



ORACLE®

MySQL & PHP

Johannes Schlüter MySQL Engineering

Johannes Selbsternanntes Urgestein der Münchner PHP UG

- Johannes Schlueter <php@schlueters.de> From:
 - phpuser-Muenchen <phpuser-muenchen@kbx7.de>
 - [phpuser-muenchen] Re: Weihnachtsfeier php User München
 - Fri, 13 Dec 2002 14:13:41 +0100

Date:

Subject:

To:





Johannes?

\$ php -r 'phpcredits();' | grep Johannes

Moriyoshi Koizumi, Xinchen Hui

Alshanetsky, Johannes Schlueter

Schlüter

Andrei Zmievski, Johannes Schlueter

tokenizer => Andrei Zmievski, Johannes Schlueter

- **CLI** => Edin Kadribasic, Marcus Boerger, Johannes Schlueter,
- MySQL driver for PDO => George Schlossnagle, Wez Furlong, Ilia
- MySQLnd => Andrey Hristov, Ulf Wendel, Georg Richter, Johannes
- **Reflection** => Marcus Boerger, Timm Friebe, George Schlossnagle,





MySQL









MySQL Optimizer

- Subquery Optimizations
- File sort optimizations with small limit
 - 4X better execution time 40s to 10s
- Index Condition Pushdown 160X Better execution time – 15s to 90ms
- Postpone Materialization of views/subqueries in FROM - 240X better execution time for EXPLAIN - 8m to 2s
- Batched Key Access and Multi Range Read 280X Better execution time – 2800s to 10s





MySQL Optimizer – Diagnostics and Debugging

- EXPLAIN
 - INSERT, UPDATE, and DELETE
 - JSON format for better readability
- Persistent Optimizer **Statistics - InnoDB**
- Optimizer Traces

query_block





Performance Schema Improvements

- Statements/Stages
 - Most resource intensive queries? Where do they spend time?
- Table/Index I/O, Table Locks
 - Which application tables/indexes cause the most load or contention?
- Users/Hosts/Accounts
 - Which application users/hosts/apps consume the most resources?
- Network I/O
 - Network loaded? How long do sessions idle?
- Summaries
 - Aggregate stats grouped by thread, user, host, account or object





dbahelper

- Views and stored procedures to work with
- Maintained by Mark Leith

cd dbahelper mysql -uroot -p < ./ps_helper_56.sql</pre>

Also for MySQL 5.5 and 5.7

INFORMATION SCHEMA and PERFORMANCE SCHEMA

git clone git@github.com:MarkLeith/dbahelper.git





<pre>mysql> select * from where event_c⁻</pre>	wait_clas lass != '-	sses_global_ idle';	by_avg_laten	су	L
event_class	events	tot_lat	min_lat	avg_lat	max_lat
<pre> wait/io/file wait/io/table wait/io/socket wait/lock/table wait/lock/table wait/synch/rwlock wait/synch/mutex +</pre>	543123 22002 79613 35409 37935 390622	44.60 s 766.60 ms 967.17 ms 18.68 ms 4.61 ms 18.60 ms	19.44 ns 148.72 ns 0 ps 65.45 ns 21.38 ns 19.44 ns	82.11 μs 34.84 μs 12.15 μs 527.51 ns 121.61 ns 47.61 ns	4.21 s 44.97 ms 27.10 ms 969.88 μs 34.65 μs 10.32 μs

(output slightly modified to fit on slide)





<pre>mysql> select * '</pre>	from innodb_buffer_sta	ts_by_table;	L	L	L	L	
object_schema	object_name	allocated	data	pages	hashed	old	cached
InnoDB System	<pre>tinytiny/ttrss_e</pre>	41.95 MiB	33.15 MiB	2685	2685	2685	28595
InnoDB System	tinytiny/ttrss_u	6.88 MiB	4.67 MiB	440	440	440	27510
InnoDB System	piwik/piwik_arch	2.66 MiB	1.36 MiB	170	170	170	3505
InnoDB System	piwik/piwik_arch	2.61 MiB	764.9 KiB	167	167	167	2201
		-					

(output slightly modified to fit, more or less, on slide)







InnoDB







InnoDB in MySQL 5.5

Performance and Scalability

- Multiple buffer pool instances
- Multiple rollback segments
- Improved purge scheduling
- Extended change buffering with delete buffering and purge buffering
- Native async I/O support on Linux
- Improved log sys mutex
- Separate flush list mutex
- Windows performance improvements
- Performance schema for InnoDB





InnoDB in MySQL 5.6 Performance and Scalability

- Split the kernel mutex
- Multi threaded purge
- Use rw locks for page hash
- Add 'page cleaner' thread to flush dirty pages
- Ibuf merge rate improvement
- Configurable data dictionary cache
- InnoDB persistent optimizer statistics
- Read Only Transactions







InnoDB in MySQL 5.6 Monitoring & Diagnostics

- InnoDB Information Schema Metrics Table
- Information schema system tables for InnoDB
- Information schema table for InnoDB buffer pool
- InnoDB: report all deadlocks





InnnoDB Disk Improvements

- SSD optimizations
- Configurable page size (4K, 8K)
- Export/Import per table tablespaces
- Separate InnoDB Undo tablespaces
- Read Only media support
- Compression improvements
- Large (over 4GB) redo logs support





Online ALTER TABLE

- ADD [FULLTEXT] INDEX
- DROP INDEX
- ADD/DROP FOREIGN KEY
- ADD/DROP COLUMN





Key-value Access to InnoDB Data



- Fast, simple access to InnoDB
 - via Memcached API
 - Use existing Memcached clients
 - Bypasses SQL parsing
- NotOnlySQL access
 - For key-value operations
 - SQL for rich queries, JOINs, FKs, etc.
 - Implementation
 - Memcached plug-in to mysqld
 - Memcached mapped to native InnoDB API
 - Shared process for ultra-low latency





InnoDB FullText Search

- Support all query types supported by MyISAM:
 - Natural language search
 - Query expansion
 - Boolean search
- Plus
 - Proximity search

Create full-text index with parallel tokenization and parallel sort





InnoDB FT: Parallel Indexing and Tokenization

		innodb_ft_so	rt_pll_degree	
	1	4	8	16
InnoDB	7 min 48.12 sec	5 min 12.06 sec	4 min 10.95 sec	3 min 33.40 sec
MyISAM		Around 11	min 30 sec	
	(innodb_	ft_sort_pll_degre	e does not affect	MyISAM)

Data Size : Approx 2.7GB, 1 million row , 238 million word Platform: Linux x86 64 bit, 8 cores, 32 GB RAM







Automatic Recovery











Crash-safe master – automatic binary log file trimming

- If the master crashes, binary log will recover automatically from incomplete file flushes.
- On recovery:
 - The active binary log is scanned and any log corruption is detected;
 - Invalid portion of the binary log file is discarded and the file is trimmed.





Crash-safe slave - Slave Info Tables

- Transactional positioning InnoDB based by default. Protection against slave crashes:
 - Automatic recovery.
- Possibility to code multi-source replication in pure SQL.
- Possibility to issue SELECTs on slave information. Automatic conversion between files and tables on startup.





Crash-safe slave - Slave Info Tables

- System tables:
 - slave master_info (master.info)
 - slave_relay_log_info (relay-log.info)



Positional info transactionally stored with the data in tables.









Automatic Switchover/Failover









- Promote any slave to be the new master.
- Automatically pick replication stream correctly:
 - Slave auto-positioning;
 - Slave automatically skips transactions that already exist in the server.







- Two major parts:

 - 2. Server State: set of applied transactions.

1. Logical Identifier: (Server UUID, Transaction Number);

Designed to work with transactional engines (InnoDB).







- The identifier:
 - Abstracts physical positions into logical identifiers human friendly;

 - It is uniquely assigned when the transaction executes; - It is preserved when the transaction is re-applied.
- The server state:
 - Enables slave auto-positioning;
 - Failover facilitator reduced administration overhead;
 - Protects against undesirable transaction re-execution.





- Persistence:
 - It is written to the binary log.
 - Precedes a collection of events that comprise a transaction.





Automatic Switchover/Failover: mysqlfailover

- Automatic failover.

000		Te	erminal — Pyt	hon — 80×24		
MySQL Replica Failover Mode	tion Fai [®] = auto	lover Util Next In	ity nterval =	Wed Apr 4 11	L:29:54 201	. 2
Master Informa	ation					
Binary Log Fi mysql-bin.0000 Replication He	le Pos 901 103! ealth Sta	ition Bin [:] 5 atus	log_Do_DB	Binlog_Ignor	re_DB	F
host	port	role	state	gtid_mode	health	
+ localhost localhost localhost localhost localhost localhost localhost localhost localhost	+ 3310 3311 3312 3313 3314 3315 3316 3317 3318 +	HASTER SLAVE SLAVE SLAVE SLAVE SLAVE SLAVE SLAVE SLAVE	+	0N 0N 0N 0N 0N 0N 0N 0N 0N	+ 0K 0K 0K 0K 0K 0K 0K	F
Q-quit R-refr	esh H-hea	alth G-GTII) Lists U-	-UUIDs L-log (entries Up	Down-scroll

Check and report health at specific intervals in seconds.







Automatic Switchover/Failover: mysqlrpladmin

- General replication administration utility: start, stop, topology health, elect, failover, switchover, gtid.
- On-demand failover or switchover.

			Terminal	l — bash — 87×33	
⁴ Discovering	g slaves –	for master	at local	host:3307	
[#] Checking pr	rivileges				
[#] Performing	switchov	er from ma	ster at l	ocalhost:3307	to slave a
Checking ca	andidate	slave prer	equisites		
[±] Waiting for	r slaves -	to catch u	p to old	master.	
Stopping sl	laves. 👘				
Performing	STOP on	all slaves			
Demoting ol	ld master	to be a s	lave to t	he new master	
Switching s	slaves to	new maste	r.		
Starting al	ll slaves				
Performing	START on	all slave	S.		
Checking sl	laves for	errors.			
	complete				
Switchover	comprete				
: Switchover : Getting hea	alth for i	master: lo	calhost:3	310.	
Switchover Getting hea S	alth for 1	master: lo	calhost:3	310.	
Switchover Getting hea Replicatior	alth for i	master: lo y Health: +	calhost:3	310.	.
Switchover Getting hea Replication host	alth for Topolog +	master: lo y Health: + role	calhost:3 + state	310. + gtid_mode	+ health
Switchover Getting hea Replicatior host localhost	alth for Topolog + port + 3310	master: lo y Health: + role + MASTER	calhost:3 + state +	310. + gtid_mode +	+ health +
Switchover Getting hea Replication host localhost localhost	alth for (Topolog) -+ port -+ 3310 3308	master: lo y Health: + role + MASTER SLAVE	calhost:3 + state + UP UP	310. + gtid_mode + ON ON	+ health + OK OK
Switchover Getting hea Replication host localhost localhost localhost	alth for (Topolog) + port + 3310 3308 3309	master: lo y Health: + role + MASTER SLAVE SLAVE	calhost:3 + state + UP UP UP	310. + gtid_mode + ON ON ON	+ health + OK OK OK OK
Switchover Getting hea Replication host localhost localhost localhost localhost	alth for Topolog + port + 3310 3308 3309 3311	master: lo y Health: + role + MASTER SLAVE SLAVE SLAVE	calhost:3 + state + UP UP UP UP	310. + gtid_mode + ON ON ON ON	+ health + OK OK OK OK
Switchover Getting hea Replication host localhost localhost localhost localhost localhost	alth for (Topolog) + port + 3310 3308 3309 3311 3312	master: lo y Health: + role + MASTER SLAVE SLAVE SLAVE SLAVE	<pre>calhost:3 + state + UP UP UP UP UP UP UP</pre>	310. + gtid_mode + ON ON ON ON ON ON	+ health + OK OK OK OK OK
Switchover Getting hea Replication host localhost localhost localhost localhost localhost localhost	alth for (Topolog) -+ port -+ 3310 3308 3309 3311 3312 3313	master: lo y Health: + role + MASTER SLAVE SLAVE SLAVE SLAVE SLAVE	calhost:3 + state + UP UP UP UP UP UP	310. + gtid_mode + ON ON ON ON ON ON ON	+ health + OK OK OK OK OK OK
Switchover Getting hea Replication host localhost localhost localhost localhost localhost localhost localhost localhost	alth for (Topolog) + port + 3310 3308 3309 3311 3312 3313 3314	master: lo y Health: + role + MASTER SLAVE SLAVE SLAVE SLAVE SLAVE SLAVE	<pre>calhost:3 + state + UP UP UP UP UP UP UP UP</pre>	310. + gtid_mode + ON ON ON ON ON ON ON ON ON	+ health + OK OK OK OK OK OK OK
Switchover Getting hea Replication host localhost localhost localhost localhost localhost localhost localhost localhost localhost	alth for Topolog port port 3310 3308 3309 3311 3312 3313 3314 3314 3315	master: lo y Health: + role + MASTER SLAVE SLAVE SLAVE SLAVE SLAVE SLAVE SLAVE	<pre>calhost:3 + state + UP UP UP UP UP UP UP UP</pre>	310. + gtid_mode + ON ON ON ON ON ON ON ON ON ON	+ health + OK OK OK OK OK OK OK
Switchover Getting hea Replication host localhost localhost localhost localhost localhost localhost localhost localhost localhost	alth for Topolog + port + 3310 3308 3309 3311 3312 3313 3314 3315 3316	master: lo y Health: + role + MASTER SLAVE SLAVE SLAVE SLAVE SLAVE SLAVE SLAVE SLAVE	<pre>calhost:3 + state + UP UP UP UP UP UP UP UP</pre>	310. + gtid_mode + ON ON	+ health + OK OK OK OK OK OK OK OK
Switchover Getting hea Replication host localhost localhost localhost localhost localhost localhost localhost localhost localhost localhost localhost	Topolog Topolog Port Port 3310 3308 3309 3311 3312 3313 3314 3314 3315 3316 3317	master: lo y Health: + role + MASTER SLAVE SLAVE SLAVE SLAVE SLAVE SLAVE SLAVE SLAVE	<pre>calhost:3 + state + UP UP UP UP UP UP UP UP</pre>	310. + gtid_mode + ON ON	+ health + OK OK OK OK OK OK OK OK

hucks-iMac:mysql-wl-6143_cbell\$

ocalhost:3310







MySQL Utitlities

- Originally MySQL Workbench Plugin
- Available under the GPLv2 license
- Written in Python
- Python library to grow solutions for common administrative problems

A collection of Python utilities for managing MySQL databases







MySQL Utilities

— ...

- - mysqldbcompare compare databases
 - mysqldbcopy copy databases between servers
 - mysqlfailover Automatic fail-over
 - mysqlrpladmin General replication administration utility
 - mysqlrplshow show a graph of your topology
 - mysqlreplicate setup replication
 - mysqlrplcheck check replication configuration
- Build your own tools on top of the core of the library

Easily administer MySQL servers from the command line









An integrated environment for managing a farm of MySQL server supporting high-availability and sharding.

http://labs.mysql.com









MySQL Fabric

- "Farm" Management System
- Distributed
- High-Availability
- Sharding
- Procedure Executor
- Extensible

Written in Python

- Early alpha
 - Long road ahead
- You can participate
 - Suggest features
 - Report bugs
 - Contribute patches
- MySQL 5.6 is focus







MySQL Fabric

- Decision logic in connector
 Shard Multiple Tables Reducing network load
- Support Transactions API to provide sharding key
- Global Updates
 - Global Tables
 - Schema updates
- Procedure Executor

- Using same key
- Sharding Functions
 - Range
 - (Consistent) Hash
- Shard Operations
 - Using built-in executor
 - Shard move
 - Shard split







PHP: mysqlnd_ms

- Plugin for PHP for doing load-balancing
- ms originally meant "master / slave"
- http://dev.mysql.com/doc/refman/5.5/en/apis-phpbook.mysqlnd-ms.html
- Uses GTIDs to ensure sending users to a slave which received data we wrote before
- 1.6.0-alpha added MySQL Fabric Support







mysqlnd_ms 1.6.0 configuration









mysqlnd_ms 1.6.0

<?php \$c = new mysqli("test", "root", "", "test");

echo "Creating global table:\n"; mysqlnd_ms_fabric_select_global(\$c, "test.fabrictest"); \$c->query("CREATE TABLE fabrictest (id INT NOT NULL)");

echo "Inserting with ID 10:\n"; mysqlnd_ms_fabric_select_shard(\$c, "test.fabrictest", 10); \$c->query("INSERT INTO fabrictest VALUES (10)"); ?>









PHP's MySQL Architecture













API Choice

- mysqli
 - Support for all MySQL features
 - Best support / stability
 - Integration with existing applications / environments

PDO MYSQL

 Simple applications supporting multiple databases

-API-compatibility is often not enough, though

 Integration with existing applications / environments







Three kinds of PHP Users

Users of existing Applications

Take existing applications (i.e. Wordpress, phpBB, moodle, ...)

Little customization only

Often hard to scale

Ad-hoc developed applications

Code often mixture of PHP with embedded SQL etc.

Hard to scale/adapt to new situations

Quick'n'Dirty Applications

Framework-based Applications

Custom applications built on i.e. Symfony, Zend Framework or CakePHP

Built on DB abstraction Layers

Abstractions hide full power of MysQL







On-Going Demands



Improve Security!

Better performance!

More Scalability!

Higher Availability!

Refactor the application!











Solving these requests











mysqlnd in depth

mysqlnd_result:fetch_row



mysqlnd net:send



mysqlnd_net:read_result

mysqlnd

mnd::allocate

mysqlnd:prepare

mysqInd:query







myslqnd Plugins

- Loaded as regular PHP extension
 - Written in C
- Hook into mysqlnd in any level
 - From high-level APIs (query, prepare fetch row, etc)
 - To low-level (network IO, memory allocation)
- Can be 100% transparent
- Can provide userspace APIs which can work with ext/mysql, mysqli and PDO mysql instances







What can Plugins be used for?

- Read/Write Splitting
- Failover Round-Robin









Existing Plugins

- Stable open source:
 - mysqlnd_uh Userspace Handler
 - mysqlnd ms Repliation and load balancing
 - mysqlnd qc Client Side Query Cache
 - mysqlnd_memcache memache mapping
- Experimental openn source:
 - mysqlnd mc multi connect
 - mysqlnd_pscache prepared statement cache
 - mysqlnd_sip SQL injection protection
- Commercial
 - MySQL Enterprise Query Analyzer









Multi-Threaded Slaves









The problem



Replication system architecture: multi-threaded master vs single threaded slave applier







The problem

- Multicore on the Master and the Slave
- Recent progress of scaling-up the Innodb engine
- Slave side needs to match an upcoming solution for group commit to the binary log to overcome deficiency of the current architecture







Facilities & Prerequisites

- The user application is logically partitioned per database
- The user application is content with local database consistency



ORACLE





Requirements

- Arbitrary manually configurable number of Worker threads
- All binary log formats supported
- All storage engines supported
- Automatic fallback to sequential execution for statements that do not comply with parallel execution Automatic recovery in face of crashes in the case of the
- Innodb engine







Overview of architecture

- SQL thread is split into Coordinator and Worker threads communicating in style of the producer-consumer model
- Coordinator activities include: read replication event, handle Skip-Until-Delay options and find an appropriate Worker to execute the rest.
- Worker activities include:
- apply the event including the commit of transaction. Coordinator and Worker cooperate on recovery provision







MySQL 5.7 Development Milestones









MySQL 5.7 Development Milestone Release 2

500,000 Queries/Second

Download Now »









MySQL 5.7 **Development Milestone Release 3 1** Million Queries/Second

Download Now »









MySQL 5.7 – OLTP_RO Point-Selects 8-tables • UNIX socket, sysbench 0.4.8

Query/sec









MySQL 5.7 – Connections per Second









MySQL 5.7 – Multi-Source-Replication



MSQL 5.7 - EXPLAIN

- Problem: A statement in a session is taking too long
- New option:
 - EXPLAIN FOR CONNECTION from other session

EXPLAIN [FORMAT=(JSON|TRADITIONAL)] FOR CONNECTION id

MySQL 5.7 – InnoDB

- Improved Online ALTER TABLE
 - Online RENAME INDEX
 - Online change VARCHAR
- Improved InnoDB Temp Tables
 - New separate tablespace for temp tables
 - Improved CREATE/DROP performance for temp tables
- Optimized temp table DML - No REDO logging, no change bffering, no locking

Summary

With MySQL 5.6 and 5.7 we want to be

- ... more scalable
- ... faster
- ... better obsersvable
- ... more stable
- ... help operations

The world's most popular open source database

Resources

- Blog from MySQL Engineering VP Tomas Ulin http://insidemysql.com/
- Group Blog of the MySQL Server Team http://mysqlserverteam.com
- Ulfs Blog (mysqlnd and other things) http://blog.ulf-wendel.de
- Benchmark Details by DimitriK http://dimitrik.free.fr

Johannes Schlüter

johannes.schlueter@oracle.com Twitter: @phperror Blog: http://schlueters.de/blog

The preceding is intended to outline our general product direction. It is intended for information purposes only, and may not be incorporated into any contract.

It is not a commitment to deliver any material, code, or functionality, and should not be relied upon in making purchasing decisions. The development, release, and timing of any features or functionality described for Oracle's products remains at the sole discretion of Oracle.

