**BIGDATA TEAM**

# Формат данных Parquet и сравнение с ORC

**Драль Алексей**, study@bigdatateam.org
CEO at BigData Team, https://bigdatateam.org
https://www.facebook.com/bigdatateam

# Dremel: Interactive Analysis of Web-Scale Datasets

Sergey Melnik, Andrey Gubarev, Jing Jing Long, Geoffrey Romer,
Shiva Shivakumar, Matt Tolton, Theo Vassilakis
Google, Inc.
{melnik,andrey,jlong,gromer,shiva,mtolton,theov}@google.com

## ABSTRACT

Dremel is a scalable, interactive ad-hoc query system for analysis of read-only nested data. By combining multi-level execution trees and columnar data layout, it is capable of running aggregation queries over trillion-row tables in seconds. The system scales to thousands of CPUs and petabytes of data, and has thousands of users at Google. In this paper, we describe the architecture and implementation of Dremel, and explain how it complements MapReduce-based computing. We present a novel columnar storage representation for nested records and discuss experiments on few-thousand node instances of the system.

## 1. INTRODUCTION

Large-scale analytical data processing has become widespread in web companies and across industries, not least due to low-cost exchanged by distributed systems, structured documents, etc. lend themselves naturally to a *nested* representation. Normalizing and recombining such data at web scale is usually prohibitive. A nested data model underlies most of structured data processing at Google [21] and reportedly at other major web companies.

This paper describes a system called Dremel[1] that supports interactive analysis of very large datasets over shared clusters of commodity machines. Unlike traditional databases, it is capable of operating on *in situ* nested data. *In situ* refers to the ability to access data 'in place', e.g., in a distributed file system (like GFS [14]) or another storage layer (e.g., Bigtable [8]). Dremel can execute many queries over such data that would ordinarily require a sequence of MapReduce (MR [12]) jobs, but at a fraction of the execution time. Dremel is not intended as a replacement for MR and is often used in conjunction with it to analyze outputs of MR pipelines or rapidly prototype larger computations.

# Пример страницы Википедии

Not logged in   Talk   Contributions   Create account   Log in

Article   Talk                                    Read   Edit   View history   Search Wikipedia

WIKIPEDIA
The Free Encyclopedia

Main page
Contents
Current events
Random article
About Wikipedia
Contact us
Donate

Contribute

Help
Learn to edit
Community portal
Recent changes
Upload file

Tools

What links here
Related changes
Special pages
Permanent link
Page information
Cite this page

## Apache Hadoop

From Wikipedia, the free encyclopedia

**Apache Hadoop** ( /həˈduːp/ ) is a collection of open-source software utilities that facilitates using a network of many computers to solve problems involving massive amounts of data and computation. It provides a software framework for distributed storage and processing of big data using the MapReduce programming model. Hadoop was originally designed for computer clusters built from commodity hardware, which is still the common use.[3] It has since also found use on clusters of higher-end hardware.[4][5] All the modules in Hadoop are designed with a fundamental assumption that hardware failures are common occurrences and should be automatically handled by the framework.[6]

The core of Apache Hadoop consists of a storage part, known as Hadoop Distributed File System (HDFS), and a processing part which is a MapReduce programming model. Hadoop splits files into large blocks and distributes them across nodes in a cluster. It then transfers packaged code into nodes to process the data in parallel. This approach takes advantage of data locality,[7] where nodes manipulate

**Apache Hadoop**
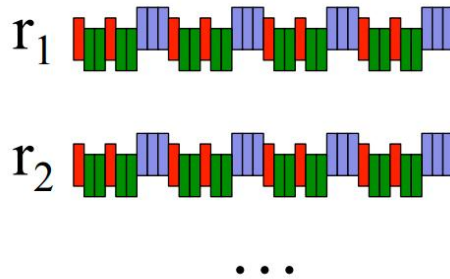
APACHE Hadoop

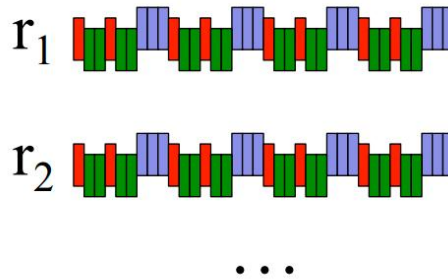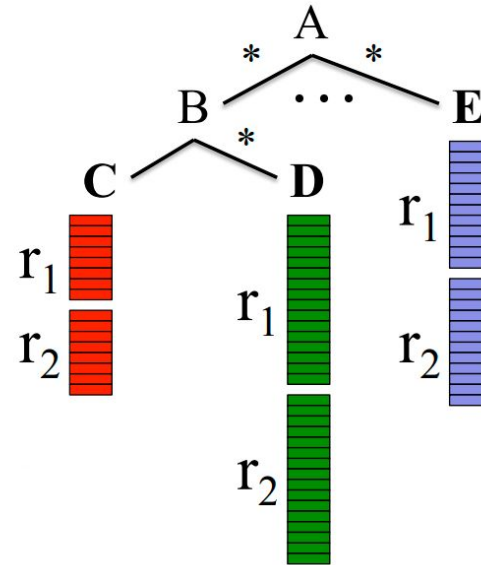| | |
|---|---|
| **Original author(s)** | Doug Cutting, Mike Cafarella |
| **Developer(s)** | Apache Software Foundation |
| **Initial release** | April 1, 2006; 15 years ago[1] |
| **Stable release** | 2.7.x  2.7.7 / May 31, 2018; 2 years ago[2] |
| | 2.8.x  2.8.5 / September 15, 2018; 2 years ago[2] |
| | 2.9.x  2.9.2 / November 9, 2018; 2 years ago[2] |

## History  [ edit ]

record-oriented

$r_1$

$r_2$

. . .

https://research.google/pubs/pub36632/

record-oriented

column-oriented

https://research.google/pubs/pub36632/

```
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'
```

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

Figure 2: Two sample nested records and their schema

```
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'
```

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

Figure 2: Two sample nested records and their schema

# Google Dremel*

Figure 2: Two sample nested records and their schema



Figure 3: Column-striped representation of the sample data in Figure 2, showing repetition levels (r) and definition levels (d)

# Лайфхаки для понимания

▶ repetition и definition levels



| Document: R = 0, D = 0 |
| --- |
| DocId: *required* |
| Links: *optional* D = 1 |
|   Backward: *repeated* R = 1, D = 2 |
|   Forward: *repeated* R = 1, D = 2 |
| Name: *repeated* R = 1, D = 1 |
|   Language: *repeated* R = 2, D = 2 |
|     Code: *required* |
|     Country: *optional* D = 3 |
|   Url: *optional* D = 2 |

| R = 0 | R = 1 | R = 2 |
| --- | --- | --- |
| Document.DocId | | |
| Document.Links | Backward | |
| Document.Links | Forward | |
| Document | Name | Language.Code |
| Document | Name | Language.Country |
| Document | Name.Url | |

*required*: same Repetition and Definition level as parent
*optional*: same Repetition level as parent, increment Definition level
*repeated*: increment both Repetition and Definition levels

| D = 0 | D = 1 | D = 2 | D = 3 |
| --- | --- | --- | --- |
| Document.DocId | | | |
| Document | Links | Backward | |
| Document | Links | Forward | |
| Document | Name | Language.Code | |
| Document | Name | Language | Country |
| Document | Name | Url | |

The striping and assembly algorithms from the Dremel paper

- File ~ блок в HDFS
- Row group ~ stripe
- Column Chunk
- Page ~ Compression Chunk

|  | ORC | Parquet |
|---|:---:|:---:|
| Поколоночное хранение | ✓ | ✓ |

|  | ORC | Parquet |
|---|---|---|
| Поколоночное хранение | ✓ | ✓ |
| Версионирование схемы | ✓ (protobuf) | ✓ (thrift) |

# Сравнение

| | ORC | Parquet |
|---|---|---|
| Поколоночное хранение | ✓ | ✓ |
| Версионирование схемы | ✓ (protobuf) | ✓ (thrift) |
| Статистики и фильтры Блума | ✓ ✓ | ✓ ✗ |

false

| | ORC | Parquet |
|---|---|---|
| Поколоночное хранение | ✓ | ✓ |
| Версионирование схемы | ✓ (protobuf) | ✓ (thrift) |
| Статистики и фильтры Блума | ✓ ✓ | ✓ ✗ |
| Типы данных | 14+ | 7 + nested |

**BIGDATA TEAM**

| | ORC | Parquet |
|---|---|---|
| Поколоночное хранение | ✓ | ✓ |
| Версионирование схемы | ✓ (protobuf) | ✓ (thrift) |
| Статистики и фильтры Блума | ✓✓ | ✓✗ |
| Типы данных | 14+ | 7 + nested |
| Ориентация на мир Hadoop и world-wide | ✓✗ | ✓✓ |