# OAP: Optimized Analytics Package for Spark Platform

Daoyuan Wang (Intel)

Yuanjian Li (Baidu)

SPARK SUMMIT 2017

# Notice and Disclaimers:

SPARK
SUMMIT
2017

# About me

Daoyuan Wang

- developer@Intel
- Focuses on Spark optimization
- An active Spark contributor since 2014

Yuanjian Li

- Baidu INF distributed computation
- Apache Spark contributor
- Baidu Spark team leader

# Agenda

- Background for OAP
- Key features
- Benchmark
- OAP and Spark in Baidu
- Future plans

# Agenda

- Background for OAP
- Key features
- Benchmark
- OAP and Spark in Baidu
- Future plans

SPARK
SUMMIT
2017

# Data Analytics in Big Data Definition

Big Data: 4V definition

Volume

Velocity

Variety

Veracity

Copyright: infoDiagram.com 2014

- People wants OLAP against large dataset **as fast as possible**.

- People wants extract information from new coming data **as soon as possible**.

# Data Analytics Acceleration is Required by Spark Users

**FEATURES USERS CONSIDER IMPORTANT**

*Respondents were allowed to select more than one feature.*

**69%**
EASE OF DEPLOYMENT

**76%**
EASE OF PROGRAMMING

**51%**
REAL-TIME STREAMING

**91%**
PERFORMANCE

**82%**
ADVANCED ANALYTICS

SPARK SUMMIT 2017

http://cdn2.hubspot.net/hubfs/438089/DataBricks_Surveys_-_Content/2016_Spark_Survey/2016_Spark_Infographic.pdf

# Emerging hardware technology

Intel® Optane™ Technology Data Center Solutions

Accelerate applications for fast caching and storage, reduce transaction costs for latency-sensitive workloads and increase scale per server. Intel® Optane™ technology allows data centers to deploy bigger and more affordable datasets to gain new insights from large memory pools.

# Our proposal – OAP

**OAP (Codename "Spinach")**

| Spark* Job Server | • Auto tuning based on periodical job history<br>• K8S Integration / AES-NI Encryption |

| Spark SQL / Structured Streaming / Core | • Indexed Data Source / Cache Aware<br>• RDMA, QAT, ISA-L, FPGA … |

Hive* Table   Parquet *   JSON *   ORC *   Redis * Connector   Cassandra * Connector

• User Customized Indices
• Columnar formats & support Parquet, ORC
• Runtime Computing V.S. Data Store

Alluxio*   Redis*   Cassandra*   HBase*

• Columnar Fine-grained Cache
• Spark Executor in-process Cache
• 3D Xpoint (APP Direct Mode)

| HDFS*   S3*   … | Storage Layer |

# Why OAP

## Low cost

- Makes full use of existing hardware
- Open source

## Good Performance

- Index just like traditional database
- Up to 5x boost in real-world

## Easy to Use

- Easy to deploy
- Easy to maintain
- Easy to learn

SPARK SUMMIT 2017

# Agenda

- Background for OAP
- Key features
- Benchmark
- OAP and Spark in Baidu
- Future plans

SPARK
SUMMIT
2017

# A Simple Example

1. Run with OAP
   **$SPARK_HOME/sbin/start-thriftserver** --package oap.jar;
2. Create a OAP table
   **beeline> CREATE TABLE src(a: Int, b: String) USING** spn;
3. Create a single column B+ Tree index
   **beeline>** CREATE SINDEX **idx_1 ON src (a) USING BTREE;**
4. Insert data
   **beeline> INSERT INTO TABLE src SELECT key, value FROM xxx;**
5. Refresh index
   **beeline>** REFRESH SINDEX **on src;**
6. Execution would automatically utilize index
   **beeline> SELECT MAX(value), MIN(value) FROM src WHERE a > 100 and a < 1000;**

# OAP Files and Fibers

OAP meta file

Column (Fiber) #1
Column (Fiber) #2
⋮
Column (Fiber) #N
RowGroup #1

RowGroup #2

RowGroup #N

OAP data files

Index meta

statistics

Index data structure (Index Fiber)

OAP index files

One Index file for every data file

Index meta

statistics

Index data structure (Index Fiber)

OAP index files

SPARK SUMMIT 2017

# OAP Internals - index

Spark predicate push down → FilteredScan

FilteredScan → Read OAP Meta

**OAP cached access**

Read OAP Meta → Available index?

Available index? — Y → Index selection

Available index? — N → Full table scan

Index selection → read statistics before use index

read statistics before use index → Full table scan

read statistics before use index → Get Local RowID from index

Get Local RowID from index → Access data file for RowIDs directly

Supports Btree Index and BitMap Index, find best match among all created indices

Supports statistics such as MinMax, PartbyValue, Sample, BloomFilter

Only reads data fibers we need and puts those fibers into cache (in-memory fiber)

# OAP compatible layer



Cache

Read row #m from parquet file

In-memory fiber contains Row #m data

**Parquet compatible layer**

Find Row group #k

Read row group and get specific rows

RowGroup #1

RowGroup #2

RowGroup #k

Parquet data file

# OAP Data locality

# Agenda

- Background for OAP
- Key features
- Benchmark
- OAP and Spark in Baidu
- Future plans

# Performance

Cluster:
1 Master + 2 Slaves

Hardware:
CPU – 2x E5-2699 v4
RAM – 256 GB
Storage – S3610 1.6TB

Data:
300GB (Compressed Parquet)
2 Billion Records

## OAP Index And Cache Performance

Query Time (seconds)

| Parquet Vectorized Read | OAP Indexed Read | OAP Indexed Read with Fiber Cache |
|---|---|---|
| 72.083 | 7.095 | 2.304 |

# Agenda

- Background for OAP
- Key features
- Benchmark
- OAP and Spark in Baidu
- Future plans

SPARK
SUMMIT
2017

# Spark In Baidu



Chart legend: Nodes, Jobs/day

Data points:
- 2014: 80
- 2015: Nodes 1000, Jobs/day 300
- 2016: Nodes 3000, Jobs/day 1500
- 2017: Nodes 6500, Jobs/day 5800

Y-axis: 0, 1000, 2000, 3000, 4000, 5000, 6000, 7000

**2014**
- Spark import to Baidu
- Version: 0.8

**2015**
- Build standalone cluster
- Integrate with in-house FS\Pub-Sub\DW
- Version: 1.4

**2016**
- Build Cluster over YARN
- Integrate with in-house Resource Scheduler System
- Version: 1.6

**2017**
- SQL\Graph Service over Spark
- OAP
- Version: 2.1

SPARK SUMMIT 2017

Baidu INF

# Baidu Big SQL



**Baidu Big SQL**

| Web UI | Restful API |

**BBS HTTPServer**

**BBS Master**

| BBS Worker | BBS Worker | BBS Worker |

| Roll Up Table Layer | Cache & Index Layer(OAP) |

**Spark Over Yarn**

API Layer:
- Meta Control API
- Job API:
Load\Export\Query\Index Control

Control Layer:
- Meta Control
- Job Scheduler
- Spark Driver
- Query Classification

Boosting Layer:
- Roll Up Table Management
- Roll Up Query Change
- Index Create\Update
- Cache Hit

# Baidu Big SQL

# Introductory Story

# Introductory Story

Get the top 10 charge sum and correspond advertiser which triggered by the query word 'flower'

```
1  --- 鼠标移出输入框后，将自动检测可查询
2  select userid, sum(charge) as charge
3  from
4  where event_day=20170104
5  and query = '鲜花'
6  group by userid
7  order by charge desc
8  limit 10
```

- **Create index** on 'userid' column
- Various index types to choose for different fields types

- ×5 speed boosting than native spark sql, ×80 than MR Job
- 3 day baidu charging log, 4TB data,70000+ files, query time in 10~15s

# Roll Up Table Layer

700+ Columns

| date | userid | searchid | baiduid | cmatch | | shows | clicks | charge |
|------|--------|----------|---------|--------|-----|-------|--------|--------|
| 1 | 1 | 1 | 10 | 2 | | 10 | 1 | 5 |
| 1 | 1 | 2 | 11 | 3 | | 10 | 1 | 5 |
| 1 | 1 | 3 | 12 | 2 | | 10 | 1 | 5 |
| 1 | 1 | 4 | 13 | 1 | | 10 | 1 | 5 |
| 1 | 1 | 5 | 14 | 1 | ... | 10 | 1 | 5 |
| 1 | 2 | 6 | 14 | 2 | ... | 10 | 1 | 5 |
| 1 | 2 | 7 | 15 | 3 | | 10 | 1 | 5 |
| 1 | 2 | 8 | 16 | 4 | | 10 | 1 | 5 |
| 1 | 2 | 9 | 17 | 5 | | 10 | 1 | 5 |

Select date,userid,shows,clicks,charge from...

99% query only use <10 columns

## Multi Roll Up Table (user-transparent)

| date | userid | shows | clicks | charge |
|------|--------|-------|--------|--------|
| 1 | 1 | 50 | 5 | 25 |
| 1 | 2 | 40 | 4 | 20 |

| date | cmatch | shows | clicks | charge |
|------|--------|-------|--------|--------|
| 1 | 1 | 20 | 2 | 10 |
| 1 | 2 | 30 | 3 | 15 |
| 1 | 3 | 20 | 2 | 10 |
| 1 | 4 | 10 | 1 | 5 |
| 1 | 5 | 10 | 1 | 5 |

Bai du INF

# OAP In BigSQL

| ... | Name | Department | Age | ... |
|---|---|---|---|---|
| ... | ... | ... | ... | ... |
| ... | John | INF | 35 | ... |
| ... | Michelle | AI-Lab | 29 | ... |
| ... | Amy | INF | 42 | ... |
| ... | Kim | AI-Lab | 27 | ... |
| ... | Mary | AI-Lab | 47 | ... |
| ... | ... | ... | ... | ... |

Normal Table Scan

Data File

Use Index

Skippable Reader

Index Build

| Sorted Age | Row Index in Data File |
|---|---|
| 27 | 3 |
| 29 | 1 |
| 35 | 0 |
| 42 | 2 |
| 45 | 4 |

| Department | Bit Array |
|---|---|
| INF | 10100 |
| AI-Lab | 01011 |

Index File

Select xxx from xxx where age > 29 and department in (INF, AI-Lab)

# OAP In BigSQL

| ... | Name | Department | Age | ... |
|-----|------|------------|-----|-----|
| ... | ... | ... | ... | ... |
| ... | John | INF | 35 | ... |
| ... | Michelle | AI-Lab | 29 | ... |
| ... | Amy | INF | 42 | ... |
| ... | Kim | AI-Lab | 27 | ... |
| ... | Mary | AI-Lab | 47 | ... |
| ... | ... | ... | ... | ... |

Data File

Load Cache

| Department | Row Index in Data File |
|------------|------------------------|
| INF | 2 |
| AI-Lab | 3 |

| Age | Row Index in Data File |
|-----|------------------------|
| 35 | 0 |
| 29 | 1 |

In Memory Cache

# BBS's Contribute to Spark

- ## Spark-4502

  Spark SQL reads unneccesary nested fields from Parquet

- ## Spark-18700

  getCached in HiveMetastoreCatalog not thread safe cause driver OOM

- ## Spark-20408

  Get glob path in parallel to reduce resolve relation time

- ## …

# Agenda

- Background for OAP
- Key features
- Benchmark
- OAP and Spark in Baidu
- Future plans

# Future plans

- Compatible with more data formats
- Explicit cache and cache management
- Optimize SQL operators (join, aggregate) with index
- Integrate with structured streaming
- Utilize Latest hardware technology, such as Intel QAT or 3D XPoint.
- Welcome to contribute!

https://github.com/Intel-bigdata/OAP

# WOMEN IN BIG DATA NETWORKING LUNCHEON

The Women in Big Data team invites you to join us for lunch, network with your peers and hear from a dynamic panel of experts. Come learn what career & growth opportunities are available in the field of big data analytics.

## Agenda:

| | |
|---|---|
| 12.20PM | Grab Lunch & Networking |
| 12:30PM-12:40PM | Women in Big Data Overview with Soumya Guptha, Marketing Manager, Intel |
| 12:40PM-12:45PM | My journey in Data Analytics & Artificial Intelligence with Ziya Ma, Intel VP & Director, Big Data Technologies |
| 12:50PM-01:40PM | Panel: Making The Best Out Of The Fast Paced Data World! |

## Panel: Making The Best Out Of The Fast Paced Data World!

Gayle Sheppard, VP, New Technology Group, Intel | Ritika Gunnar, Global VP of IBM Cloud and Cognitive, IBM | Eva Tse, Director of Big Data Services, Netflix | Jennifer Shin, CEO 8 path solutions | Soumya Guptha, Marketing Manager, Software and Solutions Group, Intel

Join us for a networking luncheon to hear from industry experts from leading companies such as IBM, Intel and others on their investments in Big Data technologies such as Spark, Machine Learning, Artificial Intelligence.

www.womeninbigdata.org/ | @DataWomen | Women in Big Data Forum | www.meetup.com/Women-in-Big-Data-Meetup/

**GRASSROOTS COMMUNITY CHAMPIONING WOMEN'S LEADERSHIP AND SUCCESS IN BIG DATA**

# Thank You.

daoyuan.wang@intel.com

liyuanjian@baidu.com