

DTCC

DATABASE TECHNOLOGY CONFERENCE CHINA 2013 大数据数据库架构与优化数据治理与分析

Internals Calvin Sun

SequeMedia InnoDB © ChinaUnix

Senior Manager, Twitter

Database BDAAS flowingdata DB2 NoSQL MySQL Oracle Big Data

# WHO AM I?

- Calvin Sun (孙春生), email: <u>csun@twitter.com</u>
- Joined Twitter in Mar 2013

2005

MS

- Senior Manager at Oracle, Feb 2008 -Mar 2013
- Manager at MySQL, Jan 2006 Jan 2008

computer science from UST(

 Team lead, product architect at Pervasive Software, Jul 1997 - Dec

## Agenda

- Introduction to InnoDB
- InnoDB Data Format
- InnoDB Logging
- InnoDB Execution
- InnoDB Online Operations
- InnoDB Monitoring & Diagnostics
- Get Involved!









# Introduction to InnoDB





























### InnoDB Timeline









#### MySQL Server Architecture



### NoSQL to InnoDB via Memcached API





#### Fast, simple access to InnoDE

- Accessed via Memcached API
- Use existing Memcached clients
- Bypasses SQL transformations
- SQL/NoSQL access
  - NoSQL for key-value operations
  - SQL for rich queries, JOINs, FKs, etc.
- Implementation
  - Memcached daemon plug-in to mysqld
  - Memcached protocol mapped to the native InnoDB API

 Shared process space for ultralow latency ueMedia

#### InnoDB Features - Transactions

- Full transaction support
  - Atomicity
  - Consistency
  - Isolation
  - Durability
- SQL-standard isolation levels
- Row-level locking
- Multi-version concurrency control (MVCC)
- Automatic deadlock detection
- Plus
  - Automatic crash recovery
  - Referential integrity







#### InnoDB Design Considerations

- Modeled on Gray & Reuter's "Transaction Processing: Concepts & Techniques"
  - Next key locking
- Also emulated the Oracle architecture
  - Multi-version concurrency control (MVCC)
  - Undo info in the database, not the logs
  - Tablespaces for data & index storage
- Added unique subsystems/features







#### InnoDB Innovative Features

- Adaptive Hash Indexes: automatically created on prefix of key for frequent queries
  - Approximates in-memory databases
- Change Buffering: buffers modifications to secondary indexes when the leaf pages are not in the buffer pool
  - Batched merges result in less random access patterns
- Doublewrite Buffer: data first written into the buffer, then flush to the datafiles

reventing partially written pages





# **InnoDB Data Format**









#### InnoDB Database Files





SequeMedia



ChinaUnix

#### InnoDB Tablespaces







#### InnoDB Page Structure

A page consists of: a page header, a page trailer, and a page body (rows or other contents).









Chinal Inix

### InnoDB Compressed Pages



- InnoDB keeps a "modification log" in each page
- Updates & inserts of small records are written to the log w/o page reconstruction; deletes don't even require uncompression
- Log also tells InnoDB if the page will compress to fit page size

• When log space runs out, InnoDB uncompresses the

**2013中国数据库技术大会**page, applies the changes DATABASE TECHNOLOGY CONFERENCE CHINA 2013 ATABASE TECHNOLOGY CONFERENCE CHINA 2013 大数据数据库架构与优化数据治理与分析 and recompresses the page Connaunix

#### InnoDB Row Format

#### InnoDB Row Format:

- Redundant: The oldest InnoDB row format
- **Compact:** The default InnoDB row format since MySQL 5.0.3. It has a more compact representation for nulls and variable-length fields
- **Dynamic**: Store long columns entirely "offpage".
- **Compress**: Compress data & index pages from normal page size (16KB) to specified compressed page size







#### InnoDB Row Structure



# **InnoDB** Logging









#### ARIES

ARIES, *Algorithms for Recovery and Isolation Exploiting Semantics*, is a recovery algorithm used by almost all modern database systems.

Three main principles lie behind ARIES:

- Write ahead logging (WAL)
- Repeating history during Redo
- Logging changes during Undo







#### Types of Logging

- Physical logging: changes to data pages and their data record are logged by their byte offsets and byte-for-byte copies of the data.
- Logical logging: page and byte locations do not matter; only the logical operations are recorded in the recovery log.

### Why Physiological Logging?

- Smaller log files
- Quick recovery
- Proven: It is the method of choice in modern database systems
- Address fundamental flaw in logical logging: operations are not atomic.
  - Example "insert t in T" requires an update to both a data and an index page. A crash might occur after t has been inserted in T but before the index has been updated
  - Example page split







### InnoDB Logging



### InnoDB Redo Log

- Physiological logging
- The redo log remembers EVERY operation on any page in the database
- Redo log record format:

| SpaceID Page | No Offset | OperationType | Changes on | that page |
|--------------|-----------|---------------|------------|-----------|
| Changes      | (only r   | edo values,   | no old     | values)   |

- except for DELETEs, which need no change notes at all
- Examples of operations
  - Insert after record at offset 5444
  - Reorganize page 1234

### InnoDB Undo Log

#### A collection of undo log records

|                   | i |                        |
|-------------------|---|------------------------|
|                   | 1 |                        |
| Brimary Kay Value |   | Old values on that row |
| Frimary Rey value |   | Uld values on that row |
|                   | 1 | 1                      |
|                   | 1 |                        |

- Primary Key Value (no page numbers, no physical addresses)
- Old transaction ID: The ID of the trx that updated that row
- The old field values of that row, which will make the old transaction ID the newest update to the row in question







# **InnoDB** Execution















ChinaUnix

#### Memory Management

- Buffer pool: data pages; index pages; undo records; adaptive hash indexes; table of lock info
- Log buffer: redo records
- Additional memory pool: cached data dictionary; open table handles
- Multiple buffer pool instancesLRU, MRU







### Threads

- User threads (MySQL server threads)
- Master thread
- IO threads
  - read io
  - write io
  - ibuf io
  - log io
- Purge threads
- Page cleaner (flush) thread
- Deadlock detection thread
- FTS, Statistics, Monitor, Drop table, Dump buffer pool, and more















#### InnoDB Transaction Handling









#### Checkpointing

- A checkpoint is a log sequence number (LSN) such that: the data pages in the files contain all changes to the database earlier than LSN
- InnoDB's redo log files have a fixed capacity
- The 'age' of the latest checkpoint must not exceed this capacity

If the checkpoint age would become DICC 2013 + BAREST AC ON OLD AND DB writes the Cold Star Consume

### Flushing

- Activity of writing dirty pages & logs to the disk.
- There are two types of flushing:
  - LRU flushing, based on LRU\_list (roughly ordered on time since last access)
  - Adaptive flushing, based on flush\_list (strictly ordered on oldest\_modification LSN)









### Flushing (cont.)

- Flushing a batch typically involves:
  - Scanning the tail of the relevant list to find victims
  - Select neighbors as candidates for flushing as well
  - Copy dirty pages to the doublewrite buffer
  - Writing double write buffer to disk
  - Sync double write buffer
  - Write to data files
  - Sync all data files







### Purging

- Purge is a type of garbage collection.
- Purge includes:
  - Remove obsolete values from indexes
  - Remove delete marked records that will not be seen by any active transaction
  - Remove the relevant undo entries from history list (a.k.a rollback segment)
- Multi threaded purge: perform purge on a periodic schedule







#### Change Buffer Merging

- Allows changes to secondary index leaf blocks to be deferred when block is not in the buffer.
- Three types of buffering: insert, delete, purge
- Merging
  - Choose random page from the buffer
  - Open a cursor on a random record on that page
  - Read buffer entries from that cursor to find at most 8 pages that should be fetched
  - Issue async IO requests; when async read IO completes, callback is done to apply deferred changes







### Prefetching

- buf\_read\_ahead\_random: before requesting a block read:
  - Count number of blocks in extent that were recently read based on position in buffer pool LRU list.
  - If more than 13 were recently read, prefetch others
- buf\_read\_ahead\_linear: may be used when
  - Accessing first or last page in extent
  - Many pages (56) in extent have been accessed
  - Access pattern was sequential
  - When used, issue read requests for extent that contains that page that follows or precedes







#### Crash Recovery

Commonly, there are 3 phases in Recovery



- There are 4 phases in InnoDB recovery
  - Recover incomplete pages from doublewrite buffer
  - Scan: read redo logs from disk and insert redo log entries into a red-black tree which is sorted on LSN
  - Redo: insert 'dirty' pages into the "flush\_list"

DICCU doulging conference china 2013 DATABASE TECHNOLOGY CONFERENCE CHINA 2013 大数据数据库架构与优化数据治理与分析

# **InnoDB** Online Operations









#### Online Operations in MySQL 5.6 ADD INDEX ADD PRIMARY KEY DROP INDEX ADD COLUMN ADD FOREIGN KEY DROP COLUMN RENAME COLUMN DROP FOREIGN KEY RENAME TABLE ALTER COLUMN NULLABLE ALTER KEY\_BLOCK\_SIZE ALTER COLUMN NOT NULL ALTER ROW FORMAT



SequeMedia



ChinaUnix

#### Type of Online Operations

- Metadata only
  - MySQL Server metadata, such as alter column default
  - MySQL Server metadata & InnoDB metadata, such as add/drop foreign key
- Metadata plus w/o rebuilding the table, such as add/drop index
- Metadata plus rebuilding the table, such as add primary index, add column.







#### How Does It Work?











#### Pre-preparation Phase

#### Server

- Determine the algorithm and concurrency level supported by the storage engine.
- Hold MDL\_SHARED\_UPGRADABLE: allow concurrent DML
- InnoDB

ha\_innobase::check\_if\_supported\_inplace\_alter()

• Check if InnoDB supports a particular alter table in-place.









#### Prepare Phase

#### Server

- Upgrade to MDL\_EXCLUSIVE: no concurrent DML allowed
- Build internal objects describing requested changes

#### InnoDB

ha\_innobase::prepare\_inplace\_alter\_table()

- Check whether the alter is legitimate
- Update internal structures
- Create temporary file(s) for change log(s) due to DMLs







#### Build Phase

#### Server

- Hold MDL\_SHARED\_UPGRADABLE: allow concurrent DML
- Let storage engine to carry out the changes requested by ALTER.

#### InnoDB

ha\_innobase::inplace\_alter\_table()

- Alter the table in-place with operations specified.
- Apply the change logs







### Final Phase

#### Server

- Update .frm and remove old table definitions
- Upgrade to MDL\_EXCLUSIVE: no concurrent DML allowed
- Notify storage engine
- Cleanup internal structures

#### • InnoDB

ha\_innobase:: commit\_inplace\_alter\_table()

- Commit or rollback the changes
  - a) Sync and delete the logs
  - b) Commit metadata changes

DICC 2013 Cleanup internal structures





### Online Add Index

CREATE INDEX index\_name ON table name (column)

|                          | Concurrent User                                 | Source (table)  | (cluster) Index  | Metadata Lock                         |
|--------------------------|---|---|--|---------------------------------------|
| Pre-<br>Prepare<br>Phase | Concurrent Select,<br>Delete, Insert,<br>Update | Check whether<br>the online DDL is<br>supported                               |  | Upgradable<br>Shared Metadata<br>Lock |
| Prepare<br>Phase         | No concurrent DML<br>allowed                    | Create temp table<br>for new index (if<br>primary)                            | Create log files;<br>Logging starts                      | Exclusive<br>Metadata Lock            |
| Build<br>Phase           | Concurrent Select,<br>Delete, Insert,<br>Update | Scan clustered index<br>Extract index entries;<br>Sort / merge index<br>build | DML Logging;<br>Apply logs at the<br>end of create index | Upgradable<br>Shared Metadata<br>Lock |
| Final<br>Phase           | No concurrent DML<br>allowed                    | Drop old table (if<br>primary)  | Update system<br>tables (metadata)                       | Exclusive<br>Metadata Lock            |







# **InnoDB Monitoring & Diagnostics**









# Overview of Monitoring & Diagnostics

- SHOW ENGINE INNODB STATUS
- InnoDB Monitors
- Performance schema for InnoDB
- Information schema tables
  - Information schema metrics table
  - Information schema for InnoDB system tables
  - Information schema for InnoDB buffer pool









#### **INNODB** Show Status and Monitors

- Typical sections of monitor output
  - BACKGROUND THREAD
  - SEMAPHORES
  - LATEST FOREIGN KEY ERROR
  - LATEST DETECTED DEADLOCK
  - TRANSACTIONS
  - FILE I/O
  - LOG
  - BUFFER POOL AND MEMORY
  - ROW OPERATIONS







#### Performance Schema in InnoDB

- 46 mutexes
- 12 rwlocks
- 7 types of threads
- 3 types of I/O (data, log, tmpfile)
- More with debug binaries









#### Performance Schema in InnoDB

Types of running InnoDB threads from THREADS table

mysql> SELECT DISTINCT(name) FROM threads WHERE name LIKE "%innodb%"; +------+ | name | +-----+ | thread/innodb/io\_handler\_thread | | thread/innodb/srv\_lock\_timeout\_thread | | thread/innodb/srv\_error\_monitor\_thread | | thread/innodb/srv\_monitor\_thread | | thread/innodb/srv\_master\_thread | | thread/innodb/srv\_purge\_thread | | thread/innodb/srv\_purge\_thread | +------+ 7 rows in set (0.00 sec)







#### Information Schema Tables

- 29 INFORMATION\_SCHEMA tables in InnoDB
- 9 Data Dictionary related
- 7 FTS related
- 6 Compression related
- 3 Buffer Pool related
- 3 on Locks / Transactions
- 1 General Statistics gold mine (metrics table)







## Information Schema Metrics Table

17 modules, 207 counters

| mysql> select DISTINCT  | subsystem | from                            | innodb_metrics   | order   | by subsyste | em ; |
|---|-----------|---------------------------------|--|---------|-------------|------|
| subsystem   |           |                                 |  |         |             |      |
| <pre>  adaptive_hash_index<br/>  buffer<br/>  buffer_page_io<br/>  change_buffer<br/>  compression<br/>  ddl<br/>  dml<br/>  file_system<br/>  icp<br/>  index<br/>  lock<br/>  metadata<br/>  os<br/>  purge<br/>  recovery<br/>  server<br/>  transaction</pre> |           | mys<br>inn<br>+<br>  c<br>+<br> | <pre>sql&gt; select cour<br/>odb_metrics;<br/>+<br/>count(*)  <br/>+<br/>207  </pre> | nt(*) f | rom         |      |
| 17 rows in set (0.02 se   | ec)       | +<br>1 r                        | +<br>ow in set (0.00   | ) sec)  |             |      |







#### InnoDB System Tables

| Tables_in_information_schema (INNODB_SYS_%) | <br>_ + |
|---|---------|
| INNODB SYS DATAFILES                        | ,       |
| INNODB SYS TABLESTATS                       |         |
| INNODB SYS INDEXES                          |         |
| INNODB SYS TABLES                           |         |
| INNODB SYS FIELDS                           |         |
| INNODB SYS TABLESPACES                      |         |
| INNODB SYS FOREIGN COLS                     |         |
| INNODB SYS COLUMNS                          |         |
| INNODB SYS FOREIGN                          |         |







# **Get Involved!**



















#### 欢迎莅临

Database BDAAS flowingdata DB2 NoSQL MySQL Oracle Big Data

### 2013中国数据库技术大会