



中立安全·赋能产业

基于Serverless的USQL数据湖分析实践

UCloud优刻得产品总监 张晓康



产品理念—客户想要的是数据分析



成本低

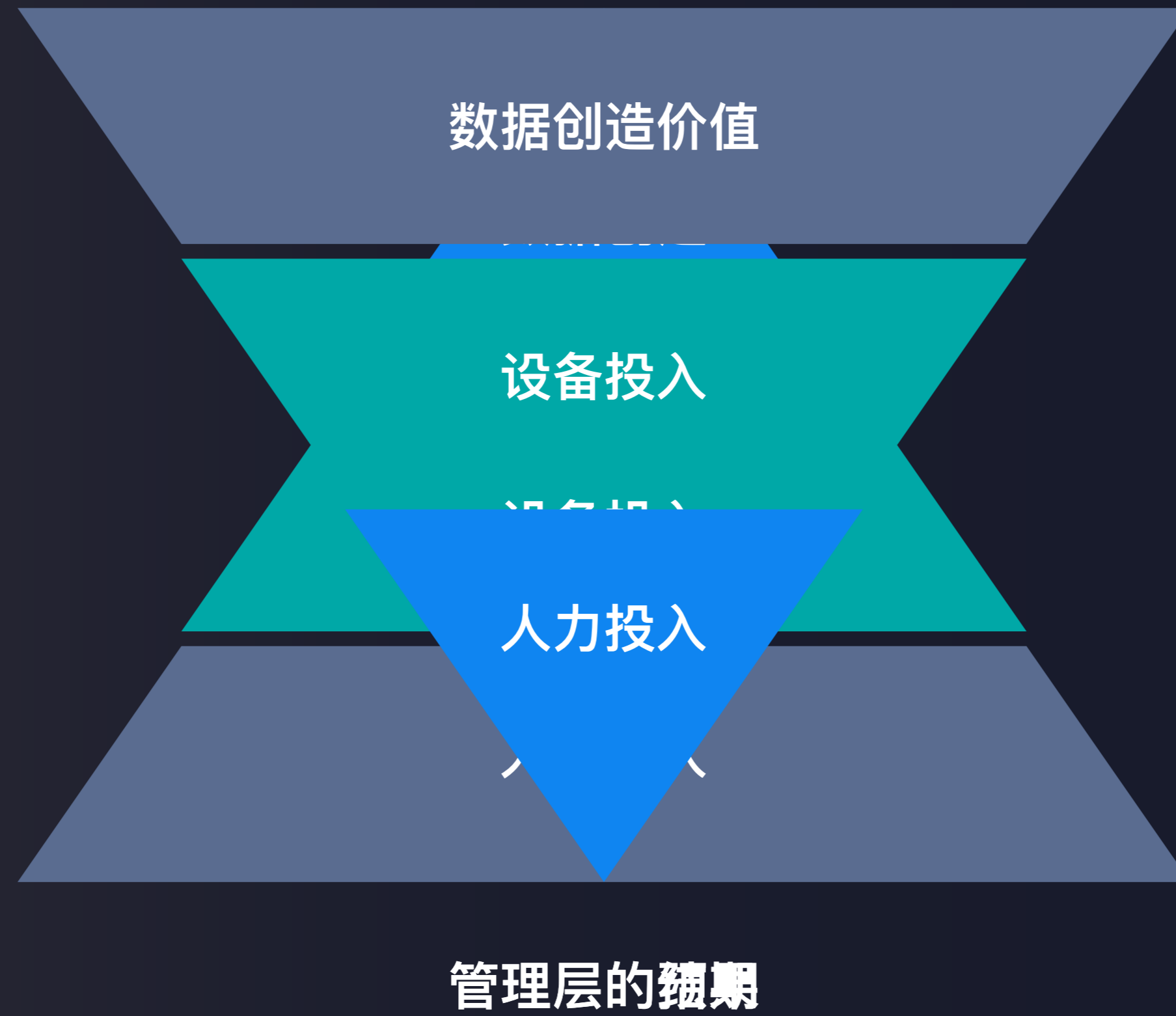


分析简单



无需运维

投入与产出



131.5元/SQL

业务数据状况

300G/日增

120T/总量

50个SQL/日运行

240万/年

配套

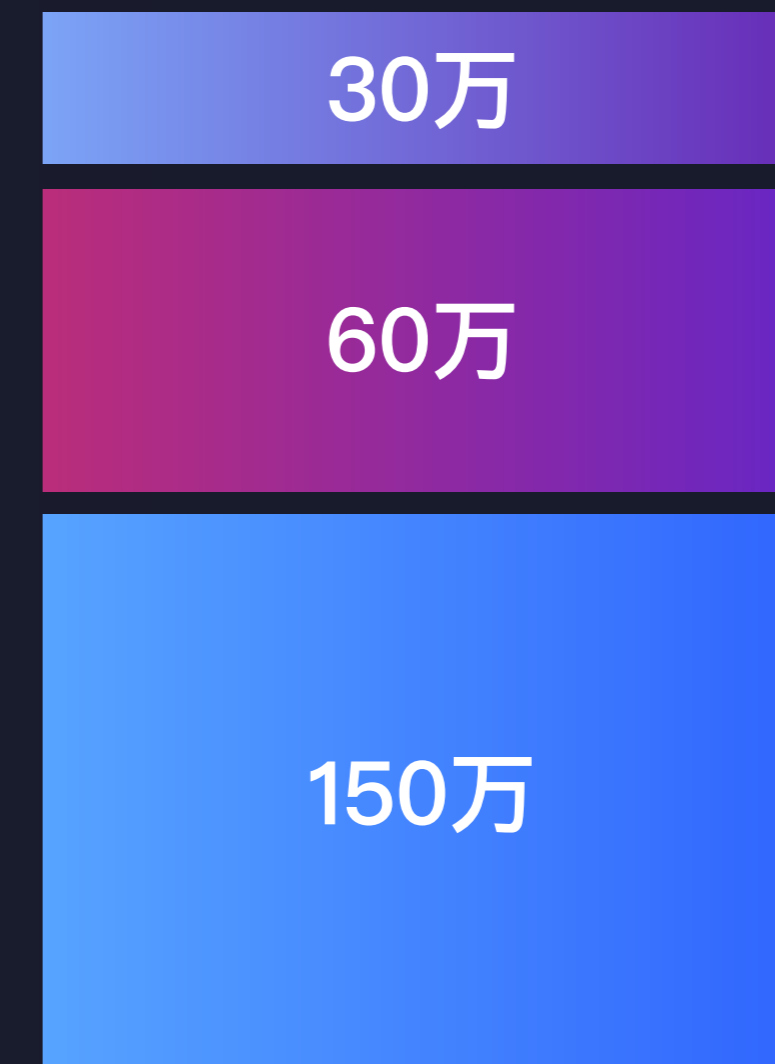
30万

设备

60万

人力

150万



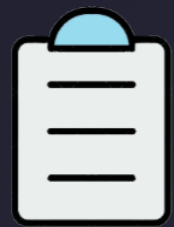
存储成本优化

积累了1年的数据

65%存储空间闲置



副本1



副本2



副本3

UFile对象存储



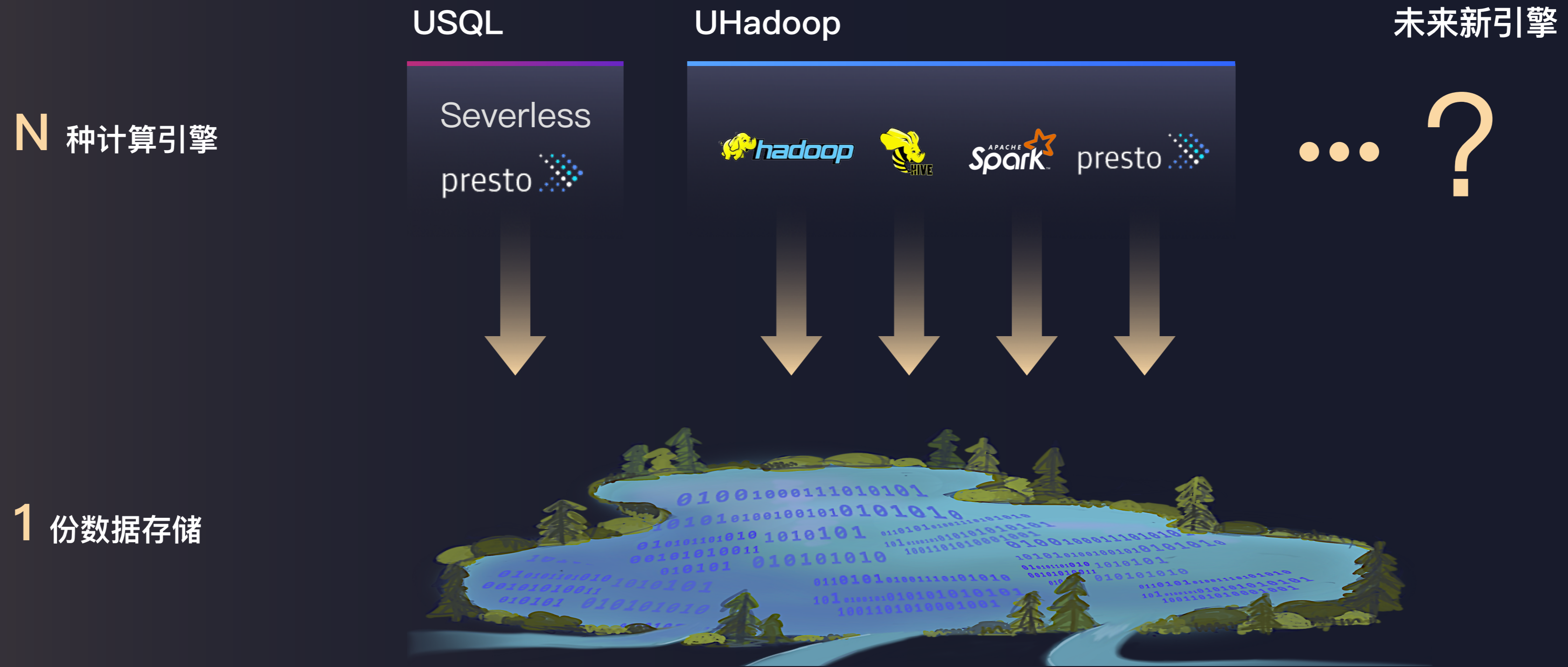
计算成本优化



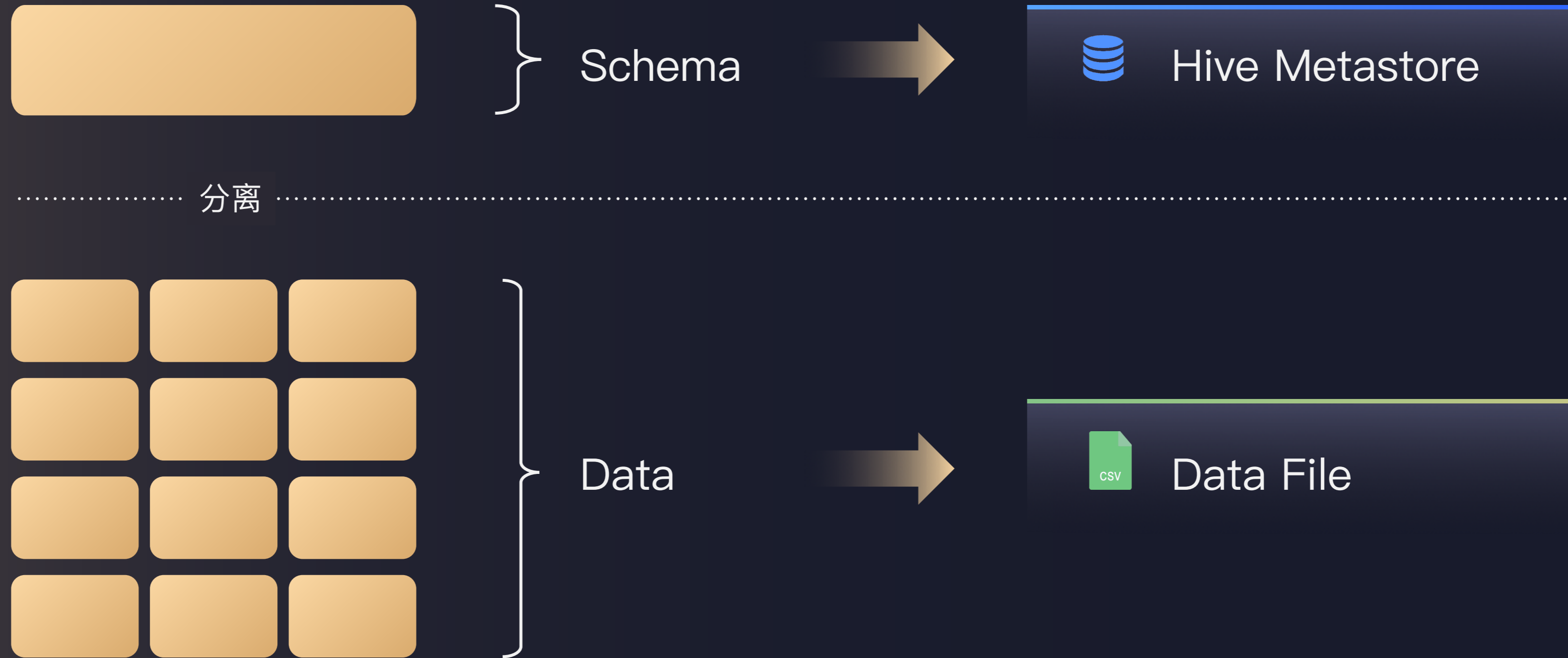
数据湖

集中存储任意规模和任意格式的数据

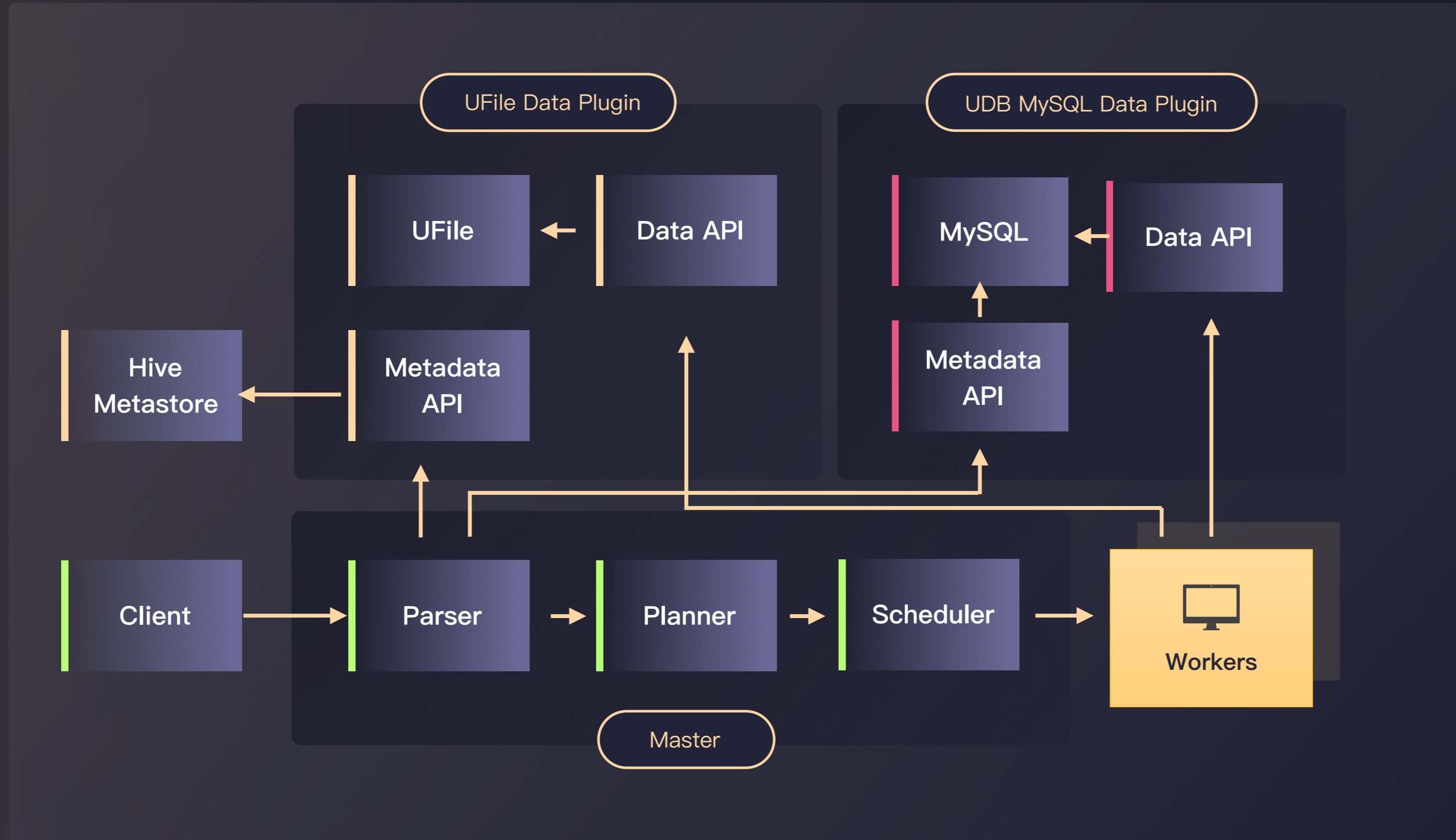
数据湖——计算存储分离



数据湖——表结构与数据分离



基于开源Presto项目



数据库分析 USQL

编辑SQL

已保存SQL

历史SQL

专家咨询

数据库

adx_offline

数据表

创建数据表

adx_ssp_req_log

bidid

sspip

delay

status

cd

ct

selfad

ssp_request

adx_pv_log

adx_click_log

adx_loss_log

adx_inmobi_log

adx_dsp_rsp_log

adx_twin_log

查询1 X

查询3 X

查询4 X

查询6 X

查询7 X

查询8 X

查询9 X

查询10 X

查询11 X

```
1 SELECT ssp_request.app.id AS appid,
2 COUNT(DISTINCT ssp_request.device.adid)-1 AS _count,0 as _count,
3 COUNT(DISTINCT ssp_request.device.ip)-1 AS ip_count,
4 COUNT(DISTINCT ssp_request.device.ip||ssp_request.device.ua)-1 AS _count
5 FROM adx_ssp_req_log
6 WHERE year=2019 AND month=05 AND day=04 AND ssp_request.device.os = 0 AND selfad = 0
7 GROUP BY ssp_request.app.id
8 HAVING COUNT(DISTINCT ssp_request.device.ip||ssp_request.device.ua) > 1
9 OR COUNT(DISTINCT ssp_request.device.adid) > 1
10 OR COUNT(DISTINCT ssp_request.device.ip) > 1
```

立即运行

示例SQL

保存

清空

运行结果 已运行: 347.331秒, 已扫描数据量: 77.07 GB, 结果总行数:111

全屏查看结果

下载结果

	appid	_count	_count	_count	_count
1	153	47197	0	69297	77432
2	420	119248	0	116695	150051

创建表结构 — 支持JSON嵌套结构

```
CREATE EXTERNAL TABLE person (
```

```
    name string,  
    age int,  
    city string,  
    friend:array<  
        struct<  
            name:string,  
            city:string,  
            birth:date  
        >  
    >  
>
```

```
)
```

```
ROW FORMAT SERDE
```

```
'org.openx.data.jsonserde.JsonSerDe'
```

```
LOCATION 'ufile://bucket_name/logs/'
```

结构映射

```
{  
  name: john,  
  age: 18,  
  city: shanghai  
  friend: [  
    {  
      name: mason,  
      city: beijing,  
      birth: 1995-03-12  
    },  
    {  
      name: julian  
      city: guangzhou,  
      birth: 1998-10-07  
    }  
  ]  
}
```

指定UFile位置

```
ufile://bucket_name/logs/2019/05/12/a.gz  
ufile://bucket_name/logs/2019/05/13/a.gz  
ufile://bucket_name/logs/2019/05/14/a.gz  
ufile://bucket_name/logs/2019/05/15/a.gz  
ufile://bucket_name/logs/2019/05/16/a.gz
```

数据湖分析 USQL

- 编辑SQL
- 已保存SQL
- 历史SQL
- 专家咨询

数据库

adx_offline

数据表

adx_ssp_req_log

bidid

sspip

delay

status

cd

ct

selfad

ssp_request

adx_pv_log

adx_click_log

adx_loss_log

adx_inmobi_log

adx_dsp_rsp_log

adx_twin_log

- 查询1 X
- 查询3 X
- 查询4 X
- 查询6 X
- 查询7 X
- 查询8 X
- 查询9 X
- 查询10 X
- 查询11 X

```
1 SELECT ssp_request.app.id AS appid,
2 COUNT(DISTINCT ssp_request.device.adid)-1 AS _count,0 as _count,
3 COUNT(DISTINCT ssp_request.device.ip)-1 AS ip_count,
4 COUNT(DISTINCT ssp_request.device.ip||ssp_request.device.ua)-1 AS _count
5 FROM adx_ssp_req_log
6 WHERE year=2019 AND month=05 AND day=04 AND ssp_request.device.os = 0 AND selfad = 0
7 GROUP BY ssp_request.app.id
8 HAVING COUNT(DISTINCT ssp_request.device.ip||ssp_request.device.ua) > 1
9 OR COUNT(DISTINCT ssp_request.device.adid) > 1
10 OR COUNT(DISTINCT ssp_request.device.ip) > 1
```

立即运行

示例SQL

保存

清空

运行结果 已运行: 347.331秒, 已扫描数据量: 77.07 GB, 结果总行数:111

全屏查看结果

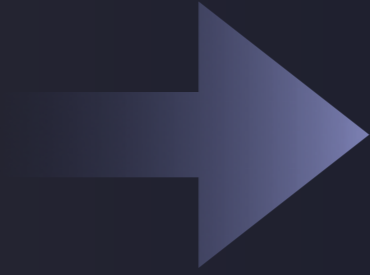
下载结果

	appid	_count	_count	_count	_count
1	153	47197	0	69297	77432
2	420	119248	0	116695	150051

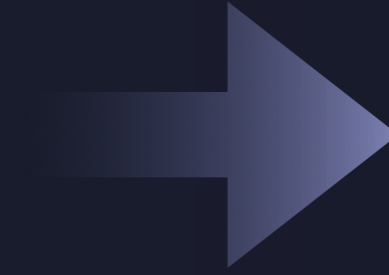
只需专注于数据分析



编写SQL



运行SQL



查看结果

数据湖分析 USQL

编辑SQL

已保存SQL

历史SQL

专家咨询

数据库

adx_offline

数据表

创建数据表

adx_ssp_req_log

bidid

sspip

delay

status

cd

ct

selfad

ssp_request

adx_pv_log

adx_click_log

adx_loss_log

adx_inmobi_log

adx_dsp_rsp_log

adx_twin_log

查询1 X

查询3 X

查询4 X

查询6 X

查询7 X

查询8 X

查询9 X

查询10 X

查询11 X

```

1 SELECT ssp_request.app.id AS appid,
2 COUNT(DISTINCT ssp_request.device.adid)-1 AS       _count,0 as       _count,
3 COUNT(DISTINCT ssp_request.device.ip)-1 AS ip_count,
4 COUNT(DISTINCT ssp_request.device.ip||ssp_request.device.ua)-1 AS       _count
5 FROM adx_ssp_req_log
6 WHERE year=2019 AND month=05 AND day=04 AND ssp_request.device.os = 0 AND selfad = 0
7 GROUP BY ssp_request.app.id
8 HAVING COUNT(DISTINCT ssp_request.device.ip||ssp_request.device.ua) > 1
9 OR COUNT(DISTINCT ssp_request.device.adid) > 1
10 OR COUNT(DISTINCT ssp_request.device.ip) > 1

```

立即运行

示例SQL

保存

清空

运行结果 已运行: 347.331秒, 已扫描数据量: 77.07 GB, 结果总行数:111

全屏查看结果

下载结果

	appid	 _count	 _count	 _count	 _count
1	153	47197	0	69297	77432

已运行: 347.331 秒, 已扫描数据量: 77.07 GB

77.07GB = 800G的JSON数据压缩

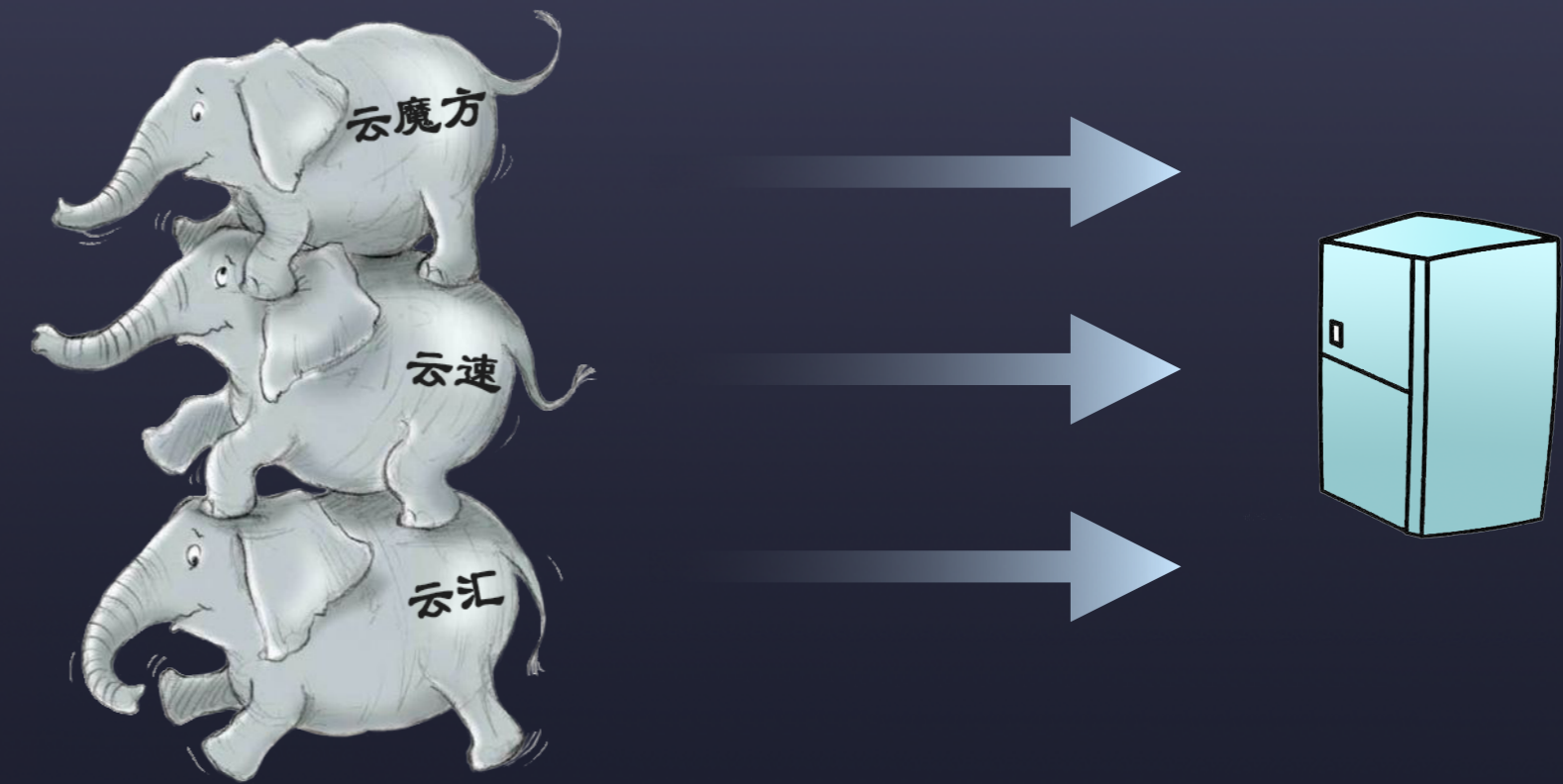
app 工厂
爱普新媒

“云魔方”互联网精准营销平台

“云速”企业移动信息服务平台

“云汇”互动广告平台

广告数据规模太过庞大、赛过大象，如何把这头大象“瘦身”装进“冰箱”真正使用起来这是爱普新媒迫切需要解决的问题



ADHOC



爱普新媒体——效果

研发人力投入

20人/日

0人/日

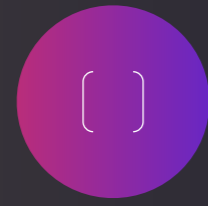
大数据月消费

降低
92.85%

需求完成周期

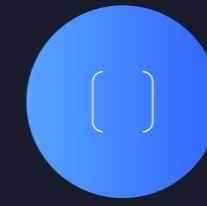
43.2小时

2小时



数据探索场景

- 数据格式不固定
- 数据规模大
- 没有元数据
- 分析目标不明确



USQL 数据湖分析

- 支持JSON, ORC, CSV
- 计算与存储分离
- 元数据与数据分离
- 快速SQL分析

高度匹配



数据探索——技术方案对比

方案名称	核心组件	数据规模	分析速度	计算成本	存储成本	运维要求	Schema后定义	数据格式支持
OLTP方案	MySQL PostgreSQL	百GB级	✓ 毫秒/秒级	中	高	中	不支持	固定
MPP方案	GreenPlum	百TB级	✓ 毫秒/秒级	高	高	高	不支持	少
HADOOP方案	Hive	✓ PB级	分钟级	高	中	高	✓ 支持	✓ 丰富
数据湖方案	USQL	✓ PB级	分钟级	✓ 极低	✓ 低	✓ 0	✓ 支持	✓ 丰富

定价

按数据扫描量计费： **0.03**元/GB

行式存储 vs 列式存储

客户ID	姓名	年龄	地址	城市	国家代码
318470519	张三	23	朝阳门8号	北京	CN
101258799	李四	35	江苏路153弄	上海	CN
4191301012	王五	19	南湖大道10号	深圳	CN

318470519|张三|23|朝阳门8号|北京|CN
Row 1

101258799|李四|35|江苏路153弄|上海|CN
Row 2

4191301012|王五|19|南湖大道10号|深圳|CN
Row 3

客户ID	姓名	年龄	地址	城市	国家代码
318470519	张三	23	朝阳门8号	北京	CN
101258799	李四	35	江苏路153弄	上海	CN
4191301012	王五	19	南湖大道10号	深圳	CN

318470519|101258799|4191301012
Column 1

张三|李四|王五
Column 2

23|35|19
Column 2

朝阳门8号|江苏路153弄|南湖大道10号
Column 4

北京|上海|深圳
Column 5

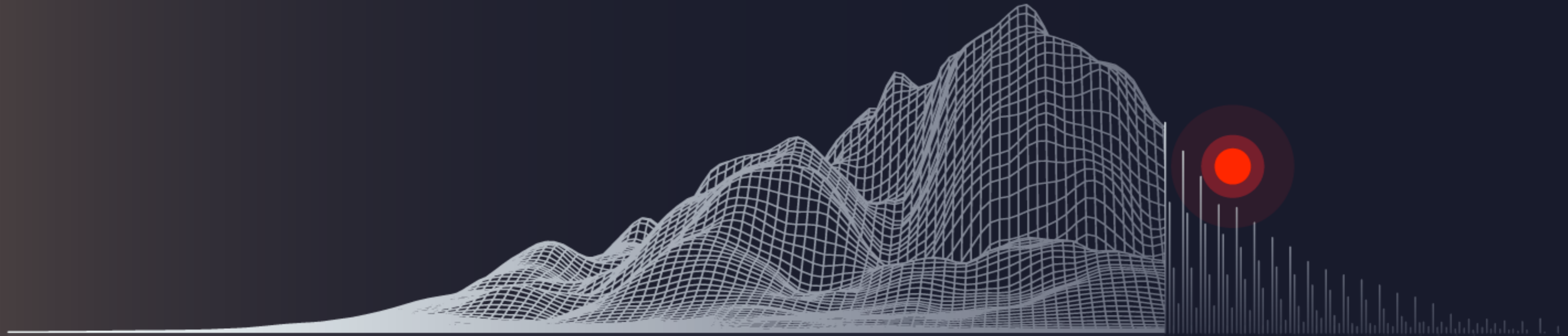
CN|CN|CN
Column 6

即将到来

支持JDBC驱动

数据加密 —— 协同UKMS (密钥管理服务)

CREATE TABLE AS SELECT (CTAS)



THANKS