# 移动数据收集及处理之迷思

## Mystery of Mobile Data Collecting & Processing

### 黄洋成

*Chief Architect@TalkingData*

# WHAT to collect

## DEVICE
IDentification

- IMEI?
- MAC?
- IDFA?

## SESSION
RECognition

## CONTEXTUAL
DATA

- Network
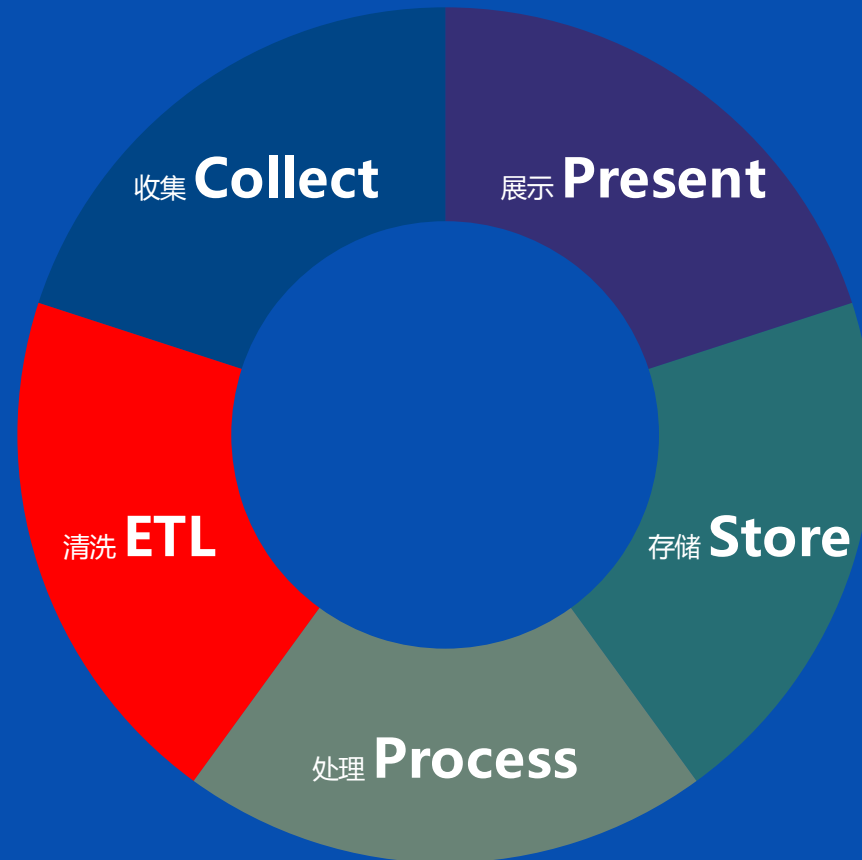- Location
- Sensors
- Voice/Image

# HOW to collect

**STORE** and
**F**orward

- Unreliable connection
- Unreliable execution

**CLOCK**
**I**ssues

**POWER**
**C**onsumption
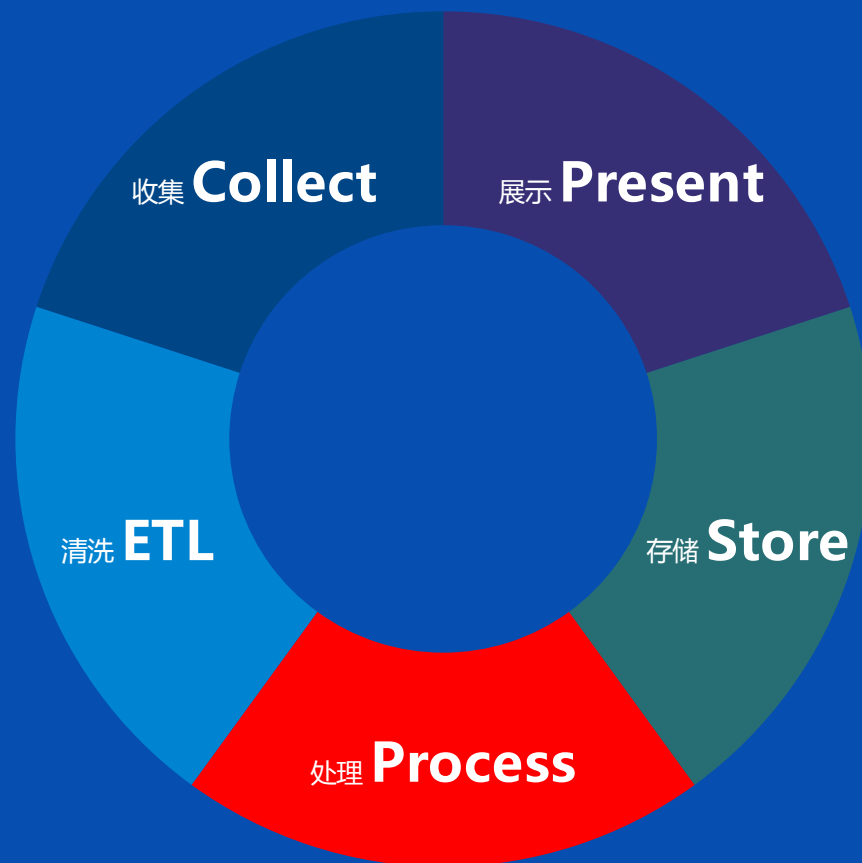
- Log vs structured
- Compress
- Keep-alive?

# ETL 数据清洗

## De-Duplication

## Matching

## Normalization

- Time
- Location
- Mobile Device
- Carrier

# REPORT REQuirement

| LATENCY | Comprehensive | Quick |
|---|---|---|
| | Report | Response |
| Action → Report in minutes | | |

# REPORT Processing

## FROM

- **Metrics** oriented
- Co-exist Stream/batch
- Pre-calc & index

## TO

- **Model**(cube) oriented
- Micro-batch
- Rollback support
- SQL-like query interface

# REPORT Achievements

## MODEL based
Processing

## QL/Script
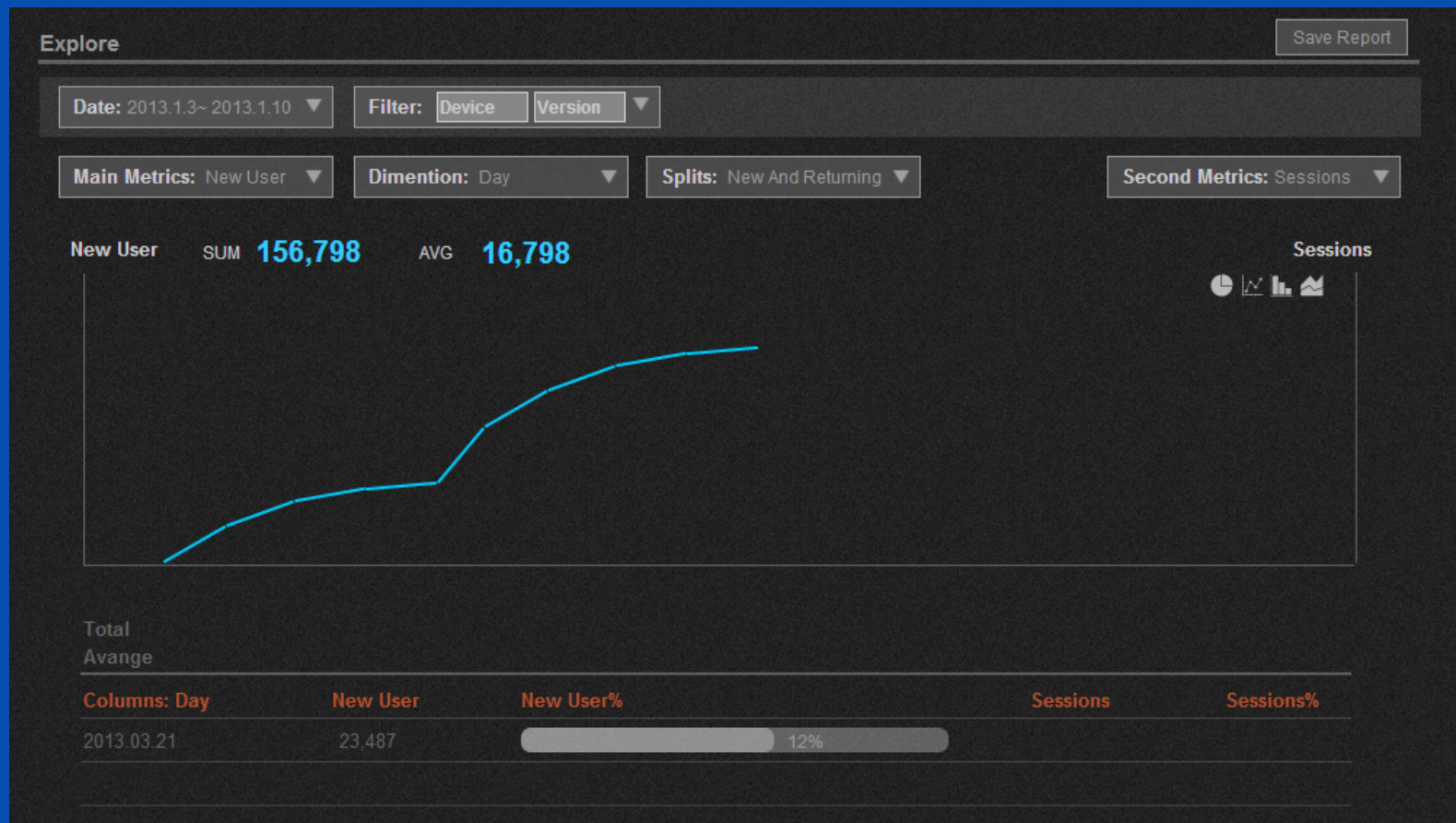Query Front

## Report
Tools

- ongoing

# Sample SCRIPT

留存 (Retention) 计算：

```
days = [20131023, 20131024, 20131025];
result = new HashMap<Integer, Integer>();
for(day in days){
    new = select user_id from new_user where product_id =
'1313930' and _day_of_rec_time = day;
    active = select user_id from active_user where product_id =
'1313930' and _day_of_rec_time = (day + 1);
    retention = intersection(new, active);
    result.put(day, retention.cardinality());
}
return result;
```

# Target Report

# 简单直接最有效
## Do it straightforward

关注视角

# DATA Mining WORKS

## Process

- Azkaban based workflow
- ETL based on Pig
- Machine Learning
  - Most on single machine
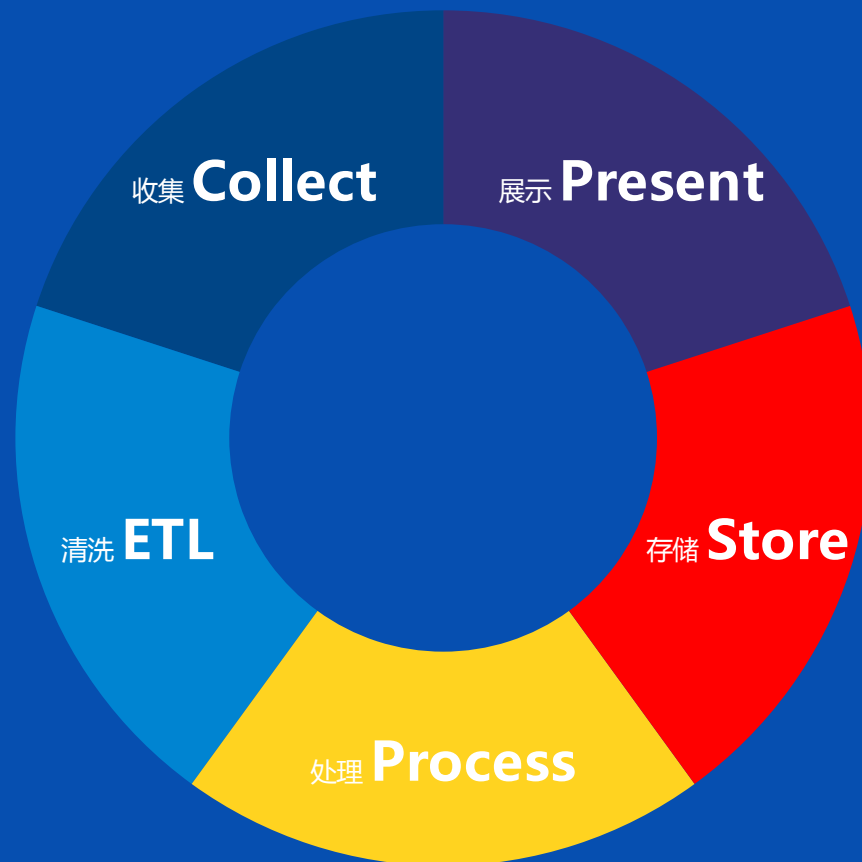  - Some on Spark

## Input

- Apps desc
- Device → Apps
  - Explore to download
  - Use
  - Pay

## Output

- ID mapping
- Tags for apps
- Tags for devices
- IP location

## Experiments

- App recommend
- Audience targeting Ads

# Key Components

## Logs

- HDFS
- RCFile/Snappy
- Archived as bzip

## META DATA

- Redis

## Report DATA

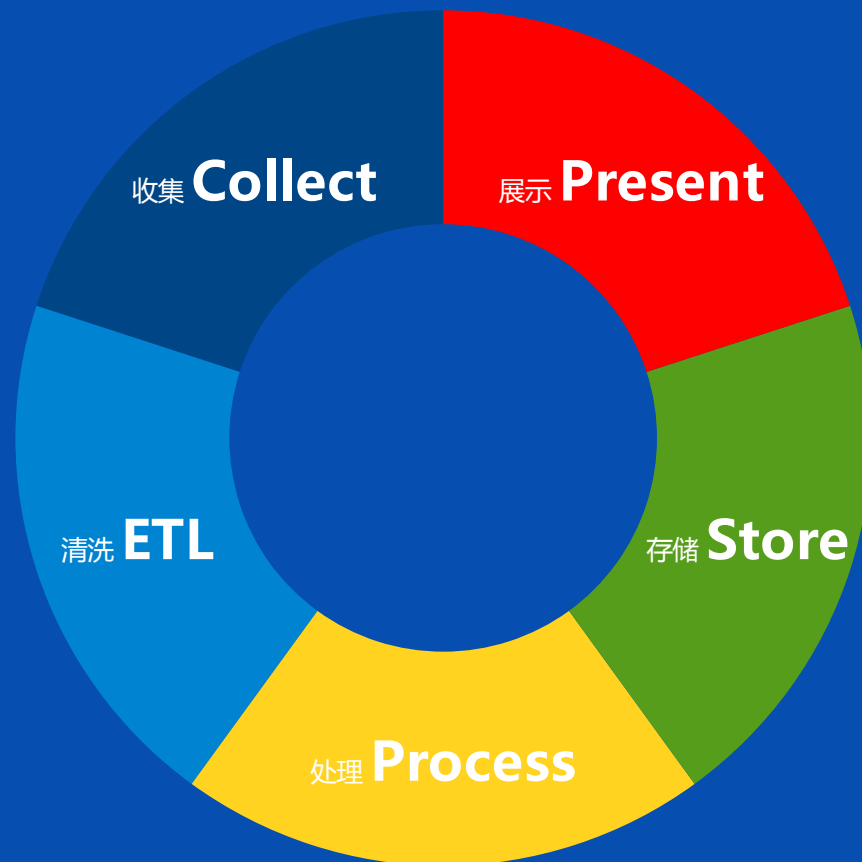- MySQL/MyISAM
- Toku

## NEAR Cache

- Off-heap cache
- Off-heap serialization lib
- Off-heap collections

# Challenges

**NETWORK**
DATA Flow Design

**Capacity**
Plan

**Partition**
Management

# PRESENTATION Challenges

## Security

- https
- Authorization check

## Complex Pages

- Components
- Rich Webapp
- Service Oriented

## Performance

- Parallel
  - Fork-join
  - Intelligent cache
  - Grid computing
  - GPU

# 谢谢！

Contact:
黄洋成 YC Huang
yc.huang@tendcloud.com
http://weibo.com/yc08
https://www.talkingdata.net/