

4

### 从BigTable说起

◎ 中南大学

- •BigTable是一个分布式存储系统
- •利用谷歌提出的MapReduce分布式并行计算模型来处理海量数据
- •使用谷歌分布式文件系统GFS作为底层数据存储
- •采用Chubby提供协同服务管理
- •可以扩展到PB级别的数据和上千台机器,具备广泛应用性、可扩展性、高性能和高可用性等特点
- •谷歌的许多项目都存储在BigTable中,包括搜索、地图、财经、打印、社交 网站Orkut、视频共享网站YouTube和博客网站Blogger等



◎ 中南大学

#### HBase简介



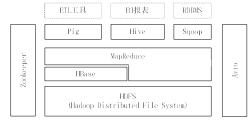
#### 表5-1 HBase和BigTable的底层技术对应关系

		HBase
	BigTable	
文件存储系统	GFS	HDFS
海量数据处理	MapReduce	Hadoop MapReduce
协同服务管理	Chubby	Zookeeper

## HBase简介

HBase是一个高可靠、高性能、面向列、可伸缩的分布式数据库,是谷歌BigTable的开源实现,主要用来存储非结构化和半结构化的松散数据。HBase的目标是处理非常庞大的表,可以通过水平扩展的方式,利用廉价计算机集群处理由超过10亿行数据和数百万列元素组成的数据表

#### Hadoop生态系统



#### 图4-1 Hadoop生态系统中HBase与其他部分的关系

## HBase简介

# ◎ 中南大学

#### 关系数据库已经流行很多年,并且Hadoop已经有了HDFS和MapReduce,为什么需要HBase?

- •Hadoop可以很好地解决大规模数据的离线批量处理问题,但是,受限于Hadoop MapReduce编程框架的高延迟数据处理机制,使得Hadoop无法满足大规模数据实时处理应用的需求
- •HDFS面向批量访问模式,不是随机访问模式
- •传统的通用关系型数据库无法应对在数据规模剧增时导致的系统扩展性和性能问题(分库分表也不能很好解决)
- •传统关系数据库在数据结构变化时一般需要停机维护;空列浪费存储空间
- •因此,业界出现了一类面向半结构化数据存储和处理的高可扩展、低写入/查询延迟的系统,例如,键值数据库、文档数据库和列族数据库(如BigTable和HBase等).目前HBase已经成功应用于互联网服务领域和传统行业的众多在线式数据分析处理系统中

◎ 中南大学

## ◎中南大学

#### HBase与传统关系数据库的对比分析

- •HBase与传统的关系数据库的区别主要体现在以下几个方面:
- (1) 数据类型:关系数据库采用关系模型,具有丰富的数据类型和存储方式,HBase则采用了更加简单的数据模型,它把数据存储为未经解释的字符串
- (2) 数据操作:关系数据库中包含了丰富的操作,其中会涉及复杂的多表连接。HBase操作则不存在复杂的表与表之间的关系,只有简单的插入、查询、删除、清空等,因为HBase在设计上就避免了复杂的表和表之间的关系
- (3) 存储模式:关系数据库是基于行模式存储的。HBase是基于列存储的,每个列族都由几个文件保存,不同列族的文件是分离的

## HBase访问接口

#### 表5-2 HBase访问接口

类型	特点	场合
Native Java API	最常规和高效的访问方式	适合Hadoop MapReduce作业 并行批处理HBase表数据
HBase Shell	HBase的命令行工具,最 简单的接口	适合HBase管理使用
Thrift Gateway	利用Thrift序列化技术, 支持C++、PHP、Python 等多种语言	适合其他异构系统在线访问 HBase表数据
REST Gateway	解除了语言限制	支持REST风格的Http API访问 HBase
Pig	使用Pig Latin流式编程语言来处理HBase中的数据	适合做数据统计
Hive	简单	当需要以类似SQL语言方式来 访问HBase的时候

## HBase与传统关系数据库的对比分析

#### •HBase与传统的关系数据库的区别主要体现在以下几个方面:

- (4) 数据索引:关系数据库通常可以针对不同列构建复杂的多个索引,以提高数据访问性能。HBase只有一个索引——行键,通过巧妙的设计,HBase中的所有访问方法,或者通过行键访问,或者通过行键扫描,从而使得整个系统不会慢下来
- (5) 数据维护: 在关系数据库中,更新操作会用最新的当前值去替换记录中原来的旧值 ,旧值被覆盖后就不会存在。而在HBase中执行更新操作时,并不会删除数据旧的版本, 而是生成一个新的版本,旧有的版本仍然保留
- (6) 可伸缩性:关系数据库很难实现横向扩展,纵向扩展的空间也比较有限。相反, HBase和BigTable这些分布式数据库就是为了实现灵活的水平扩展而开发的,能够轻易地 通过在集群中增加或者减少硬件数量来实现性能的伸缩



◎ 中南大学



## 数据模型概述

◎ 中南大学

◎ 中南大学

- •HBase是一个稀疏、多维度、排序的映射表,这张表的索引是行键、列族、列限定符和时间戳
- •每个值是一个未经解释的字符串,没有数据类型
- •用户在表中存储数据、每一行都有一个可排序的行键和任意多的列
- •表在水平方向由一个或者多个列族组成,一个列族中可以包含任意多个列,同一个 列族里面的数据存储在一起
- •列族支持动态扩展,可以很轻松地添加一个列族或列,无需预先定义列的数量以及 类型,所有列均以字符串形式存储,用户需要自行进行数据类型转换
- •HBase中执行更新操作时,并不会删除数据旧的版本,而是生成一个新的版本,旧有的版本仍然保留(这是和HDFS只允许追加不允许修改的特性相关的)

#### 数据坐标

• HBase中需要根据行键、列族、列限定符和时间戳来确定一个单元格,因此,可以视为一个"四维坐标",即[行键, 列族, 列限定符, 时间戳]

键 值 ["201505003", "Info", "email", 1174184619081] "xie@qq.com" ["201505003", "Info", "email", 1174184620720] "you@163.com"



## 数据模型相关概念

- •表:HBase采用表来组织数据,表由行和列组成, 列划分为若干个列族
- 行: 每个HBase表都由若干行组成,每个行由行键 (row key)来标识。
- 列族: 一个HBase表被分组成许多"列族"(Column Family)的集合,它是基本的访问控制单元
- <mark>列限定符:</mark>列族里的数据通过列限定符(或列)来<sub>行键</sub> 定位
- <mark>单元格:</mark> 在HBase表中,通过行、列族和列限定符确定一个"单元格"(cell),单元格中存储的数据没有数据类型,总被视为字节数组byte[]
- •时间**戳**:每个单元格都保存着同一份数据的多个版本,这些版本采用时间戳进行索引



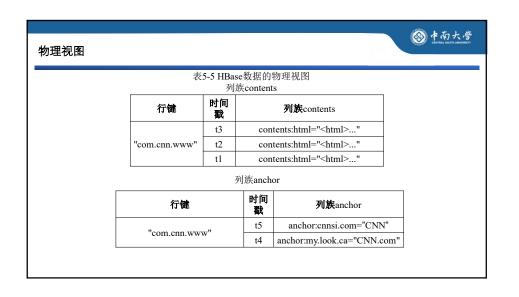
## 概念视图

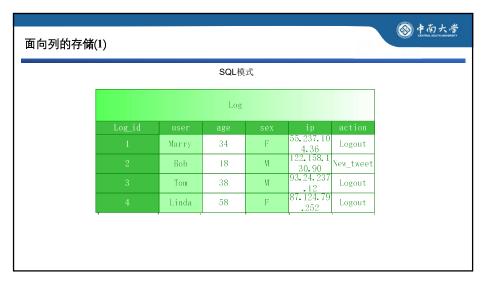
#### 表5-4 HBase数据的概念视图

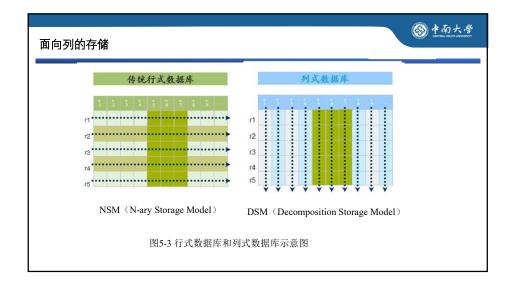
行键	哲回蠽	列族contents	列族anchor
	t5		anchor:cnnsi.com="CNN"
	t4		anchor:my.look.ca="CNN.com"
"com.cnn .www"	t3	contents:html="< html>"	
	t2	contents:html="< html>"	
	t1	contents:html="< html>"	

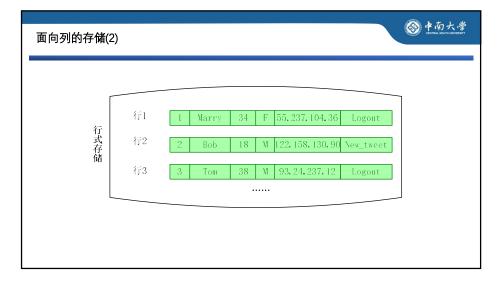


◎ 中南大学

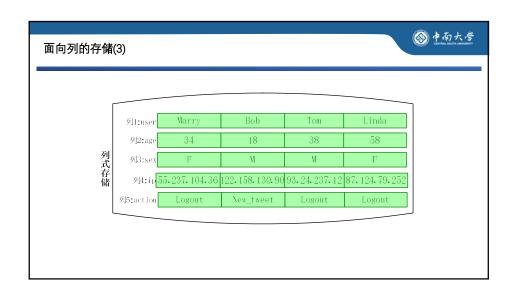


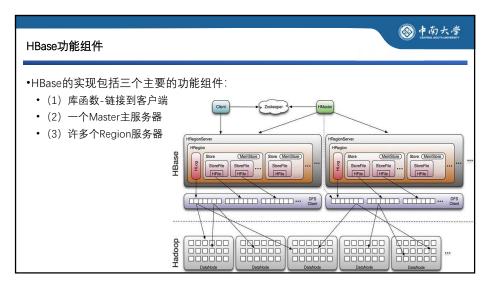






\_





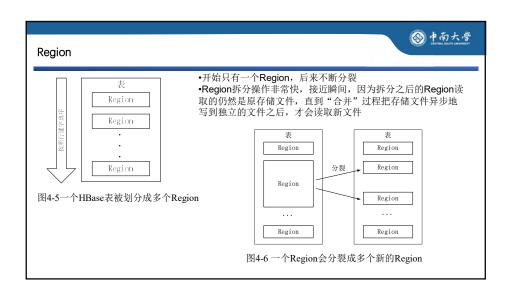


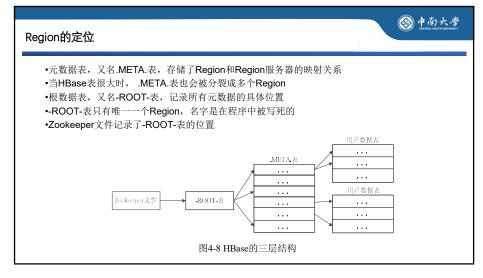
## HBase功能组件

- ◎ 中南大学
- •主服务器Master负责管理和维护HBase表的分区信息,维护Region服务器列表,分配 Region,负载均衡
- •Region服务器负责存储和维护分配给自己的Region,处理来自客户端的读写请求
- •客户端并不是直接从Master主服务器上读取数据,而是在获得Region的存储位置信息后 ,直接从Region服务器上读取数据
- •客户端并不依赖Master,而是通过Zookeeper来获得Region位置信息,大多数客户端甚至 从来不和Master通信,这种设计方式使得Master负载很小

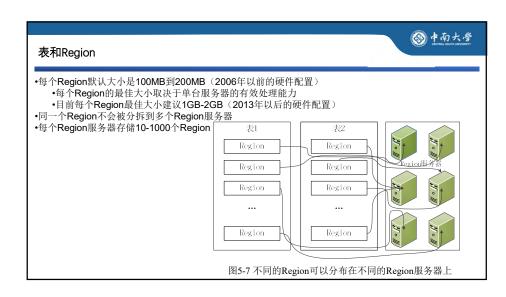
\_

◎ 中南大学





Region的定位





## 4.4.3 Region的定位

◎ 中南大学

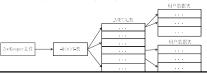
•假设.META.表的每行(一个映射条目)在内存中大约占用1KB,并且每个Region限制为128MB,那么,上面的三层结构可以保存的用户数据表的Region数目的计算方法是:

(-ROOT-表能够寻址的.META.表的Region个数)×(每个.META.表的 Region可以寻址的用户数据表的Region个数)

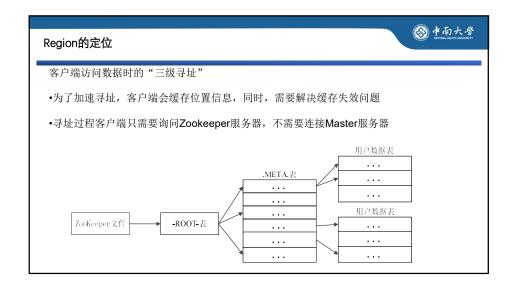
•一个-ROOT-表最多只能有一个Region,也就是最多只能有128MB,按照每行(一个映射条目)占用 1KB内存计算,128MB空间可以容纳128MB/1KB=2<sup>17</sup>行,也就是说,一个-ROOT-表可以寻址2<sup>17</sup> 个.META.表的Region。

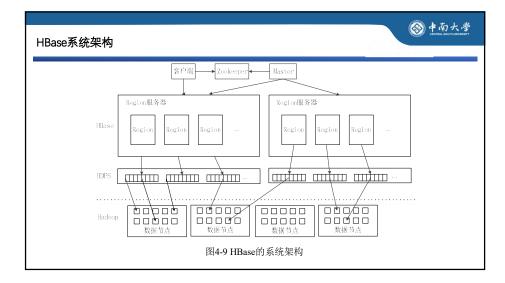
•同理,每个.META.表的 Region可以寻址的用户数据表的Region个数是128MB/1KB=217。

•最终,三层结构可以保存的Region数目是(128MB/1KB) × (128MB/1KB) = 234个Region

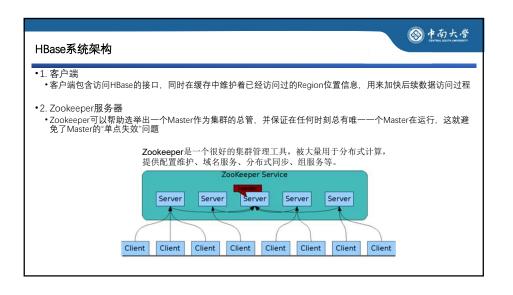


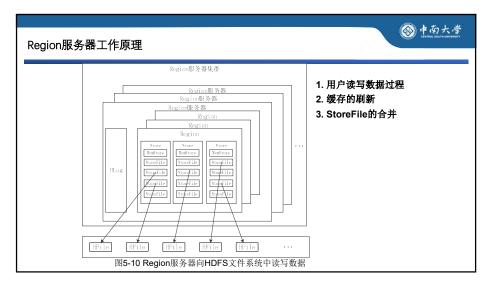






\_





#### HBase系统架构

**一个中央大学** 

- 3. Master主服务器,主要负责表和Region的管理工作:
  - 管理用户对表的增加、删除、修改、查询等操作
  - 实现不同Region服务器之间的负载均衡
  - 在Region分裂或合并后,负责重新调整Region的分布
  - 对发生故障失效的Region服务器上的Region进行迁移
- 4. Region服务器
  - Region服务器是HBase中最核心的模块,负责维护分配给自己的Region,并响应用户的读写请求



## Region服务器工作原理



#### 1. 用户读写数据过程

- •用户写入数据时,被分配到相应Region服务器去执行
- •用户数据首先被写入到MemStore和Hlog中
- •只有当操作写入Hlog之后,commit()调用才会将其返回给客户端
- •当用户读取数据时,Region服务器会首先访问MemStore缓存,如果找不
- 到,再去磁盘上面的StoreFile中寻找



## Region服务器工作原理

# ◎ 中南大学

## 2. 缓存的刷新

- •系统会周期性地把MemStore缓存里的内容刷写到磁盘的StoreFile文件中, 清空缓存,并在Hlog里面写入一个标记
- •每次刷写都生成一个新的StoreFile文件,因此,每个Store包含多个 StoreFile文件
- •每个Region服务器都有一个自己的HLog 文件,每次启动都检查该文件,确认最近一次执行缓存刷新操作之后是否发生新的写入操作;如果发现更新,则先写入MemStore,再刷写到StoreFile,最后删除旧的Hlog文件,开始为用户提供服务



## Store工作原理



- •Store是Region服务器的核心
- •多个StoreFile合并成一个
- •单个StoreFile过大时,又触发分裂操作,1个父Region被分裂成两个子Region

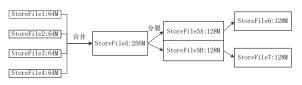


图4-11 StoreFile的合并和分裂过程



## Region服务器工作原理

# ◎ 中南大学

## 3. StoreFile的合并

- •每次刷写都生成一个新的StoreFile,数量太多,影响查找速度
- •调用Store.compact()把多个合并成一个
- •合并操作比较耗费资源,只有数量达到一个阈值才启动合并



## HLog工作原理



- •分布式环境必须要考虑系统出错。HBase采用HLog保证系统恢复
- •HBase系统为每个Region服务器配置了一个HLog文件,它是一种预写式日志 (Write Ahead Log)
- •用户更新数据必须首先写入日志后,才能写入MemStore缓存,并且,直到 MemStore缓存内容对应的日志已经写入磁盘,该缓存内容才能被刷写到磁 盘



## HLog工作原理



HBASE

- Zookeeper会实时监测每个Region服务器的状态,当某个Region服务器发生故障时, Zookeeper会通知Master
- Master首先会处理该故障Region服务器上面遗留的HLog文件,这个遗留的HLog文件中包含了来自多个Region对象的日志记录
- 系统会根据每条日志记录所属的Region对象对HLog数据进行拆分,分别放到相应 Region对象的目录下,然后,再将失效的Region重新分配到可用的Region服务器中,并 把与该Region对象相关的HLog日志记录也发送给相应的Region服务器
- Region服务器领取到分配给自己的Region对象以及与之相关的HLog日志记录以后,会重新做一遍日志记录中的各种操作,把日志记录中的数据写入到MemStore缓存中,然后,刷新到磁盘的StoreFile文件中,完成数据恢复
- 共用日志优点:提高对表的写操作性能;缺点:恢复时需要分拆日志

## HBase实际应用中的性能优化方法



#### ·行键(Row Key)

行键是按照**字典序**存储,因此,设计行键时,要充分利用这个排序特点,将经常一起读取的数据存储到一块,将最近可能会被访问的数据放在一块。

举个例子:如果最近写入HBase表中的数据是最可能被访问的,可以考虑将时间戳作为行键的一部分,由于是字典序排序,所以可以使用Long.MAX\_VALUE - timestamp作为行键,这样能保证新写入的数据在读取时可以被快速命中。

# の5 HBASE 应用方案

## HBase实际应用中的性能优化方法



#### InMemory

创建表的时候,可以通过HColumnDescriptor.setInMemory(true)将表放到Region服务器的缓存中,保证在读取的时候被cache命中。

#### Max Version

创建表的时候,可以通过HColumnDescriptor.setMaxVersions(int maxVersions)设置表中数据的最大版本,如果只需要保存最新版本的数据,那么可以设置setMaxVersions(1)。

#### Time To Live

创建表的时候,可以通过HColumnDescriptor.setTimeToLive(int timeToLive)设置表中数据的存储生命期,过期数据将自动被删除,例如如果只需要存储最近两天的数据,那么可以设置setTimeToLive(2\*24\*60\*60)。

## HBase实际应用中的性能优化方法



#### InMemory

创建表的时候,可以通过HColumnDescriptor.setInMemory(true)将表放到Region服务器的缓存中,保证在读取的时候被cache命中。

#### Max Version

创建表的时候,可以通过HColumnDescriptor.setMaxVersions(int maxVersions)设置表中数据的最大版本,如果只需要保存最新版本的数据,那么可以设置setMaxVersions(1)。

#### Time To Live

创建表的时候,可以通过HColumnDescriptor.setTimeToLive(int timeToLive)设置表中数据的存储生命期,过期数据将自动被删除,例如如果只需要存储最近两天的数据,那么可以设置setTimeToLive(2\*24\*60\*60)。



#### HBase实际应用中的性能优化方法

# ◎ 中南大学

#### ·写表操作中的HTable参数设置

#### Auto Flush

创建表的时候,可以通过HColumnDescriptor.setInMemory(true)将表放到Region服务器的缓存中,保证在读取的时候被cache命中。

#### Max Version

创建表的时候,可以通过HColumnDescriptor.setMaxVersions(int maxVersions)设置表中数据的最大版本,如果只需要保存最新版本的数据,那么可以设置setMaxVersions(1)。

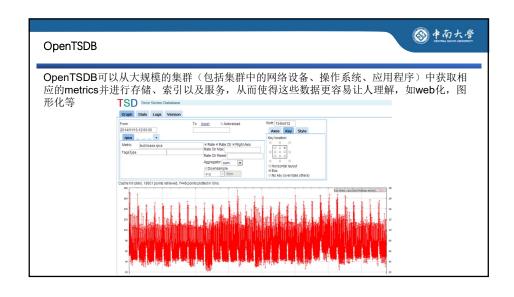
#### •Time To Live

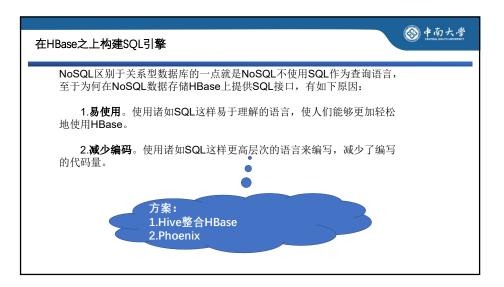
创建表的时候,可以通过HColumnDescriptor.setTimeToLive(int timeToLive)设置表中数据的存储生命期,过期数据将自动被删除,例如如果只需要存储最近两天的数据,那么可以设置setTimeToLive(2 \* 24 \* 60 \* 60)。











## 在HBase之上构建SQL引擎



#### 1.Hive整合HBase

Hive与HBase的整合功能从Hive0.6.0版本已经开始出现,利用两者对外的API接口互相通信,通信主要依靠hive\_hbase-handler.jar工具包(Hive Storage Handlers)。由于HBase有一次比较大的版本变动,所以并不是每个版本的Hive都能和现有的HBase版本进行整合,所以在使用过程中特别注意的就是两者版本的一致性。

#### 2.Phoenix

Phoenix由Salesforce.com开源,是构建在Apache HBase之上的一个SQL中间层,可以让开发者在HBase上执行SQL查询。

## 构建HBase二级索引

# ◎ 中南大学

#### ·Coprocessor构建二级索引

- Coprocessor提供了两个实现: endpoint和observer, endpoint相当于关系型数据库的存储过程,而observer则相当于触发器
- observer允许我们在记录put前后做一些处理,因此,而我们可以在插入数据时同步写入索引表

#### 优点:

非侵入性:引擎构建在HBase之上,既没有对HBase 进行任何改动,也不需要上层应用做任何妥协

•缺点: 每插入一条数据需要向索引表插入数据,即耗时是双倍的,对HBase的集群的压力也是双倍的

#### 构建HBase二级索引

◎ 中南大学

HBase只有一个针对行健的索引,访问HBase表中的行,只有三种方式:

- 通过单个行健访问
- 通过一个行健的区间来访问
- 全表扫描

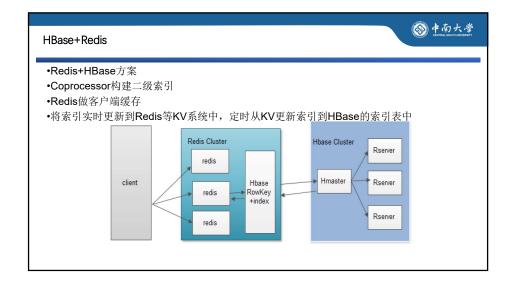
HBase作为列族数据库最经常被人诟病的特性包括:无法轻易建立"二级索引",难以执行求和、计数、排序等操作

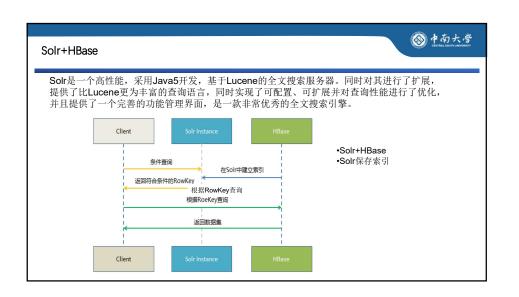
使用其他产品为HBase行健提供索引功能:

- •Hindex二级索引
- •HBase+Redis
- •HBase+solr

原理: 采用HBase0.92版本之后引入的Coprocessor特性(协处理器)













## HBase常用Shell命令

◎ 中南大学

put: 向表、行、列指定的单元格添加数据

一次只能为一个表的一行数据的一个列添加一个数据

scan: 浏览表的相关信息

例子2:继续向表tempTable中的第r1行、第 "f1:c1"列,添加数据值为 "hello,dblab"

```
hbase(main):005:0> put 'tempTable', 'r1', 'f1:c1', 'hello, dblab'
0 row(s) in 0.0240 seconds
hbase(main):006:0> scan 'tempTable'
ROW COLUMN+CELL
r1 column=f1:c1, timestamp=1430036599391, value=hello, dblab
1 row(s) in 0.0160 seconds
```

在添加数据时,HBase会自动为添加的数据添加一个时间戳,当然,也可以在添加数据时人工指定时间戳的值

#### HBase常用Shell命令

•enable/disable: 使表有效或无效

•drop: 删除表

例子4: 使表tempTable无效、删除该表

hbase(main):016:0> disable 'tempTable'
0 row(s) in 1.3720 seconds
hbase(main):017:0> drop 'tempTable'
0 row(s) in 1.1350 seconds
hbase(main):018:0> list
TABLE
testTable
wordcount

2 row(s) in 0.0370 seconds

#### HBase常用Shell命令



get: 通过表名、行、列、时间戳、时间范围和版本号来获得相应单元格的值

- (1) 从tempTable中,获取第r1行、第 "f1:c1" 列的值
- (2) 从tempTable中, 获取第r1行、第 "f1:c3" 列的值

备注: f1是列族, c1和c3都是列

```
hbase(main):012:0> get 'tempTable', 'r1', {COLUMN=>'f1:c1'}
COLUMN
CELL
f1:c1
timestamp=1430036599391, value=hello, dblab
1 row(s) in 0.0090 seconds
hbase(main):013:0> get 'tempTable', 'r1', {COLUMN=>'f1:c3'}
COLUMN
CELL
0 row(s) in 0.0030 seconds
```

从运行结果可以看出: tempTable中第r1行、第 "f1:c3" 列的值当前不存在

## 作业



◎ 中南大学

- 安装HBASE
- •利用Hadoop 提供的Shell 命令完成下列人物
- •列出HABSE中所以表的相关信息,如: 创建时间
- •设计表student, 并建立表
- 先student表添加和删除指定的列族或列

