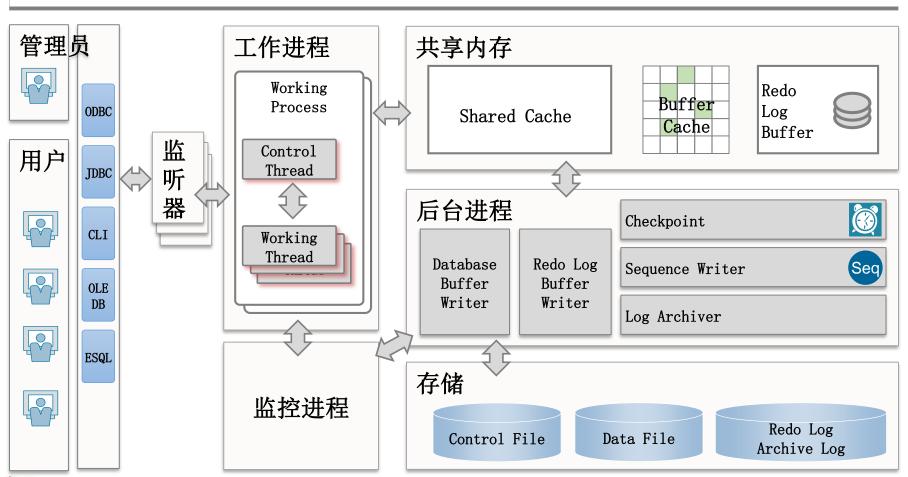


雏形

架构图

▶ K-DB数据库一开始便在工作进程上采用了多进程多线程架构,但是在后台进程、 日志、数据文件管理方式上还是更多地遵循了知名数据库的成熟架构。

K-DB RDBMS Instance



监控方法 - 日志跟踪

▶ 早期的 K-DB数据库并不具备足够多的监控事件及统计工具。初期的数据库监控及问题分析,基本依靠数据库日志及Trace日志进行分析并解决。

DBMS log

```
K-DB started up. protocol version 2.8
    MTHR(tbsvr) is ready.
                                            proc_idx(tid) :0, ospid :8338
    CTHR(tbsvr_WT001) is ready. proc_idx(tid) :1, ospid :8349
                                                  proc_idx(tid) :2, ospid :8350
    CTHR(tbsvr_WT002) is ready.
                                                  proc_idx(tid) :3, ospid :8351
    CTHR(tbsvr_WT003) is ready.
    AGENT(tbsvr_SEQW) is ready. proc_idx(tid):4, ospid:8352
    LGWR(tbsvr_LGWR) is ready. proc_idx(tid) :5, ospid :8353
    LARC(tbsvr_LOGA) is ready. proc_idx(tid) :6, ospid :8354
CKPT(tbsvr_CKPT) is ready. proc_idx(tid) :7, ospid :8355
    DBWR(tbsvr_BLK001) is ready. proc_idx(tid) :8, ospid :8356
    DBWR(tbsvr_BLK002) is ready. proc_idx(tid) :9, ospid :8357
    DB NAME - kdb
    CONTROL FILE - /home/kdb/database/kdb/c1.ctl
    BLOCK SIZE - 8192
    NLS CHARACTER SET - UTFS
NOT NULL, /* ?? ID */
         EMP_NAME
                                VARCHAR (20)
                                                                            /* ??
        HIREDATE
                                                                             /* ????
         SALARY
                                NUMBER (8, 3)
                                                                             /* ??
                                                                                                */
        BONUS
                                NUMBER (8, 3)
         DEPT OF
                                VARCHAR (4)
                                                                             /* ????
                             VARCHAR (8)
        MANAGER
                                                                         /* ?? */
        CONSTRAINT EMPLOYEE_PK
             PRIMARY KEY (EMP_NO)
             USING INDEX
             PCTFREE
         TABLESPACE USERS
 2015/09/07 22:23:32.917
                                                    21 Internal Error with condition 'cvi->cv->seqno != TCBLK_SEQNO_APP
 2015/09/07 23:26:02.603
                                                       21 Internal arrow with conduction of 17-00-3600; internal arrow with conduction of 4 slow DB BLOCK I/O ( 26ms per block) at 2015/09/07 22:16:24.808 4 slow LOG BLOCK I/O ( 66ms per block) at 2015/09/08 00:01:15.845 4 slow LOG BLOCK I/O ( 26ms per block) at 2015/09/09 00:01:40.665 4 slow LOG BLOCK I/O ( 60ms per block) at 2015/09/09 00:01:10.287
2015/09/07 20:66:07.058 [CLD] [0] 2015/09/08 00:06:07.058 [CLD] [0] 2015/09/08 00:06:07.058 [CLD] [0] 2015/09/09 00:07:47.363 [CLD] [0] 2015/09/09 10:32:09.768 [FRM] [0] 2015/09/09 10:32:09.838 [CLK] [0] 2015/09/09 10:32:09.845 [CLK] [0] 2015/09/09 10:32:09.845 [CLK] [0] 2015/09/09 10:32:09.858 [CLC] [0] 2015/09/09 10:32:56.953 [FRM] [0]
 2015/09/08 00:06:07.058
                                                     11 get shutdown request msg. mode=1
                                                        7 CHECKPOINT done 3ms, 6 blocks at 0000.000c7b7c/3.335759.26.
                                                        7 CHECKPOINT done
                                                                                                      4 blocks at 0000.000c7b7c/3.335759.26.
                                                                                          Oms,
                                                        7 CKPT write disabled
                                                       7 LOG FILE CLOSE: thread=0, current log=3, rba=3.335758
7 REDO THREAD #0 CLOSED: 0000.000c7b7c(0000.000c7b7c), 0 opens
                                                        7 DB CLOSE: C 0000.000c7b7c, $ 0000.000c7b7c, L 3.335758/0000.000c7b
                                                        0 server shut down.
                                                        0 show tip [kdb.tip]
 # tip file generated from /home/kdb/config/tip.template (Thu Aug 6 16:58:04 CST 2015)
```

Trace log

SET AUTOT[RACE] {OFF | ON | TRACE[ONLY]}
[EXP[LAIN]] [STAT[ISTICS]]

```
BEGIN
      DBMS_OUTPUT. ENABLE:
  FOR PSTSR50T_CUSRSOR_REC IN PSTSR50T_CUSRSOR LOOP
      V_REQ_NO := PSTSR5OT_CUSRSOR_REC. REQ_NO
     V_ASSIGN_SEQ := PSTSR50T_CUSRSOR_REC.ASSIGN_SEQ;
V_SERVE_SEQ := PSTSR50T_CUSRSOR_REC.SERVE_SEQ;
V_PROD_SEQ := PSTSR50T_CUSRSOR_REC.PROD_SEQ;
      V_LICENSE_KIND := PSTSR50T_CUSRSOR_REC. LICENSE_KIND;
      V_LICENSE_CLS := PSTSR50T_CUSRSOR_REC.LICENSE_CLS;
      CURSOR COUNT : - PSTSR50T CUSRSOR%ROWCOUNT:
      DBMS_OUTPUT.PUT_LINE(' ---- 整办導動?PSTSR5OT COUNT: ---' || CURSOR_COUNT);
      DBMS_OUTPUT. PUT_LINE (' REQ_NO
      DBMS_OUTPUT. PUT_LINE ('ASSIGN_SEQ
                                                               V ASSIGN SEO)
     DBMS_OUTPUT.PUT_LINE('SERVER.SEQ : 'V_SERVE_SEQ);
DBMS_OUTPUT.PUT_LINE('PROD_SEQ : 'V_PROD_SEQ);
DBMS_OUTPUT.PUT_LINE('LICENSE_KIND : 'V_LICENSE_KIND)
DBMS_OUTPUT.PUT_LINE('LICENSE_CLS : 'V_LICENSE_CLS);
--PINSTIOT 聯灣拡聯机權幹想制 駒れ權敵检應減€ 20061024 ~ 20061104 輕?聯吗楞 競办溥敬办漖
--LICENSEKIND 鞮€ LICENSECLS 蠢?勒猿   氘專維勒?UPDATE 頃懷雄.
         SET LICENSE_KIND - V_LICENSE_CLS
     SHI DICENSE_KIND = V_LICENSE_KIND

WHERE REQ_NO = V_REQ_NO

AND ASSIGN_SEQ = V_ASSIGN_SEQ

AND SERVE_SEQ = V_SERVE_SEQ
        AND PROD_SEQ =V_PROD_SEQ;
     DBMS_OUTPUT. PUT_LINE ('UPDATE SUCCESS:' || CURSOR_COUNT);
    DBMS_OUTPUT. PUT_LINE(' 職動油 競办溥載?鞋?' || CURSOR_COUNT);
END UPDATE PSTSR50T LICENSE:
08/06 17:11:01.107756 [DDL][0] 20 ddl_fram:207 DDL execution succeeded 08/06 17:11:01.109107 [DDL][0] 20 ddl_fram:185 Executing DDL: CREATE OR REPLACE PROCEDURE UPDATE_RECRUIT_STATE
   P_REQUEST_NO RCSCHOOT. REQUEST_NO%TYPE,
  P_REQUEST_SEQ RCSCHOOT. REQUEST_SEQWTYPE,
                   OUT VARCHAR2
   V_REQUEST_NO RCSCHOOT. REQUEST_NO%TYPE := P_REQUEST_NO
```

监控方法 - kdBinary

GENERAL	Parameter	Infomatio	n n							į		
1 - Inst	=======		==							į		
Current Undo	Parameter	Name		1010000001111	Value							
	CWS_RSBTBL				80							i
SEGMENT ID IA		IZE			35782	6560					Х	ACTS
Tablespac	ce Infomation											0
Tablespace	e Name	Bytes(MB)	Used(MB)	Percent(%)	Free(M	в)	Free(%) Maxi	Bytes	(MB)	¦	0
SYSSUB		20	18	91.25		2	8.	75	32	,768		0
	Count (Archiv		80	80.00		20	20.	00	32	.768		
===================================		erog only)	\ \									
LOG HISTORY (Since Last Mo	onth)		\ _								
Hourly 00 01	02 03 04	05 06 07	08 09 10	11 12 1	3 14 15	16	17 1	8 19	20	21	22	23
12/27 1 0 12/29 0 0		0 0 0 0 0 0	0 0 0		0 0 0	0 0	0	0 0	0	0	0	0
[ARCHIVED LOG				<u></u>	770							
undo segment	_HASH_JOIN	_BUCKET_SI	ZE_V		16					i		
Undo Usage T					1					ŀ	 	j
	_RCACHE_BL	CKET_COUNT		A	256					į		
5. Savel	_WLOCK_BUC	KETSET_CNT L_INIT_SIZ	E		64 400					i		
	_WTHR_PER_	PROC			10							

备份/恢复

▶ K-DB数据库在早期,随支持逻辑备份及物理备份两种模式,但是在最重要的物理 备份领域,只支持简单的冷备及基于操作系统命令实现的增量备份。

逻辑备份

\$ kdexport/kdimport username=testuser
password=testpassword sid=test file=backup.dat full=y
\$ kdexport/kdimport cfgfile=backup.cfg

表模式

Table definitions
Table data
Owner grants
Owner indexes
Table constraints
Table triggers

用户模式

Table definitions
Table data
Owner grants
Owner indexes
Table constraints
Table triggers
Views
Private synonyms
Sequences
PSM

全库模式

Table definitions
Table data
Owner grants
Owner indexes
Table constraints
Table triggers
Views
Private synonyms
Sequences
PSM
Roles
All synonyms

物理备份

- 停机备份
 - 硬拷贝数据库相关文件
- 在线备份
 - 数据库运行状态下备份

运维方案

正常运维

状态检查

实例状态检查 连接状态检查 日志信息检查 锁信息检查 空间信息检查

系统维护

参数修改 重做日志维护 表空间维护

系统备份

逻辑备份 Export 物理备份 全库冷备 在线备份

维护报告

运维日报 运维周报 运维月报

紧急运维

问题出现

业务报错 业务性能低下 数据库无反应 数据库宕机

错误信息收集

实例状态信息 连接状态信息 日志信息收集 锁信息收集 空间信息收集

紧急恢复

实例恢复 介质恢复 逻辑恢复 性能优化

分析问题原因

实例状态分析 连接状态分析 日志信息分析 锁信息分析 空间信息分析

解决方案应用

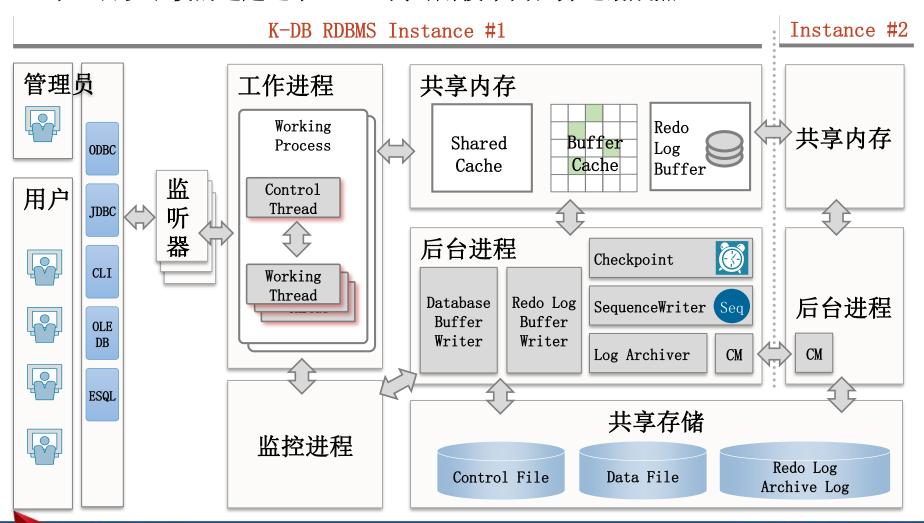
SOL优化 表结构优化 索引优化 参数优化

高可用

架构图

INSPUT 浪潮

▶ K-DB早期版本还未来得及站稳市场,基于 RAC架构的高可用架构早已成为行业标准,日以继夜的追赶之下,K-DB高可用技术面世并逐渐成熟。

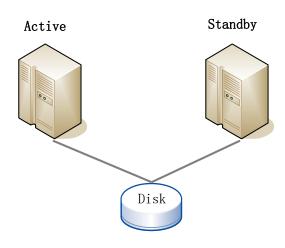


架构升级 - Active-Standby

▶ K-DB首先实现了基于共享存储的 H/A切换及基于日志同步的容灾架构。尤其,实现了容灾模式下的 Standby只读功能后,最后的难关便只剩下 RAC架构的实现。

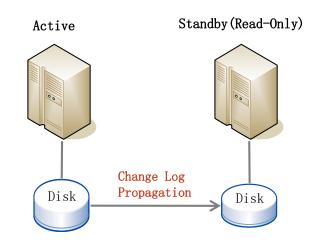
Active-Standby (Shared)

- 基于共享磁盘的 Active-Standby架构
- 同时有且只有一个节点处于运行状态
- 基于第三方HA软件进行切换



Active-Standby (Replication)

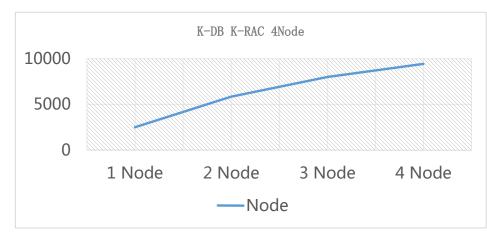
- 基于日志同步的 Active-Standby架构
- 同时有且只有一个节点以主节点模式运行, 其他节点以 Standby模式(必要时可读)运行
- 基于日志同步插件(类似于 DataGuard) 实现数据同步

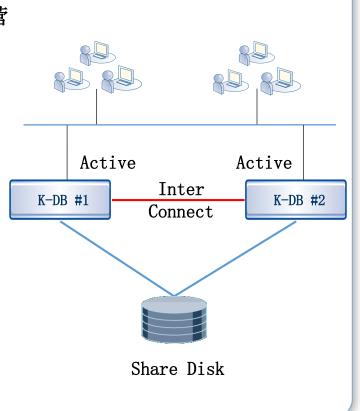


架构升级 - RAC

▶ K-DB首先实现了基于共享存储的 H/A切换及基于日志同步的容灾架构。尤其,实现了容灾模式下的 Standby只读功能后,最后的难关便只剩下 RAC架构的实现。

- ▶ DB 和系统发生故障时保证运营系统无中断运营
- ▶ 特定节点发生故障时快速恢复到正常节点
- ➤ 实现 Global Cache将 Disk I/O最小化
- ▶ 水平的添加Active节点,实现平滑扩容
- ▶ 局限: 扩容节点数量存在限制





监控方法 - APM报告

▶ 基于架构的升级及改造,大大提升了K-DB的监控能力。基于快照实现的统计报告功能推出,大大减轻了分析压力(早期为 TXT版本)。

*** APM (Automatic Performance Monitoring) Report ***	2. Worklos	2. Workload Summary						
DB Name : kdb TAC : YES Instance Cnt : 2	Call		Per Second	Per TX	Per Exec	P		
Instance cnt : 2 Release : 1 SP1 r110806 HOST CPUs : 128 (cores: 8, sockets: 64)		B Time(s):	226.60	0.37	0.02			
Interval condition: 2015-10-19 16:47:14 ~ 2015-10-19 16:57:14 (#=1 Report Snapshot Range: 2015-10-19 16:47:14 ~ 2015-10-19 16:57:14	reported) 0.02	Redo Size:	4, 662, 465. 56	7, 553. 75	398.87			
Report Instance Cnt : 1 (from Instance NO. 'ALL') Elapsed Time : 10.00 (mins)	1.1	.cal Reads:	169, 796. 53	275.09	14.53			
DB Time: 2,266.02 (mins) Avg. Session #: 3.00	11	ck Changes:	52, 392. 18	84.88	4.48			
	Physi	.cal Reads:	527.81	0.86	0.05			
** Report Overview **	11	al Writes:	5, 787. 25	9.38	0.50			
	ii t	Jser Calls:	13, 203. 25	21.39	1.13			
1. System Overview Section	ii ii	Parses:	598.39	0.97	0.05			
		ard Parses:	0.00	0.00	0.00			
1.1 CPU Usage	0.00	Logons:	0.00	0.00	0.00			
F/G Host CPU Usage(%)	0.00	Executes:	11,689.14	18.94	1.00			
DB CPU Instance# Usage(%) Busy User Sys Idle IOwait	0,00	Rollbacks:	33.82	0.05	0.00			
0 6.3 8.7 6.7 2.0 91.3 12.6	0.05	msactions:	617.24	1.00	0.05			
1.2 Memory Usage	2.1 Worklo	ad Stats				:====		
1.2 Memotry Usage	[
HOST Mem Size : 1.E+06M	DB				DB			
Total SHM Size : 71,680M Buffer Cache Size : 40,960M	Time(%)	Category		Sta	t Time	e(ms)		
Avg. Shared Pool Size : 2.5E+04M Avg. DD Cache Size : 1.29M	1200 (07							
Avg. PP Cache Size : 3.45M Max DD, PP Cache Size : 1.8E+04M	Request Ser	vice Time -			- 135,961	ι, 257		
DB Block Size : 8K Log Buffer Size : 1,024M	100.00			SQL processing tim	e 124,349), 180		
	91. 46 1111111111 6. 99		SQL process	ing (batchupdate) tim	e 9,506	ò, 149		
2. Workload Overview Section	11			reply msg etc tim	e 1,609	9, 267		

K-APM

基本信息查询

- CPU、内存使用率
- 负载状态显示
- 集群状态显示
- TOP 5 等待事件
- I/0运行状态
- PGA状态
- TOP N SQL

详细信息查询

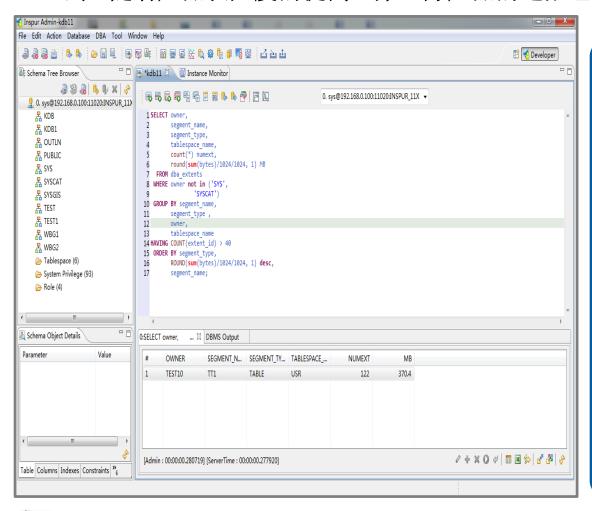
- OS 状态
- 共享池及各种缓存状态
- 等待事件详细描述
- 锁状态
- TOP N SQL及执行计划

其他

• 数据库运行参数

操作管理 - kdAdmin

▶ 作为数据库开发与管理的便利性考虑,第一款图形化数据库开发工具 kdAdmin面 试。随着产品认知度的提高,第三方产品的适配也逐渐开始。



kdAdmin功能

- 开发
 - 用户及对象信息查询
 - SQL语句编辑与执行
 - 存储过程编辑与执行
- 监控
 - 会话监控及管理
 - 实例监控
 - 事务监视
- 管理
 - 参数管理
 - 表空间管理
 - 数据库导入与导出
- 分析及优化
 - 分析管理器
 - 图形化统计报告

备份/恢复 - RMGR & Flashback

▶ RMGR - 数据库自动化备份与恢复工具。我们对其满怀期望,但是开发初期,只有简单的在线全备与增量备份功能。闪回,也只能实现对于执行语句的回退。

RMGR备份恢复工具

- RMGR 基本功能
 - Online Full Backup
 - Incremental Backup
 - Automatic Recovery

闪回功能

Flashback Query

运维方案

状态检查

RAC 状态检查 主备状态检查 CDC 状态检查 实例状态检查 连接状态检查 日志信息检查 空间信息检查

系统维护

参数修改 重做日志维护 表空间维护

系统备份

逻辑备份 Export RMGR 备份 全库冷备 在线备份

维护报告

运维日报 运维周报 运维月报

问题出现

业务报错 业务性能低下 集群失效 数据库无反应 数据库宕机

错误信息收集

RAC 状态检查 主备状态检查 CDC 状态检查 实例状态信息 连接状态信息 话信息收集 空间信息收集

紧急恢复

RAC Fail-Over 主备切换 RMGR 恢复 实例恢复 逻辑恢复 性能优化

分析问题原因

RAC 状态分析 主备状态分析 CDC 状态分析 APM 分析报告 日志信息分析

解决方案应用

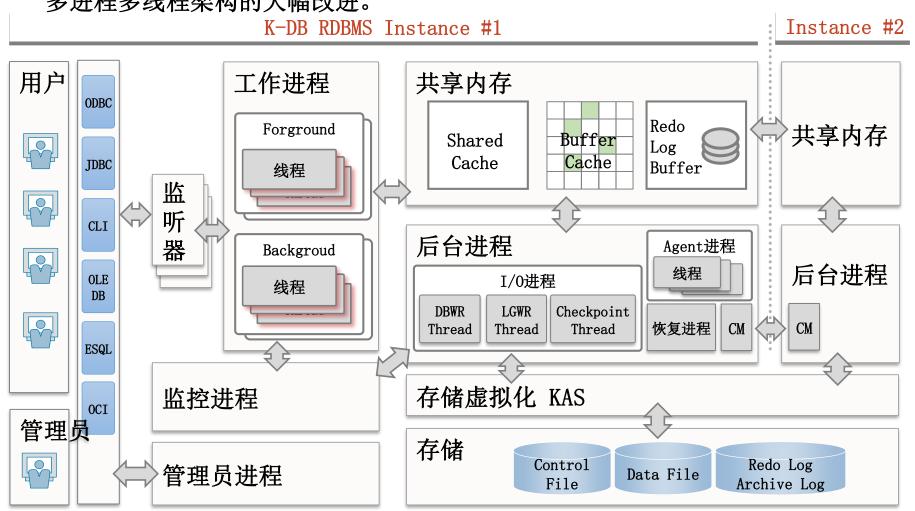
SOL优化 表结构优化 索引优化 参数优化 修改备份策略 集群架构优化



现在

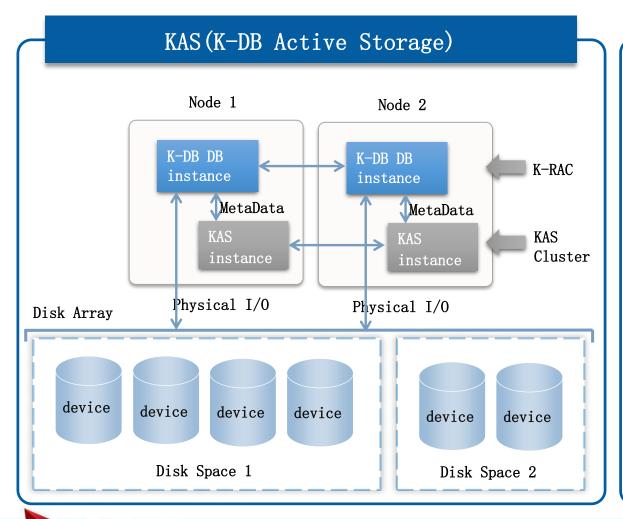
架构图

▶ K-DB 成熟版本的体系架构图,增加了KAS存储虚拟化、而后台进程也进行了基于 多进程多线程架构的大幅改进。



架构升级 - 存储虚拟化

▶ 存储虚拟化技术,让 K-DB彻底摆脱了第三方集群管理软件的束缚,具备了自主构建高效 RAC集群的能力。



KAS功能特点

Mirroring

•利用2-way或3-way镜像保 障磁盘故障时的可用性

Striping

•在多个磁盘中分布存储数据,实现并行磁盘I/0

Rebalancing

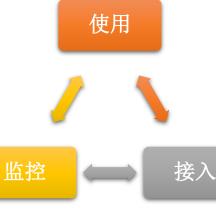
•将数据平均存储到所有磁盘中,提高并行磁盘I/0 效果

架构升级 - 三权分立/精细化审计

▶ 与其他数据库产品单独提供安全版数据库不同,K-DB基于安全性考虑,在最新版产品中将具体实现数据库管理的三权分立思想,丰富审计功能及体系。

三权分立

- 系统管理员
 - 系统管理员负责数据库的日常管理与使用
- 审计管理员
 - 审计管理员负责数据库的审计、内在风险分析、异常状态感知等功能
- 安全管理员
 - 安全管理员负责管理强制访问(VPD)、其他用户的数据 库访问及访问权限等功能

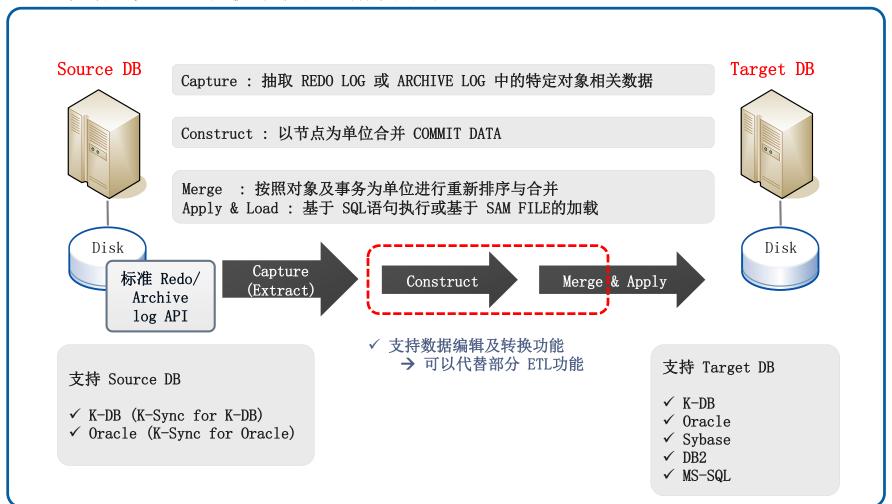


精细化审计

- 实时告警
 - 对于特定的敏感操作 记录相关日志,并通 过邮件或短信等方式 通知客户
- 潜在问题分析
 - 对于审计日志进行自 动化分类与统计,并 输出分析报告
- 异常状态感知
 - 对于数据库的非法访问(如短时间内基于同一用户频繁访问敏感信息)进行自动化监测

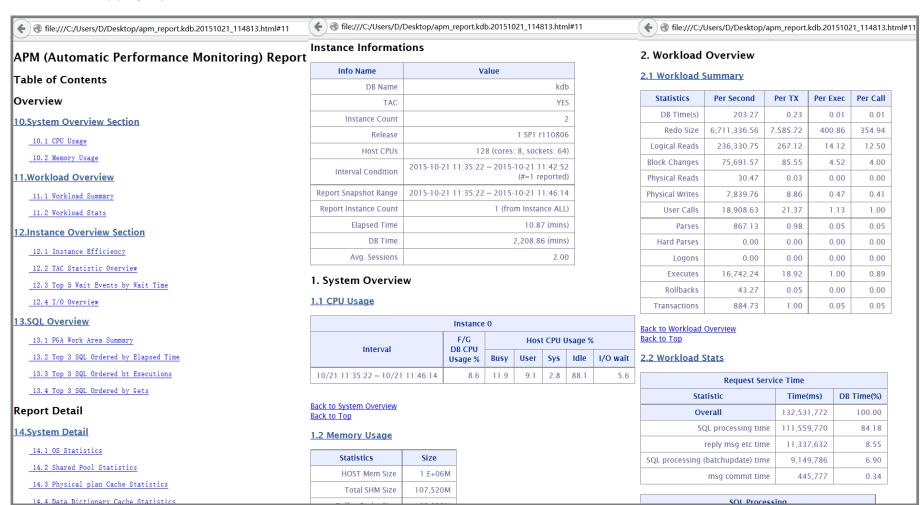
架构升级 - CDC数据同步

▶ 基于 CDC工具的高性能异构数据库同步工具,更是对于客户混合数据库运营环境 下数据同步问题提供了优秀的解决方案。



监控方法 - APM报告(HTML)

▶ 基于网页版统计报告的实现,大大提升了报告的可读性,DBA不用再盯着古老的 黑白屏幕了。



监控方法 - 图形化实时监控

inspur 浪潮

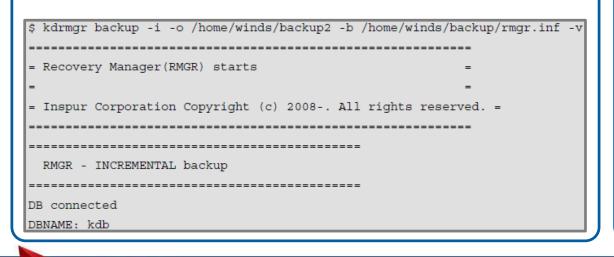
K-APM 数据库系统 实时监控平台

备份/恢复 - RMGR & Flashback

▶ RMGR终于不只是架子,不仅能够实现数据库的灵活备份与恢复,进一步提供了数据块改动检测功能。而闪回功能,已经实现了数据库级别的恢复。

RMGR备份恢复工具

- RMGR 基本功能
 - Online Full Backup
 - Incremental Backup
 - Block Change Tracking
 - Automatic Recovery
 - Tablespace/Datafile级别备份及恢复
 - 恢复 Auxiliary Database



闪回功能升级 Flashback Database Flashback Query

运维方案

统计分析

TOP Session

TOP SQL

TOP Wait

TOP Program

TOP Module

发现问题

业务性能低下

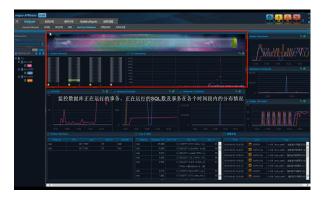
触发阈值

TOP User

Lock Tree



自动化分析



自动化监控



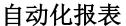
系统优化

SOL优化 表结构优化 索引优化 参数优化 修改备份策略 集群架构优化



实时监控

集群状态监控 会话状态监控 告警信息监控 锁信息监控 空间信息监控 Disk I/O监控 CPU、内存监控



巡检报告

阶段统计报告 自定义报告



新的挑战

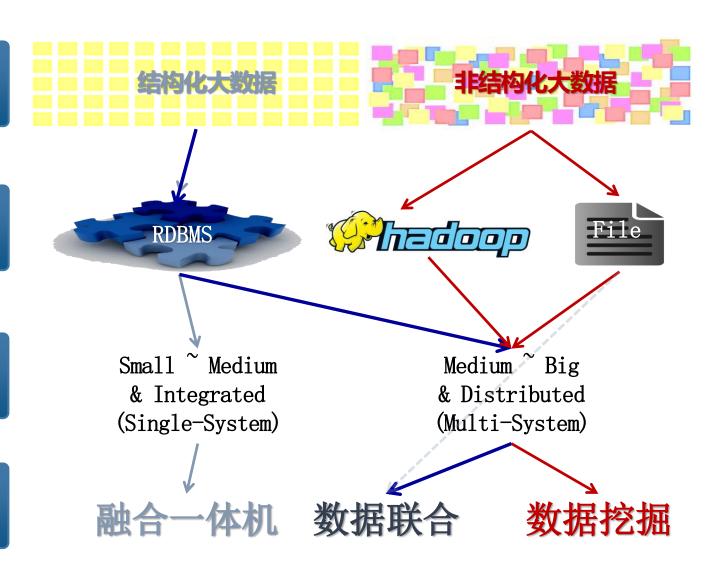
大数据环境下的新挑战

Data Source Structure

Data Store System

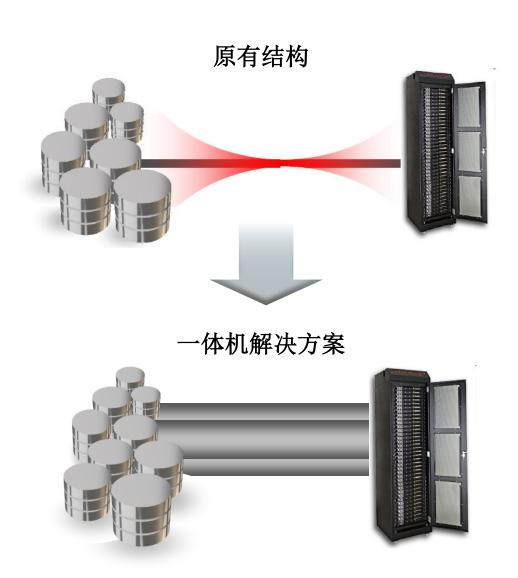
Data Size & Consolidation type

Solution

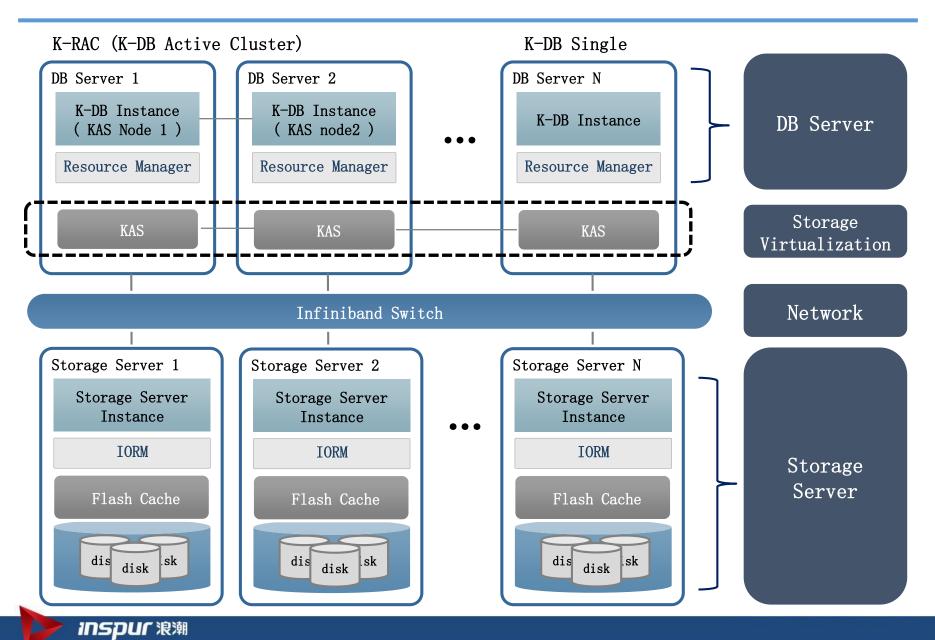


融合一体机

传统数据库瓶颈与一体机解决方案

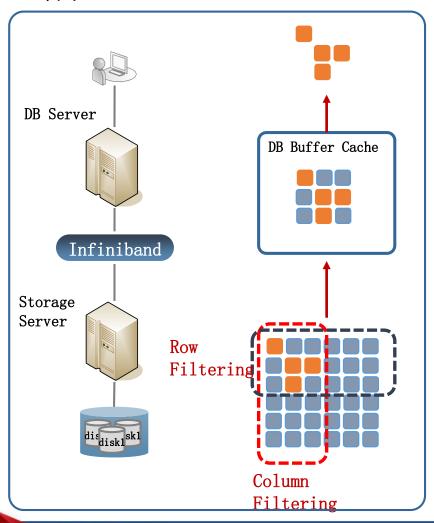


New Architecture



Smart Filtering

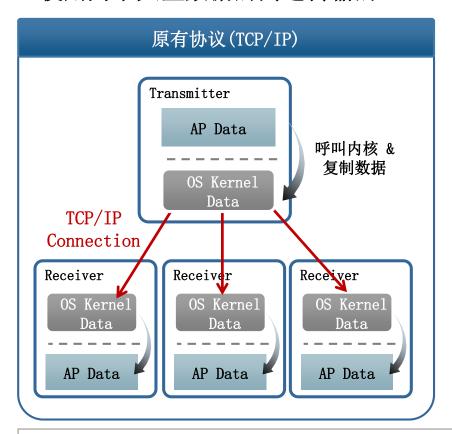
▶ 从存储服务器中只提取必要的数据传送到DB,减少DB的作业负担,缩短数据处理时间

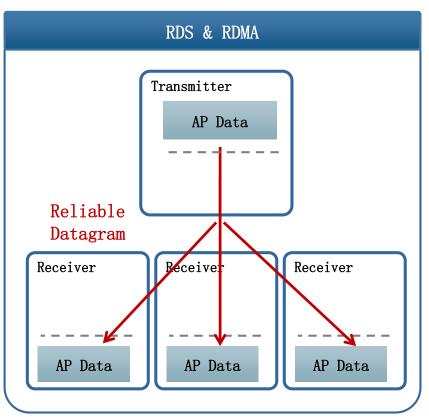


技术	内容
Smart Scan	在存储设备中只抽取符合SQL条件Row,和Column 传输到DB
Storage Index	把储存在磁盘的数据概要信息实时进行管理 → 忽略不符合SQL条件的Disk I/O

RDS & RDMA

▶ 使用为了大量数据的高速传输的 RDS&RDMA 协议→ Low Overhead低开销

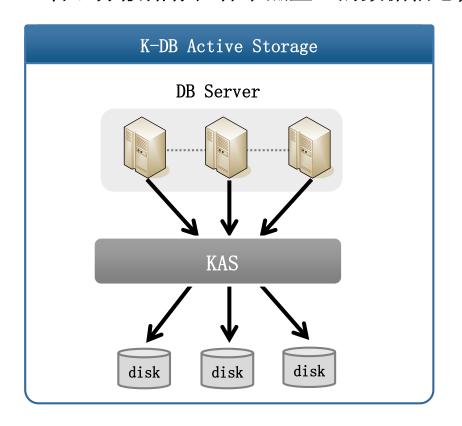


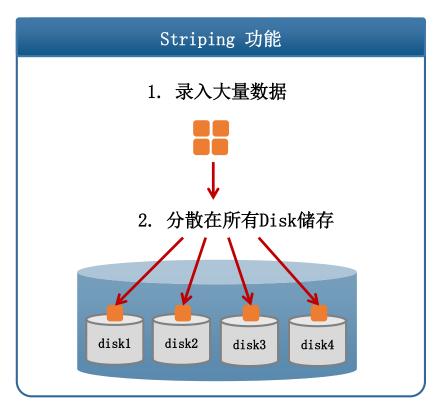


- RDS (Reliable Datagram Socket): 可信的 1:N传输协议,并且对比TCP/IP 减少负荷
- RDMA (Remote Direct Memory Access)
 - :直接参照 AP领域的内存调用,避免OS间接访问内核 →减少CPU负荷并且可以简短Latency (延时)

KAS (K-DB Active Storage)

▶ 管理分散储存在各个磁盘上的数据信息,通过 Disk I/O并行处理提高性能

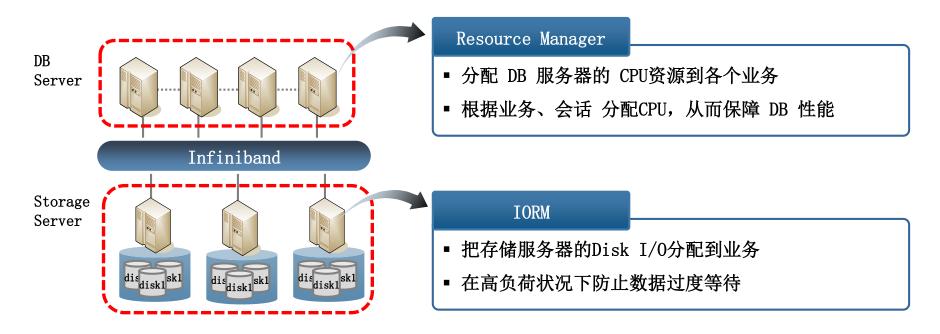




- KAS:管理分散储存在各个磁盘上的Clustered File System
 - → 并行Disk I/O: 把数据分散储存在各个Disk提高I/O 性能

Resource Management

▶ 根据业务分配CPU和Disk I/O 的资源,业务负荷大的状况下,可以保障性能和有效使用资源

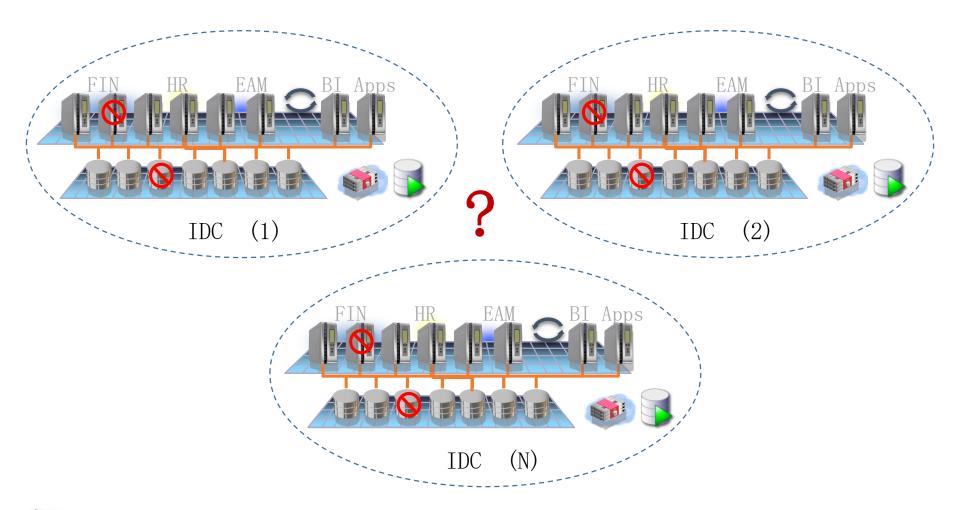


- 把DB服务器的CPU和存储服务器的Disk I/O分配到各个业务上保障性能
- 重要的业务在DB和存储服务器高负荷时也能充分保障性能

数据联合

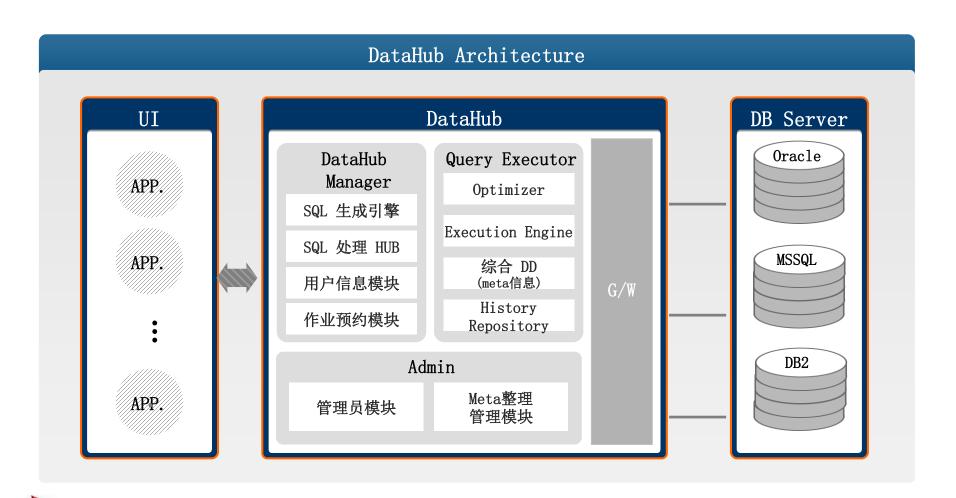
分布式数据中心的数据联合

▶ 在地域分布、上百上千的数据中心需要联合分析计算的场景下,需要寻求新的数据架构



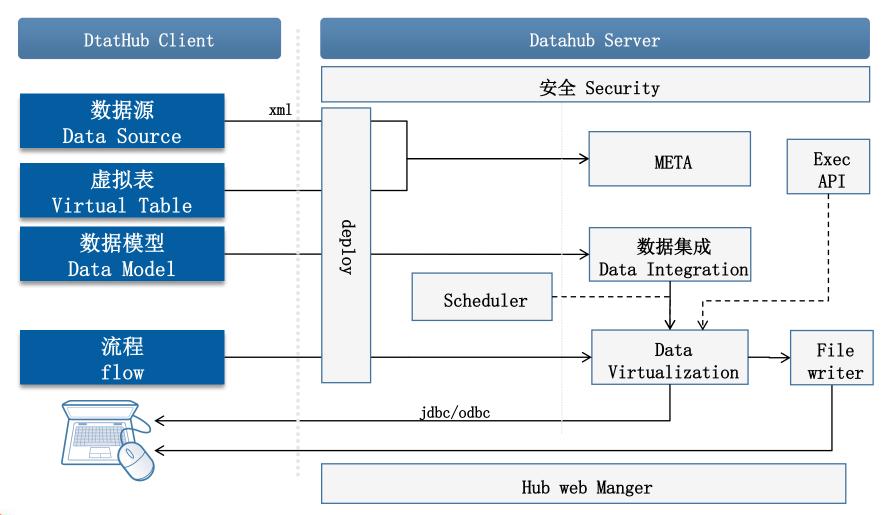
Data Hub架构

▶ Data Hub可以通过单纯连接已有的多个异构数据库系统,便能够轻松构建跨数据库、跨地域的集群数据分析系统。



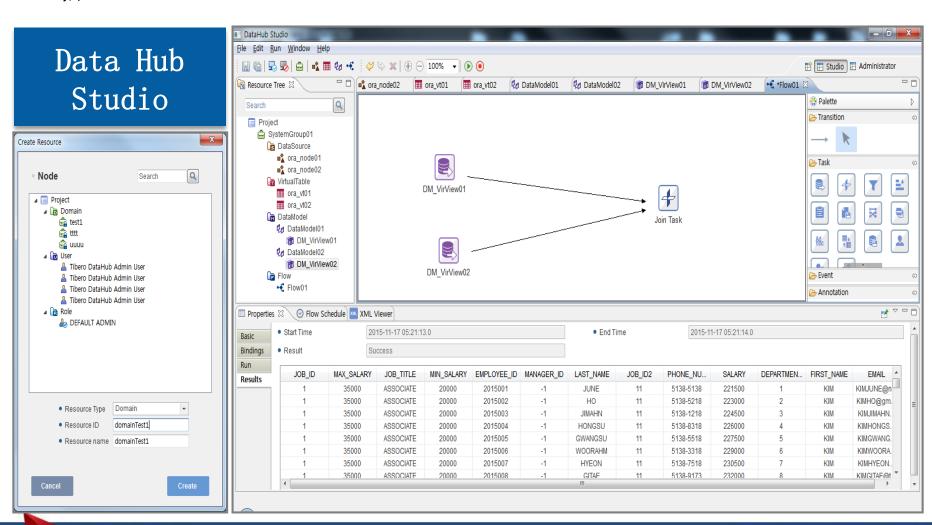
Data Hub架构

▶ 通过数据源镜像、数据对象建模、开发业务流这种层层递进的开发模式,Data Hub可以实现对于后端海量数据库集群的模拟与数据抽取。



Data Hub 业务开发工具

▶ 基于 GUI 拖拽方式的业务开发模式,使后期数据管理及分析业务的开发无比轻 松。



全新架构下的运维

全新架构特点



全新运维需求









