CarbonData: A New Hadoop File Format For Faster Data Analysis

www.huawei.com

陈亮

Email: chenliang613@huawei.com

Weichat:chenliang2007



HUAWEI TECHNOLOGIES CO., LTD.

Outline

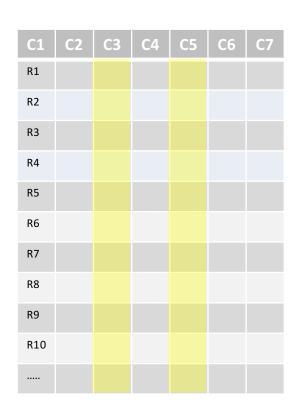
- Use Case & Motivation : Why introducing a new file format?
- CarbonData File Format Deep Dive
- Framework Integrated with CarbonData
- Demo & Performance Comparison
- Future Plan



Use case: Sequential scan

- Full table scan
 - Big scan(all rows, no filter)
 - Only fetch a few columns of the table

- Common usage scenario:
 - ETL job
 - Log Analysis

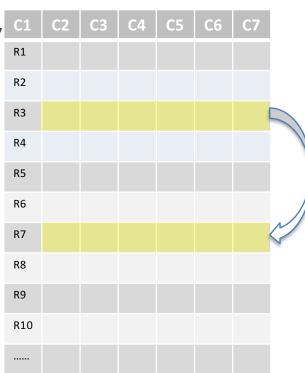




Use case: Random Access

- Predicate filtering on many columns(point query
 - Row-key query (like HBase)
 - Narrow scan but might fetch all columns
 - Requires second/sub-second level low-latency

- Common usage scenario:
 - Operational query
 - User profiling

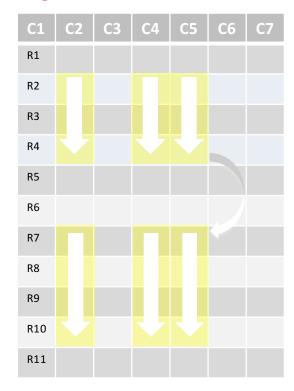




Use case: OLAP-Style Query

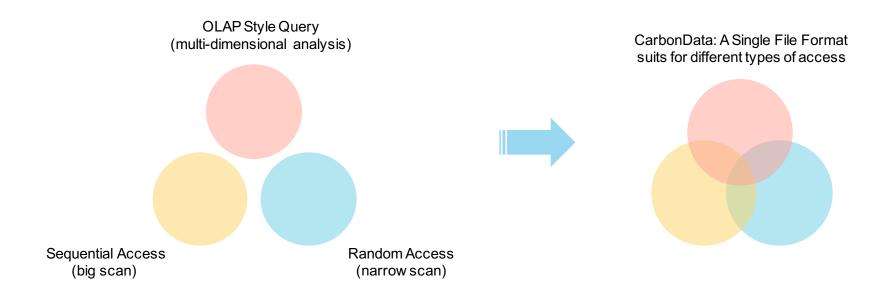
- Interactive data analysis for any dimensions
 - Involves aggregation / join
 - Roll-up, Drill-down, Slicing and Dicing
 - Low-latency ad-hoc query

- Common usage scenario:
 - Dash-board reporting
 - Fraud & Ad-hoc Analysis





Motivation





Why CarbonData

Based on the below requirements, we investigated existing file formats in the Hadoop eco-system, but we could not find a suitable solution that can satisfy all the requirements at the same time, so we start designing Carbon Data.

- •Support big scan & only fetch a few columns
- Support primary key lookup response in sub-second.
- Support interactive OLAP-style query over big data which involve many filters in a query, this type of workload should response in seconds.
- Support fast individual record extraction which fetch all columns of the record.
- Support HDFS so that customer can leverage existing Hadoop cluster.

When we investigated Parquet/ORC, it seems they work very well for R1 and R5, but they does not meet for R2,R3,R4. So we designed CarbonData mainly to add following differentiating features:

- •Stores data along with index: it can significantly accelerate query performance and reduces the I/O scans and CPU resources, where there are filters in the query. CarbonData index consists of multiple level of indices, a processing framework can leverage this index to reduce the task it needs to schedule and process, and it can also do skip scan in more finer grain unit (called blocklet) in task side scanning instead of scanning the whole file.
- •Operable encoded data: Through supporting efficient compression and global encoding schemes, can query on compressed/encoded data, the data can be converted just before returning the results to the users, which is "late materialized".
- •Column group: Allow multiple columns to form a column group that would be stored as row format. This reduces the row reconstruction cost at query time.
- •Supports for various use cases with one single Data format: like interactive OLAP-style query, Sequential Access (big scan), Random Access (narrow scan).

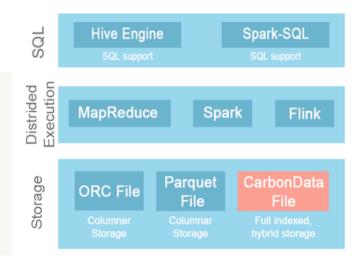


Design Goals

- Low-Latency for various types of data access pattern
- Allow fast query on fast data
- Ensure Space Efficiency
- General format available on Hadoop-ecosystem



- Read-optimized columnar storage
- Leveraging multi-level Index for low-latency
- Support column group to leverage the benefit of row-based
- Enables dictionary encoding for deferred decoding for aggregation
- Broader Integration across Hadoop-ecosystem





Outline

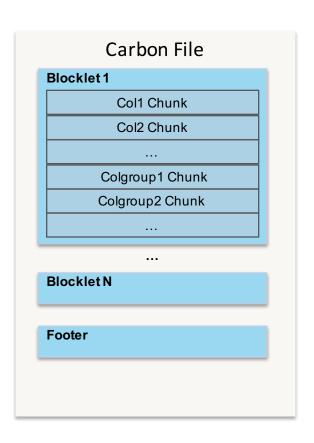
- Use cases & Motivation: Why introducing a new file format?
- CarbonData File Format Deep Dive
- Framework Integrated with CarbonData
- Demo & Performance Comparison
- Future Plan



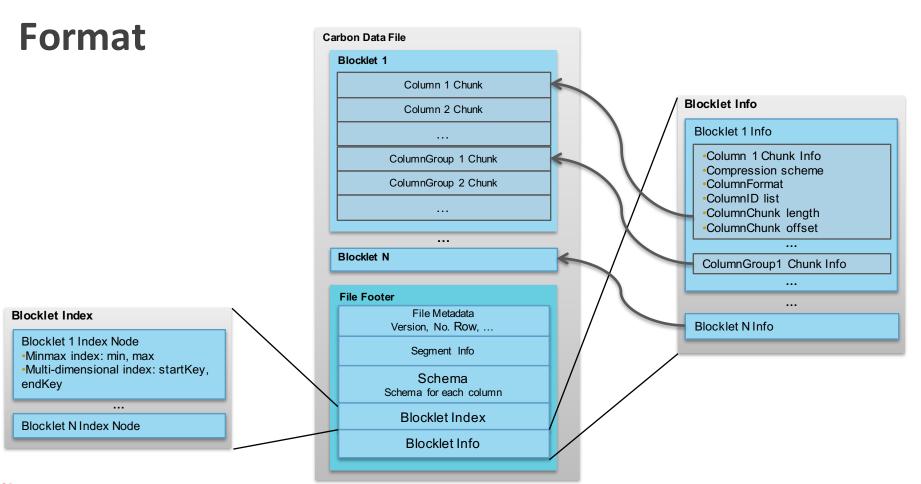
CarbonData File Structure

- Blocklet: A set of rows in columnar format
- Column chunk : Data for one column/column group in a Blocklet
 - Allow multiple columns forms a column group & stored as row-based
 - Column data stored as sorted index
- Footer: Metadata information
 - File level metadata & statistics
 - Schema
 - Blocklet Index & Blocklet level Metadata

Remark: One CarbonData file is a HDFS block.







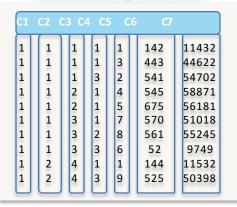


Blocklet

- Data are sorted along MDK (multi-dimensional keys)
 - data stored as index in columnar format

Years	Quarters	Months	Territory	Country	Quantity	Sales
2003	QTR1	Jan	EMEA	Germany	142	11,432
2003	QTR1	Jan	APAC	China	541	54,702
2003	QTR1	Jan	EMEA	Spain	443	44,622
2003	QTR1	Feb	EMEA	Denmark	545	58,871
2003	QTR1	Feb	EMEA	Italy	675	56,181
2003	QTR1	Mar	APAC	India	52	9,749
2003	QTR1	Mar	EMEA	UK	570	51,018
2003	QTR1	Mar	Japan	Japan	561	55,245
2003	QTR2	Apr	APAC	Australia	525	50,398
2003	QTR2	Apr	EMEA	Germany	144	11,532

Blocklet Logical View





Sorted MDK Index

[1,1,1,1,1]: [142,11432] [1,1,1,1,3]: [443,44622] [1,1,1,3,2]: [541,54702] [1,1,2,1,4]: [545,58871] [1,1,2,1,5]: [675,56181] [1,1,3,1,7]: [570,51018] [1,1,3,2,8]: [561,55245] [1,1,3,3,6]: [52,9749] [1,2,4,1,1]: [144,11532] [1,2,4,3,9]: [525,50398]

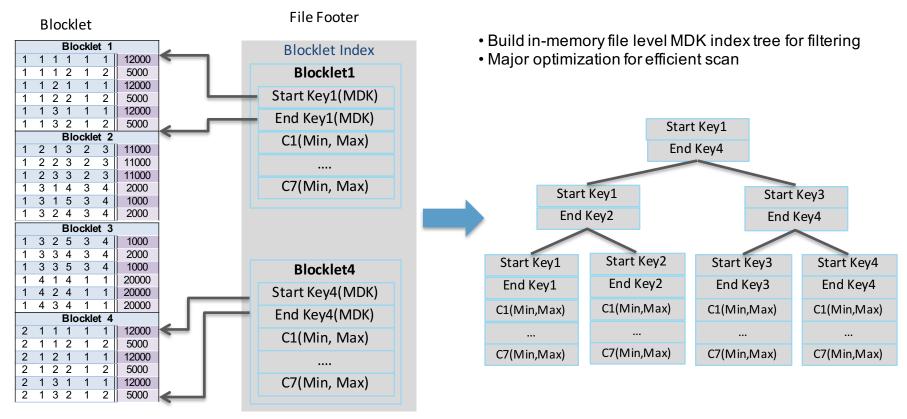


Encoding

[1,1,1,1,1]: [142,11432] [1,1,1,3,2]: [541,54702] [1,1,1,1,3]: [443,44622] [1,1,2,1,4]: [545,58871] [1,1,2,1,5]: [675,56181] [1,1,3,3,6]: [52,9749] [1,1,3,1,7]: [570,51018] [1,1,3,2,8]: [561,55245] [1,2,4,3,9]: [525,50398] [1,2,4,1,1]: [144,11532]



File Level Blocklet Index

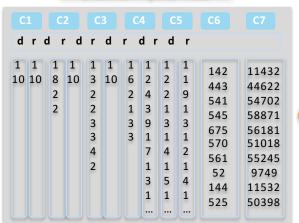




Inverted Index

- Optionally store column data as inverted index within column chunk
 - Very good benefit for low cardinality column for Better compression
 - Fast predicate filtering

Blocklet Physical View



Column chunk Level inverted Index

Blocklet (sort column within column chunk)

[1 1]	:[1 1]	:[1 1]	:[1 1]	:[1 1]	: [142]:[11432]
[1 2]	:[1 2]	:[1 2]	:[1 2]	:[1 9]	: [443]:[44622]
[1 3]	:[1 3]	:[1 3]	:[1 4]	:[2 3]	: [541]:[54702]
[1 4]	:[1 4]	:[2 4]	:[1 5]	:[3 2]	: [545]:[58871]
	:[1 5]				: [675]:[56181]
[1 6]	:[1 6]	:[3 6]	:[1 9]	:[5 5]	: [570]:[51018]
[1 7]	:[1 7]	:[3 7]	:[2 7]	:[6 8]	: [561]:[55245]
	:[1 8]				: [52]:[9749]
	:[2 9]				:[144]:[11532]
[1 10	:[2 10]:[4 1)]:[3 1	0] :[9 10]	: [525]:[50398]



Run Length Encoding & Compression

	Columnar Store	Blockl	et Rows		Measure1 Measure2
Dim1 Block 1(1-10)	Dim2 Block 1(1-8) 2(9-10)	Dim3 Block 1(1-3) 2(4-5) 3(6-8) 4(9-10)	Dim4 Block 1(1-2,4-6,9) 2(7) 3(3,8,10)	Dim5 Block 1(1,9) 2(3) 3(2) 4(4) 5(5) 6(8) 7(6) 8(7) 9(10)	Block Block [142]:[11432] [443]:[44622] [541]:[54702] [545]:[56181] [570]:[51018] [561]:[55245] [52]:[9749] [144]:[11532] [525]:[50398]



Column Group

- Allow multiple columns form a column group
 - stored as a single column chunk in rowbased format
 - suitable to set of columns frequently fetched together
 - saving stitching cost for reconstructing row

Blocklet 1					
C1	C2	C3	C4	C5	C6
Col Chunk	Col Chunk	Col Chunk	Co Ch	l unk	Col Chunk
10	2	23	23	38	15.2
10	2	50	15	29	18.5
10	3	51	18	52	22.8
11	6	60	29	16	32.9
12	8	68	32	18	21.6



Nested Data Type Representation

Arrays

- Represented as a composite of two columns
- One column for the element value

Array<Ph_Number>

[192,191]

[198,787]

[121,345,333]

One column for start_index & length of Array

Name	Array [start,len]	Ph_Number
John	0,2	192
Sam	2,3	191
Bob	5,2	121
		345
		333
		198
		787

Struts

- Represented as a composite of finite number of columns
- Each struct element is a separate column

Name	Info Strut <age,gender></age,gender>	
John	[31,M]	
Sam	[45,F]	
Bob	[16,M]	

	Name	Info.age	Info.gender
	John	31	М
	Sam	45	F
	Bob	16	М



Name

John

Sam

Bob

Encoding & Compression

- Efficient encoding scheme supported:
 - DELTA, RLE, BIT_PACKED
 - Dictionary: table level global dictionary
- Compression Scheme: Snappy

Compression ratio: 1/3



Big Win:

- Speedup Aggregation
- •Reduce run-time memory footprint
- Enable deferred decoding
- Enable fast distinct count

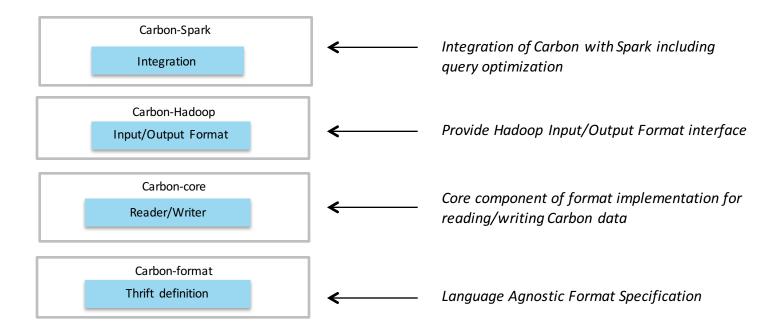


Outline

- Use Case & Motivation: Why introducing a new file format?
- CarbonData File Format Deep Dive
- Framework Integrated with CarbonData
- Demo & Performance Comparison
- Future Plan

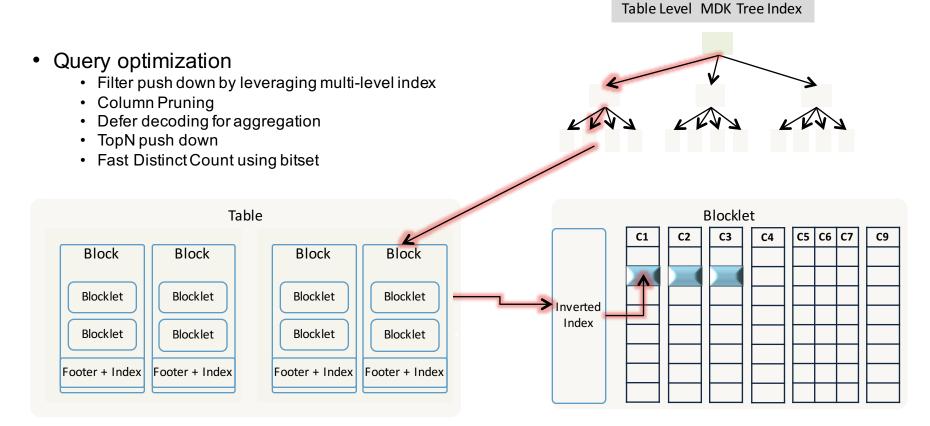


CarbonData Modules(to understand more)





Spark Integration





Spark Integration

- Query CarbonData Table
 - DataFrame API
 - Spark SQL Statement

```
CREATE TABLE [IF NOT EXISTS] [db_name.]table_name [(col_name data_type [COMMENT col_comment], ...)] [COMMENT table_comment] [PARTITIONED BY (col_name data_type [COMMENT col_comment], ...)] STORED BY 'org.carbondata.hive.CarbonHanlder' [TBLPROPERTIES (property_name=property_value, ...)] [AS select_statement];
```



Outline

- Use Case & Motivation: Why introducing a new file format?
- CarbonData File Format Deep Dive
- Framework Integrated with CarbonData
- Demo & Performance Comparison
- Future Plan



DEMO and Performance Comparison

数据量超过1亿行记录,性能比较会更明显。

因个人便携配置有限,本例子只以100万行纪录为demo演示:

1.Spark SQL里查询csv格式数据性能:

```
import org.apache.spark.sql.catalyst.util._
import org.apache.spark.sql.SQLContext
var df = sqlContext.read.format("com.databricks.spark.csv").option("header", "true").option("inferSchema", "true").load("./carbondata/hzmeetup.csv")
benchmark { df.filter($"country" == "china" and $"name" == "hangzhou" and $"seg" < 100000).count }
2.Spark SQL里查询parquet格式数据性能:
import org.apache.spark.sql.{CarbonContext, DataFrame, Row, SaveMode, SQLContext}
df.write.mode(SaveMode.Overwrite).parquet("./carbondata/parquet")
val parquetdf = sqlContext.parquetFile("./carbondata/parquet")
benchmark { parquetdf.filter($"country" === "china" and $"name" === "hangzhou" and $"seg" < 100000).count }
3. Spark SQL里查询CarbonData格式数据性能:
请参照CarbonDatagithub例子,初始化CarbonContext。
cc.sql("CREATE TABLE meetupTable (seq Int, name String, country String,age Int) STORED BY 'org.apache.carbondata.format'")
cc.sql("LOAD DATA LOCAL INPATH './carbondata/hzmeetup.csv' INTO TABLE meetupTable")
val carbondf = cc.read.format("org.apache.spark.sgl.CarbonSource").option("tableName", "meetupTable").load()
benchmark { carbondf.filter($"country" === "china" and $"name" === "hangzhou" and $"seq" < 100000).count }
```

注:以上代码因拷贝原因,双引号和单引号可能被转了格式,在Spark shell里需要修正。具体步骤,请参看github中wiki的quick start



Outline

- Motivation: Why introducing a new file format?
- CarbonData File Format Deep Dive
- Framework Integrated with CarbonData
- Demo & Performance Comparison
- Future Plan



Future Plan

- Upgrade to Spark 2.0
- Integrate with BI tools
- Add append support
- Support pre-aggregated table
- Broader Integration across Hadoop-ecosystem



Community

- CarbonData is open sourced & will become Apache Incubator project
- Welcome contribution to our Github @: https://github.com/HuaweiBigData/carbondata

- Main Contributors:
 - Jihong MA, Vimal, Raghu, Ramana, Ravindra, Vishal, Aniket, Liang Chenliang, Jacky Likun, Jarry Qiuheng, David Caiqiang, Eason Linyixin, Ashok, Sujith, Manish, Manohar, Shahid, Ravikiran, Naresh, Krishna, Babu, Ayush, Santosh, Zhangshunyu, Liujunjie, Zhujing (Huawei)
 - Jean-Baptiste Onofre (Talend, ASF member), Henry Saputra (eBay, ASF member),
 Uma Maheswara Rao G(Intel, Hadoop PMC)



Thank you

www.huawei.com

Copyright@2014 Huawei Technologies Co., Ltd. All Rights Reserved.

The information in this document may contain predictive statements including, without limitation, statements regarding the future financial and operating results, future product portfolio, new technology, etc. There are a number of factors that could cause actual results and developments to differ materially from those expressed or implied in the predictive statements. Therefore, such information is provided for reference purpose only and constitutes neither an offer nor an acceptance. Huawei may change the information at any time without notice.