

Как класть Parquet

Федор Лаврентьев @ Front Tier

t.me/fediq

План

- Что за паркет?
- Зачем всё это?
- Сильные и слабые стороны
- Что можно покрутить
- Как делать нельзя
- А есть что-нибудь получше?

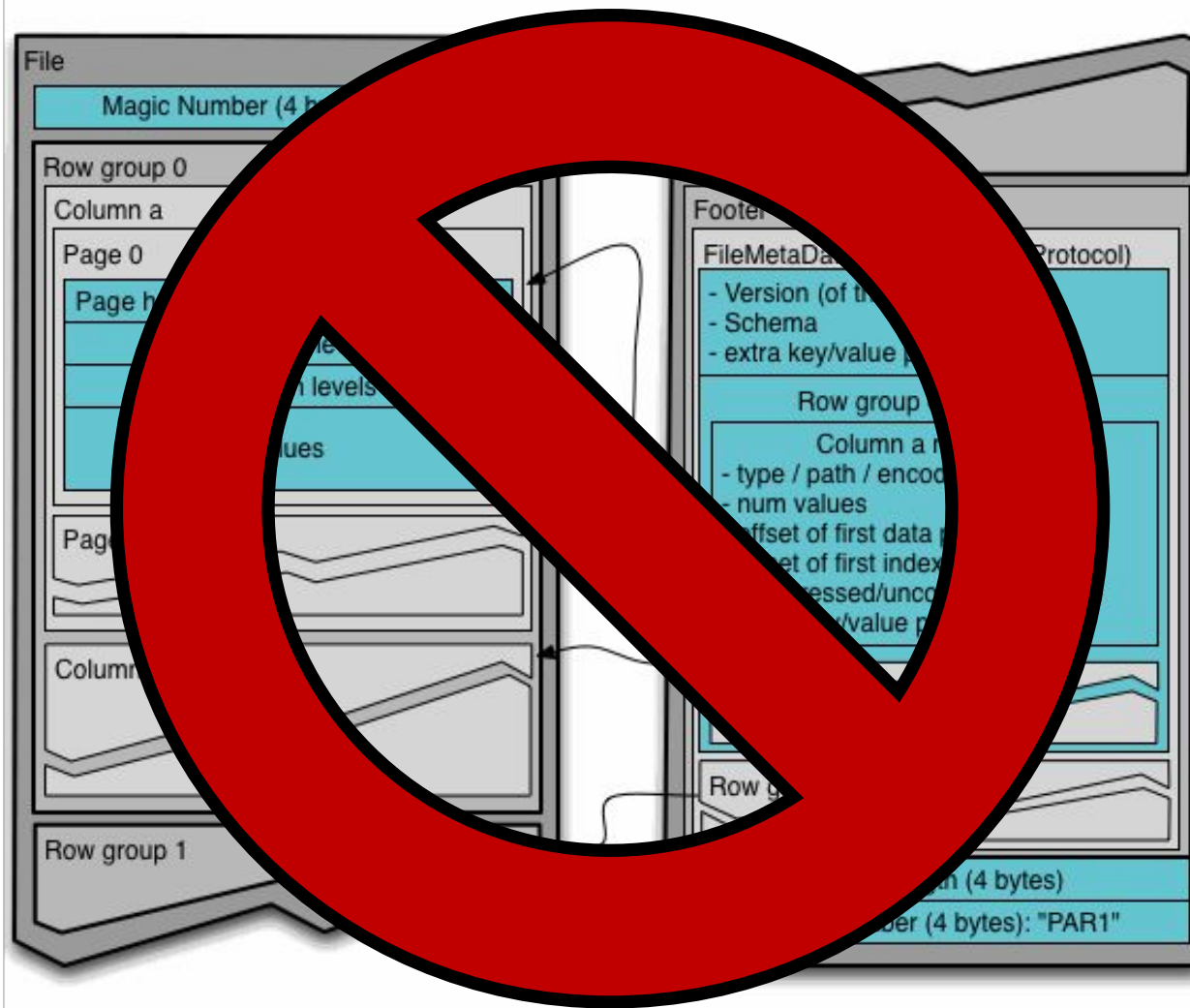
Apache Parquet

is a **columnar storage format**

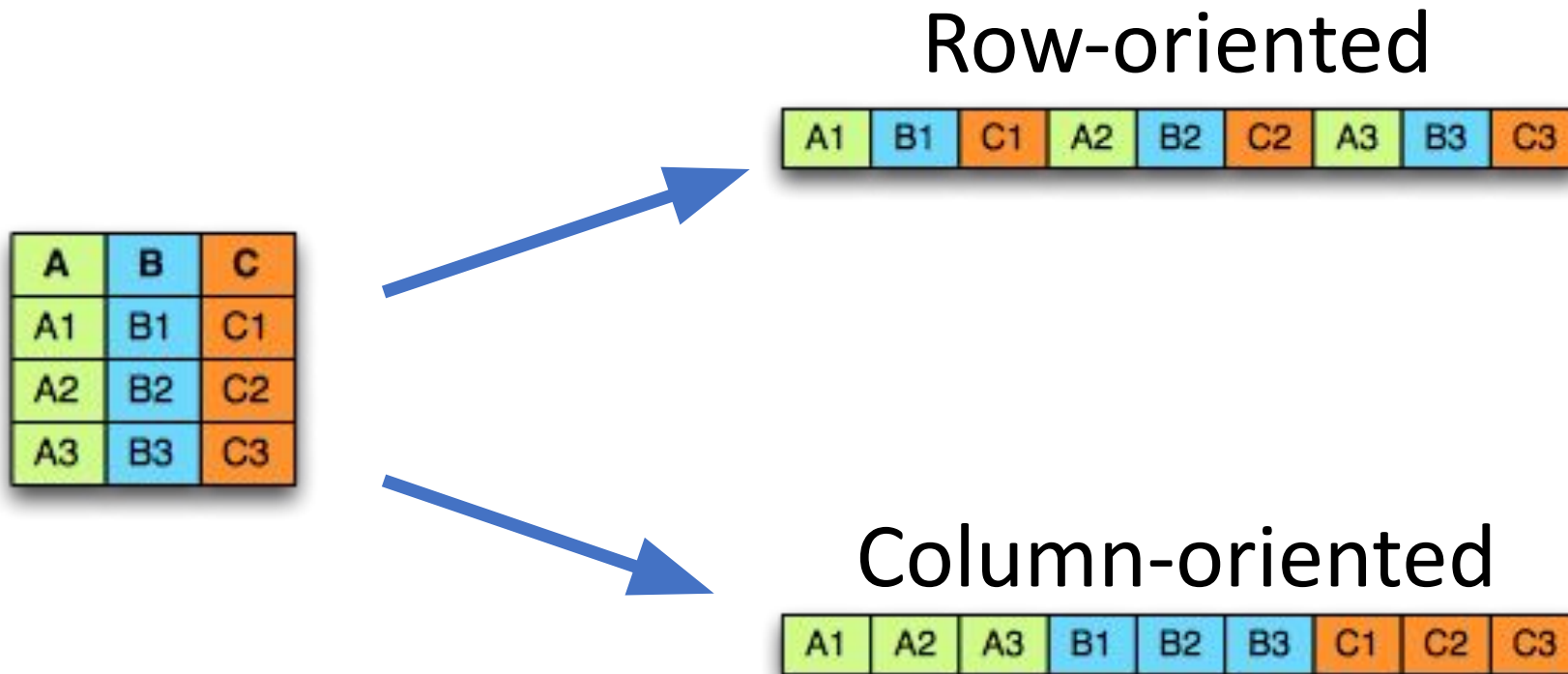
available to any project in the Hadoop ecosystem

regardless of the choice of data processing framework,
data model or programming language.

Как устроен



Колоночное хранение



Вложенные структуры

```
message Document {
  required int64 DocId;
  optional group Links {
    repeated int64 Backward;
    repeated int64 Forward; }
  repeated group Name {
    repeated group Language {
      required string Code;
      optional string Country; }
    optional string Url; }}
```

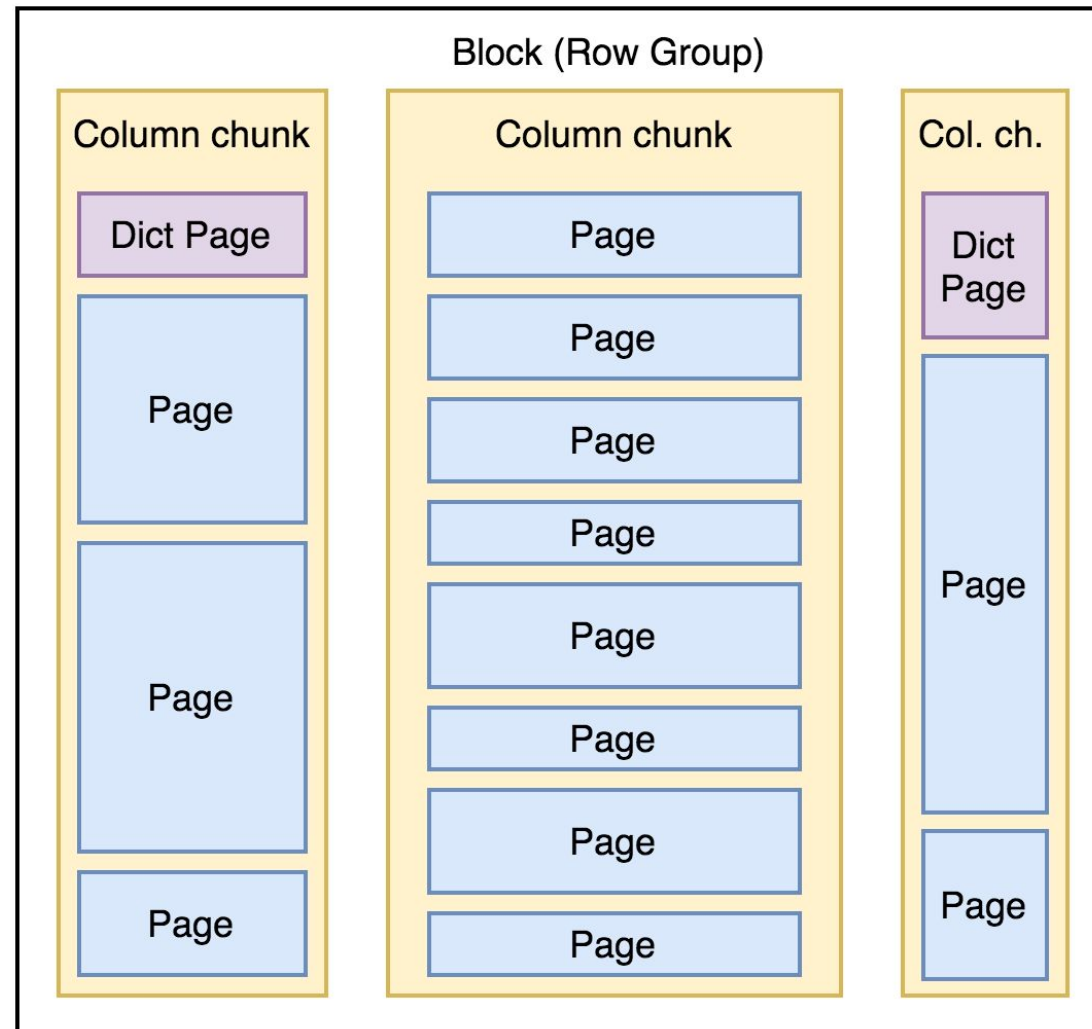
DocId: 10	r₁
Links	
Forward: 20	
Forward: 40	
Forward: 60	
Name	
Language	
Code: 'en-us'	
Country: 'us'	
Language	
Code: 'en'	
Url: 'http://A'	
Name	
Url: 'http://B'	
Name	
Language	
Code: 'en-gb'	
Country: 'gb'	

DocId: 20	r₂
Links	
Backward: 10	
Backward: 30	
Forward: 80	
Name	
Url: 'http://C'	

DocId	Name.Url	Links.Forward	Links.Backward					
value	r	d	value	r	d	value	r	d
10	0	0	20	0	2	NULL	0	1
20	0	0	40	1	2	10	0	2
			60	1	2	30	1	2
			80	0	2			

Name.Language.Code	Name.Language.Country				
value	r	d	value	r	d
en-us	0	2	us	0	3
en	2	2	NULL	2	2
NULL	1	1	NULL	1	1
en-gb	1	2	gb	1	3
NULL	0	1	NULL	0	1

Колоночное хранение в Parquet



И что?

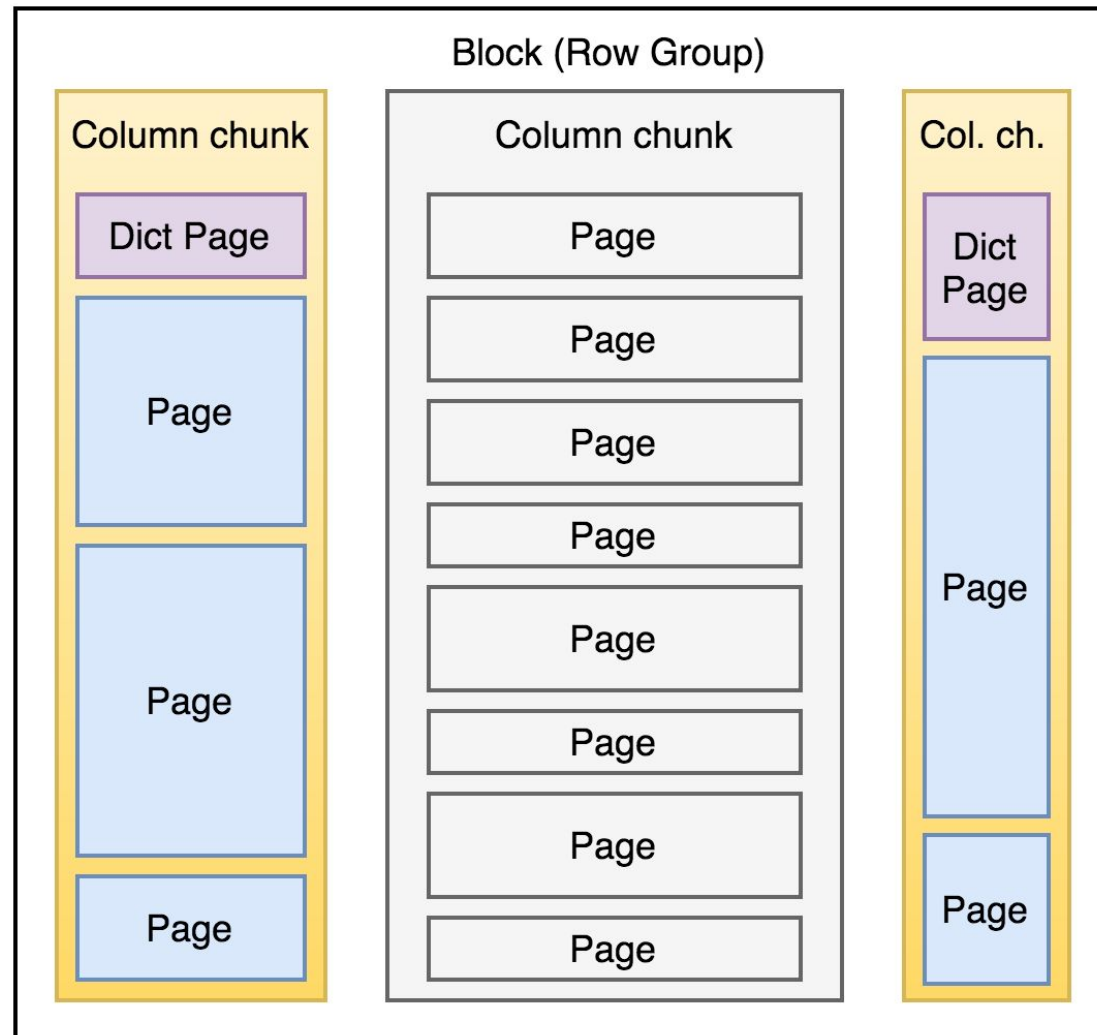
Pages МОЖНО кодировать!

- Bit packing
- Run-length encoding (RLE)
- Delta encoding
- Dictionary encoding
- ТОДО картинки

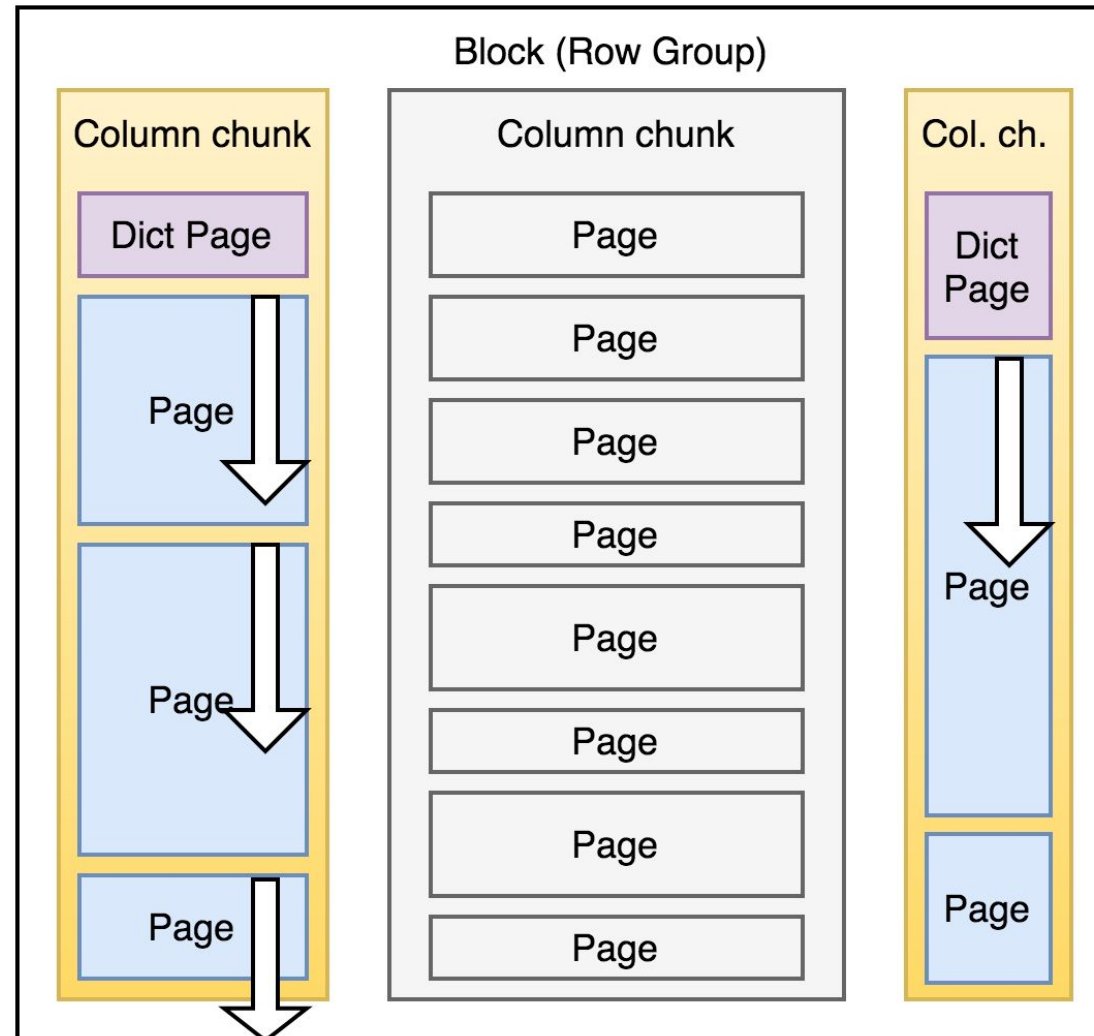
Кодированные Pages МОЖНО сжимать!

- GZIP
- LZO
- Snappy

Ненужные колонки можно не читать!



Нужные колонки МОЖНО ЧИТАТЬ параллельно!



Срочно в продакшен!..

Выбрасываем легаси!..

- Postgres VS Parquet
- Cassandra VS Parquet
- MongoDB VS Parquet
- Excel VS Parquet

Write once!

- Append only
- Read only
- Только батчевая запись
- Нет транзакций. Никак. Совсем.

Подходит для аналитики

- Только чтение
- Большие range scan'ы
- Сложные фильтры
- Группирующие запросы

Большие Range Scan'ы

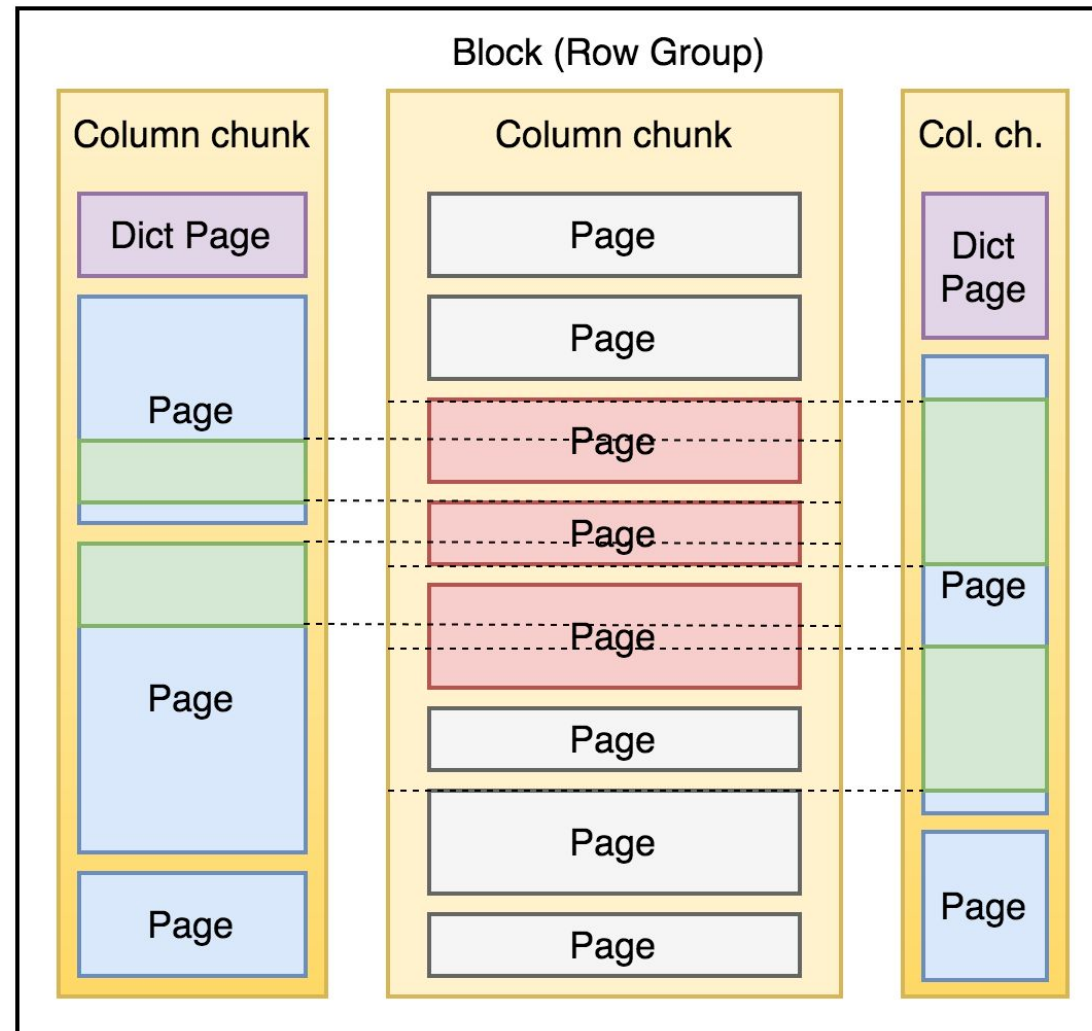
⇒ Партиционируйте данные!

/dataset/2018/07/30/

/dataset/2018/08/01/

/dataset/2018/08/02/

Сложные фильтры



Predicate Pushdown

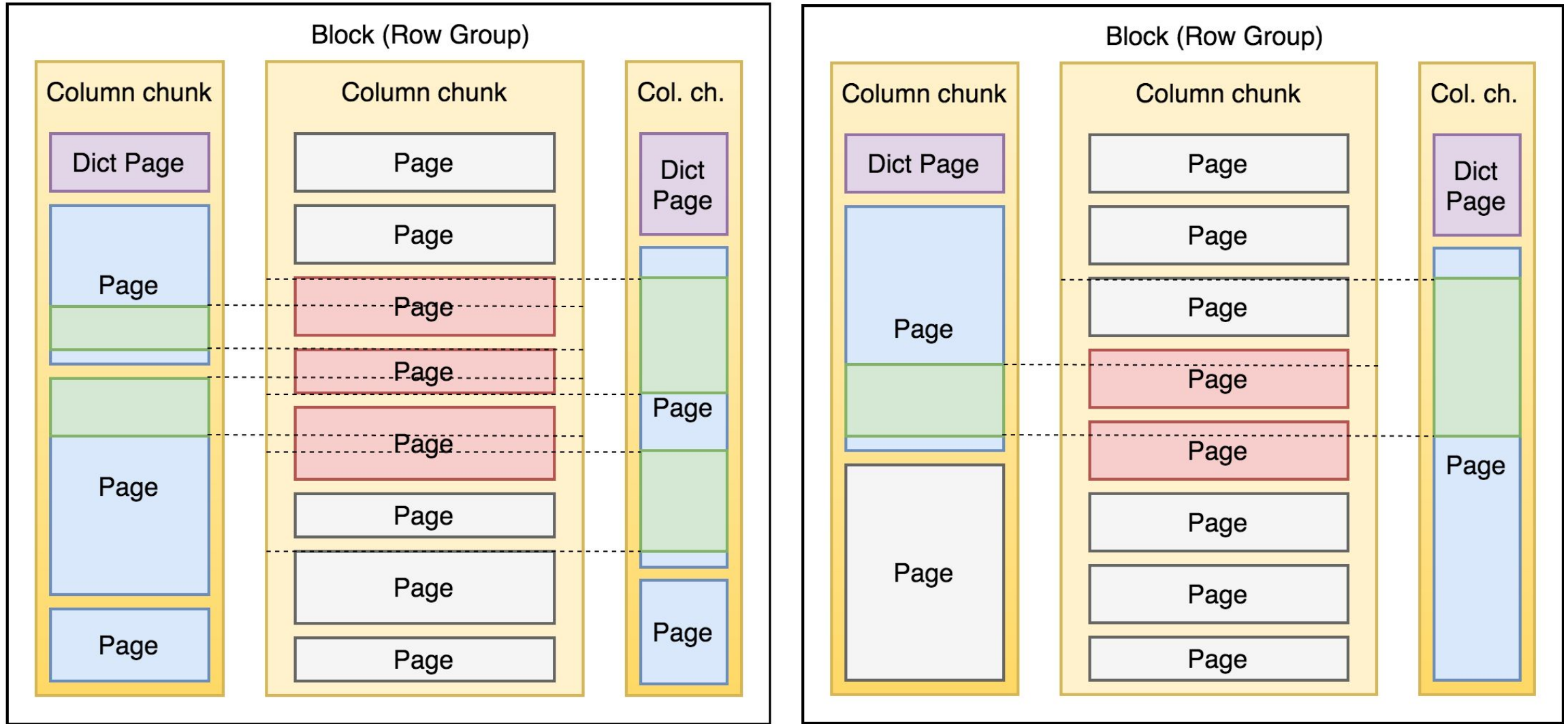
- Только простые условия (<, >, ==, IN, null)
- Заранее заданные константы
- Можно комбинировать логически (OR, AND, NOT)
- + Специальные условия для строк (StartsWith)

Predicate Pushdown

Не эффективен для чтения одной строки!

```
select * from table where id = 1234
```

Pred. Pushdown + Encoding + Sorting



**А теперь можно в
продакшен?**

Не всё так просто!

Приходится:

- Сортировать и партиционировать
- Оптимизировать типы
- Контролировать размеры

Оптимизация типов

- Чем меньше тип – тем лучше
- XML, JSON -> Infer schema -> Struct
- Plain text -> Parsed struct
- float -> int (?)

parquet.block.size

- Больше => лучше сжатие
- Больше памяти при записи
 - Требует x2-x3 памяти
- Должен уместиться в HDFS блоке

64 МБ – 1 ГБ (128 МБ)

Несколько блоков в файле?

- Формат позволяет
- Но нарушается граница HDFS блоков

Один файл – один блок!

parquet.block.size == dfs.block.size

- Делайте repartition перед записью
- Держите 10-20% запас

Repartition - ДО

df

// 200 tasks

.write

.parquet(path)

!hdfs dfs -ls -h \$path | tail -1

180.9M

/dataset/part-00199-hash.parquet

Repartition – после

df

.repartition(320)

.write

.parquet(path)

!hdfs dfs -ls -h \$path | tail -1

118.3M

/dataset/part-00319-hash.parquet

parquet.page.size

- Больше => лучше сжатие
- Меньше => эффективнее Predicate Pushdown
- Читается целиком в память

8 кб – 1 МБ (1 МБ)

parquet.dictionary.page.size

- Одна страница dictionary на колонку
- Держится в памяти целиком при чтении
- Увеличивайте при работе с повторяющимися строками

(1 МБ)

**Теперь-то всё будет
хорошо!..**

Spark Streaming – Append

stream

.write

.mode(Append)

.parquet(path)

Много маленьких файликов

⇒ Много HDFS блоков

- Неэффективное использование DataNode
- Высокая нагрузка на NameNode

⇒ Много Spark Tasks

- Большой оверхед на старт
- Нагрузка на мастер

1. Убирайте партиции

stream

.coalesce(1)

.write

.mode(Append)

.parquet(landPath)

2. Перекладывайте потоки

spark

.read.parquet(landPath)

.repartition(partitions, key)

.sortWithinPartitions(keys)

.write.parquet(path)

Spark VS Impala VS Hive

- Кто быстрее?
- Кто совместим?!

Имплементации Parquet

1. `apache/parquet-mr` (Java)
2. `apache/parquet-cpp` (C++)
3. Spark Catalyst (Scala)
4. `Dask/fastparquet` (Python)

Decimal

- Не читается

spark.sql.parquet.writeLegacyFormat

<https://issues.apache.org/jira/browse/IMPALA-2494>

<https://issues.apache.org/jira/browse/SPARK-10400>

<https://issues.apache.org/jira/browse/SPARK-6777>

<https://issues.apache.org/jira/browse/SPARK-20937>

<https://issues.apache.org/jira/browse/SPARK-20297>

<https://stackoverflow.com/questions/44279870/>

Timestamp

- Не читается
- Теряет таймзону

spark.sql.parquet.writeLegacyFormat

spark.sql.parquet.int96AsTimestamp

spark.sql.parquet.outputTimestampType
pe

spark.sql.parquet.int64AsTimestampMillis

<https://github.com/apache/spark/blob/master/sql/catalyst/src/main/scala/org/apache/spark/sql/internal/SQLConf.scala>

<https://issues.apache.org/jira/browse/HIVE-12767>

<https://issues.apache.org/jira/browse/HIVE-13534>

<https://issues.apache.org/jira/browse/SPARK-12297>

JSON/BSON

- Не нужен
- Но там всё равно что-то не работает

Spark Legacy Format

- Был дефолтным до Spark 1.5.9
 - Примерно 2015ый
 - Примерно Parquet 1.6.x
- Не задокументирован
 - SPARK-20937 «Может, всё-таки задокументируем?»

spark.sql.parquet.writeLegacyFormat

==

!spark.sql.parquet.followParquetFormatSpec

Impala?

- Ищет колонки по номерам, а не по именам
- Не поддерживает LZ0
- Не поддерживает binary
- Проблемы с Decimal и Timestamp

spark.sql.parquet.writeLegacyFormat
spark.sql.parquet.binaryAsString

Sqoop?

Не используйте Parquet в Sqoop!

- Не умеет repartitioning
- OOM'ы

ColumnWriter.writePage()

- Проверки, что пора писать очередной Page:
 - initial-page-run check
 - next-page-size check
- Когда можно ошибиться:
 - Строки большие с самого начала
 - Строки маленькие, но есть несколько больших

<https://issues.apache.org/jira/browse/PARQUET-99>

https://www.cloudera.com/documentation/enterprise/release-notes/topics/cdh_rn_sqoop_ki.html

Если совсем нельзя обойти

parquet.page.size.row.check.min

parquet.page.size.row.check.max

parquet.page.size.check.estimate

Schema Merging

- `_metadata, _common_metadata`
- Можно просто ОТКЛЮЧИТЬ

parquet.enable.summary-metadata=false

spark.sql.parquet.mergeSchema=false

А теперь – дьявольщина!

- Parquet буферизируется в памяти
- Контрольные суммы не предусмотрены
- Память может биться

⇒ Могут биться файлы

А нормального ничего нет?

Напомним, зачем нам Parquet?

- Экономия по месту
- Быстрая фильтрация
- Чтение по частям
- HDFS-native

- Очень дорогая запись

Может, CSV?

- Человекочитаемый
- Нет оверхеда для текстов
- Поддерживает append

- Бейзлайн по ужасности
- Плоская структура

Может, JSON?

- Человекочитаемый
- Schema-free
- Еще более медленный и жирный

А XML?



Avro?

- Поддерживает append
- HDFS-native
- Продвинутая эволюция схем
- Менее производительен, чем Parquet

Parquet VS ORC – всё сложно

- По объему и скорости однозначного лидера нет
- Hive отстаёт в поддержке Parquet
- Spark отстаёт в поддержке ORC
- Impala игнорирует существование ORC

MPP?

Для структурированных данных
специализированные MPP-системы
на порядок быстрее Spark + HDFS + Parquet.

Ясно, понятно...

Parquet

- Прекрасный формат для исторических данных
- Для Spark, особенно на CDN – альтернатив нет
- Имеет массу тонкостей
 - Эффективность варьируется на порядок

Главные тонкости

- Пиши один раз, читай много
- Структурируй данные
- Партиционируй и перекладывай потоки
- Оптимизируй размеры
- Не верь в интероперабельность

Спасибо!

Федор Лаврентьев

t.me/fediq