



Australian
National
University

程序代写代做 CS编程辅导



NoSQL Databases – Part 3

WeChat: cstutorcs

Column-oriented Data Stores

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导¹ Column-oriented Data Stores¹



- Inspired by Google's
- Store **data grouped** (rather than rows) and may have a very large number of columns



QQ: 749389476

- Other column-oriented data stores
 - Hbase
 - Hypertable

¹ Figure source: S. Harizopoulos, D. Abadi and P. Boncz, Column-Oriented database systems, VLDB 2009



程序代写代做 CS编程辅导

Google's Bigtable - Problem Analysis



- Used by over 60 projects at Google as of 2006, including Web indexing, Google Earth, Google Maps, Orkut, Google Docs, etc.
- Data types vary** from text to web pages to satellite imagery.
- Latency requirements vary** from backend bulk processing to real-time data processing.
- Infrastructures vary** from a handful to thousands of servers.

WeChat: cstutorcs

Assignment Project Exam Help

- Need to scale to a very large size** such as petabytes of data across thousands of commodity servers.

Email: tutorcs@163.com

QQ: 749389476

- Most applications **require only single-row transactions**.

<https://tutorcs.com>





程序代写代做 CS编程辅导

Google's Bigtable - Problem Analysis



Key questions:

- 1 How to represent data? (**expressiveness**)

Key-value pairs are useful but limited

- 2 How to store data? (**scalability**)

Data needs to be distributed across multiple servers

- 3 How to process data? (**efficiency**)

Join on distributed tables needs to be avoided

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



程序代写代做 CS编程辅导

Google's Bigtable - Problem Analysis



- Key questions:

- 1 How to represent (expressiveness)

Key-value pairs are useful but limited

- 2 How to store data (scalability)

Data needs to be distributed across multiple servers

- 3 How to process data? (efficiency)

Join on distributed tables needs to be avoided

- Solution:

QQ: 749389476

One big table
<https://tutors.com>

in which **both rows and columns** can be split over multiple servers,
according to **their relatedness**.



Google's Bigtable - Data Structure



- A (big) table is a **multi-dimensional sparse sorted map**.

WeChat: cstutorcs

(**row key**, **column key**, **timestamp**) \mapsto **value**

Assignment Project Exam Help

- The map is **indexed by** a row key, a column key, and a timestamp.

- Each value in the map is **an uninterpreted array of bytes**.

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



Google's Bigtable - Data Structure



- A (big) table is a **multidimensional sparse sorted map**.

(**row key**, **column key**, **timestamp**) \mapsto **value**

- Example:** a (big) table that stores Web pages

ROW KEY	COLUMN	COLUMN	COLUMN	...
	CONTENTS	ANCHOR: CNN.COM	ANCHOR: MY.LOOK.CA	...
com.cnn.www	$\langle \text{html} \rangle \langle \text{body} \rangle \text{Home} \dots \leftarrow t_1$ 404 Page not found $\leftarrow t_2$ $\langle \text{html} \rangle \langle \text{body} \rangle \text{Inter} \dots \leftarrow t_3$	CNN $\leftarrow t_9$	CNN.com $\leftarrow t_8$...
com.cnn.weather
com.cnn.live
...

<https://tutorcs.com>

- ("com.cnn.www", "CONTENTS:", t_1) \mapsto " $\langle \text{html} \rangle \langle \text{body} \rangle \text{Home} \dots$ "
- ("com.cnn.www", "ANCHOR:MY.LOOK.CA", t_8) \mapsto "CNN.com"



程序代写代做 CS编程辅导

Google's Bigtable - Data Structure (Row Key)



- Row keys are strings of up to 1024 KB size.
- Row keys are sorted in a lexicographical order.

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutores@163.com

QQ: 749389476

<https://tutores.com>

- Every read or write of data under a single row key is atomic (regardless of the number of different columns being read or written in the row).



程序代写代做 CS编程辅导

Google's Bigtable - Data Structure (Row Key)



- A table is **dynamically divided into tablets** (each approximately 100-200 MB in size). A tablet can be regarded as a horizontal partition in a table.
- Tablets are **the basic units of distribution and load balancing**, served by **tablet servers**.

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

	ROW KEY	
tablet ₁	com.cnn.www	...
	com.cnn.weather	...
	com.cnn.live	...

tablet ₂	nz.ac.otago.www	...
	nz.ac.otago.cs	...

...



Google's Bigtable - Data Structure (Row Key)



- Question: Why are *URLs* used as row keys?

	Row Key	...
tablet1	com.cnn.www	...
	com.cnn.weather	...
	com.cnn.live	...
tablet4	nz.ac.otago.www	...
	nz.ac.otago.cs	...

WeChat: tutormcs

Assignment Project Exam Help

Email: tutormcs@163.com

QQ: 749389476

https://tutormcs.com

Google's Bigtable - Data Structure (Row Key)



- Applications need to choose row keys
 - The ordering of row keys affects partitioning of rows into tablets.
 - Row ranges with lexicographical distance are split into fewer tablets (good for reads).

WeChat: cstutores

Assignment Project Exam Help

Email: tutores@163.com

QQ: 749389476

<https://tutores.com>

	Row KEY	...
	com.cnn.www	...
tablet ₁	com.cnn.weather	...
	com.cnn.live	...

	nz.ac.otago.www	...
tablet ₂	nz.ac.otago.cs	...

...

- As a result, reads of short row ranges are efficient and typically require communication with only a small number of machines.



Google's Bigtable - Data Structure (Column)



- Columns are **grouped into column families**, i.e., a column family contains columns of related data. A column is named as **family:qualifier**, e.g.,

WeChat: cstutorcs

COLUMN FAMILY 1	COLUMN FAMILY 2		
CONTENTS:	ANCHOR.UNNSI.COM	ANCHOR.MY.LOOK.CA	...
...

Email: tutorcs@163.com

QQ: 749389476

- Question:** Why are columns grouped into column families?

<https://tutorcs.com>



Google's Bigtable - Data Structure (Column)



- Some properties

- Column families are the **basic unit of access control**, discerning privileges to read, modify, create column-families, etc.
- They can be vertically partitioned into different files.
- Column families need to be defined in the schema (before data can be stored) but **columns within a family can be dynamically changed**.

Email: tutorcs@163.com

COLUMN FAMILY 1	COLUMN FAMILY 2		
CONTENTS:	ANCHOR:CNNS.COM	ANCHOR:MY.LOOK.CA	...
...

<https://tutorcs.com>

- The number of column families should be small (in the hundreds at most).



程序代写代做 CS编程辅导

Google's Bigtable - Data Structure (Timestamp)



- Each cell can contain **multiple versions** of the same data, indexed by timestamp.

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

- Each cell version is **a string**, i.e., a scalar value.
- Stored in decreasing timestamp order, and thus the most recent version can be **read first**.

QQ: 749389476

<https://tutorcs.com>



Google's Bigtable - Read Operations



```
Scanner scanner(T
ScanStream *stream
stream = scanner.l
InFamily("anchor");
stream->SetReturnAllVersions();
scanner.Lookup("com.cnn.www");
for ( ; !stream->Done(); stream->Next()) {
    printf("%s %s %11d %s\n",
        scanner->RowName(),
        stream->ColumnName(),
        stream->Microtimestamp(),
        stream->Value());}
```

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

ROW KEY	COLUMN	COLUMN	COLUMN	...
	CONTENTS:	ANCHOR:CNNSI.COM	ANCHOR:MY.LOOK.CA	...
com.cnn.www	<html><body>Home for u... 404 Page not found ← t_2 </html></body>Inter... ← t_3	CNN ← t_9	CNN.com ← t_8	...
com.cnn.weather
com.cnn.live
...



Google's Bigtable - Write Operations



```
# Open the table
Table *T = OpenOrigTable("table/web/webtable");
```

```
# Write a new anchor and delete an old anchor
```

```
RowMutation r1(T, "com.cnn.www");
r1.Set("anchor:www.c-span.org", "CNN");
```

```
r1.Delete("anchor:my.look.ca");
```

```
Operation op;
```

```
Apply(&op, &r1);
```

Email: tutorcs@163.com

ROW KEY	COLUMN CONTENTS	COLUMN ANCHOR:CNNSI.COM	COLUMN ANCHOR:MY.LOOK.CA	...
com.cnn.www	$\langle \text{html} \rangle \langle \text{body} \rangle \text{Home} \dots \leftarrow t_1$ 404 Page not found $\leftarrow t_2$ $\langle \text{html} \rangle \langle \text{body} \rangle \text{Inter} \dots \leftarrow t_3$	CNN $\leftarrow t_9$	CNN.com $\leftarrow t_8$...
com.cnn.weather
com.cnn.live
...



程序代写代做 CS编程辅导

Google's Bigtable - Infrastructure Dependencies



- Bigtable is built upon the following components:

- **Google file system**: a highly scalable distributed file system
 - e.g., store table data and log.
- **Chubby lock service**: a highly-available and persistent distributed lock service
 - e.g., handles master election, manage metadata, etc.
- **MapReduce programming model**: a parallel computing model
 - Google's batch processing tool of choice
- **Cluster scheduling system**: a cluster management system
 - e.g., handles failover, monitoring, etc.
- ...

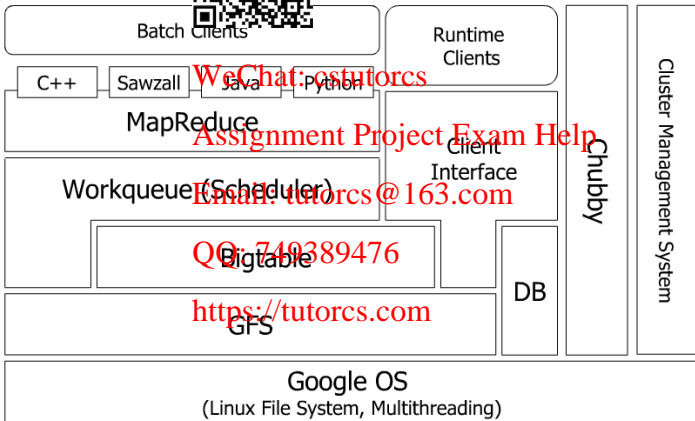
<https://tutorcs.com>

- Similar components are being made available as Open Source by the Apache project Hadoop.



程序代写代做 CS编程辅导 Google's Overall Architecture

- Use **shared-nothing** architecture, consisting of thousands of commodity machines.





程序代写代做 CS编程辅导 Google's Bigtable - Summary



- Uses a **shared-nothing architecture** to provide scalability over massive data sets:
 - **Horizontal partitioning** by range of row keys.
 - **Vertical partitioning** by column families
- **Replication**: eventual-consistency replication across datacenters, between multiple BigTable serving setups (master/slave & multi-master)
- Supports **single-row transactions**.
- Supports **only simple queries**.
- Does **not support secondary indices**.