# How to Avoid Pitfalls in Schema Upgrade

with Galera

**February 2, 2020**

**Sveta Smirnova**

PERCONA

# Sveta Smirnova



- MySQL Support engineer
- Author of
  - MySQL Troubleshooting
  - JSON UDF functions
  - FILTER clause for MySQL
- Speaker
  - Percona Live, OOW, Fosdem, DevConf, HighLoad…

# Table of Contents

- TOI

- RSU

- pt-online-schema-change (pt-osc)

PERCONA

# Introduction

- Galera Replication Library
  - Provides synchronous replication for MySQL

PERCONA

# Introduction

- Galera Replication Library
  - Provides synchronous replication for MySQL
- Galera Clusters
  - MariaDB Galera Cluster
  - Percona XtraDB Cluster
  - Galera Cluster for MySQL

PERCONA

# How Galera works

- Data modification happens on a node
- Optimistic locking control

# How Galera works

- Data modification happens on a node
- Optimistic locking control
- At the `COMMIT` time
  - Broadcasts write set for the cluster
  - Waits confirmation of the successful update
    - From all other nodes

Yes Commits transaction locally
No Rollbacks transaction

PERCONA

# Data Updates

- Committed on all nodes or nowhere
- Safe

PERCONA

# Challenges of DDL

- Replicated independently from storage engine

PERCONA

# Challenges of DDL

- Replicated independently from storage engine
- Changes may affect query results
  - Adding/removal of `UNIQUE` keys
  - Adding/removal columns
  - Changing column definition

PERCONA

# Challenges of DDL

- Replicated independently from storage engine
- Changes may affect query results
- Modification can happen on any node
  - The schema must be upgraded before DML
  - There is no way to rollback schema upgrade
  - MDLs are set only on one node
    - Not across the cluster
    - Not possible to rely on them for all nodes
    - Additional control required

PERCONA

TOI

# Total Order Isolation (TOI)

- DDL changes are replicated in the same order regarding other transactions
- All nodes are in the absolutely same state at any point of time

PERCONA

# TOI: Illustration

- 3-nodes cluster
  - Node A
  - Node B
  - Node C

PERCONA

# TOI: Illustration

- Initial state

| Node A | Node B | Node C |
|--------|--------|--------|
| `INSERT(103)` | `UPDATE(101)` | `SELECT(100)` |
| `UPDATE(104)` | `INSERT(102)` | `INSERT(112)` |
| `ALTER(105)` | `DELETE(108)` | `SELECT(113)` |
| | `UPDATE(109)` | `UPDATE(114)` |

PERCONA

# TOI: Illustration

- Queries status

**Node A**
- ► `INSERT(103)`
- ► `UPDATE(104)`
- 🕐 `ALTER(105)`

**Node B**
- ► `UPDATE(101)`
- ► `INSERT(102)`
- 🕐 `DELETE(108)`
- 🕐 `UPDATE(109)`

**Node C**
- ► `SELECT(100)`
- 🕐 `INSERT(112)`
- ↻ `SELECT(113)`
- 🕐 `UPDATE(114)`

PERCONA

# TOI: Illustration

- ALTER in progress

**Node A**
▶ `ALTER(105)`

**Node B**
🕐 `DELETE(108)`
🕐 `UPDATE(109)`

**Node C**
🕐 `INSERT(112)`
↻ `SELECT(113)`
🕐 `UPDATE(114)`

PERCONA

# TOI: Illustration

- ALTER finished

| Node A | Node B | Node C |
|--------|--------|--------|
| | ► DELETE(108) | ► INSERT(112) |
| | ► UPDATE(109) | ► SELECT(113) |
| | | ► UPDATE(114) |

**PERCONA**

# PROCESSLIST: DML before ALTER

```
DML node> select DB, COMMAND, TIME, STATE, INFO from information_schema.processlist WHERE DB='sbtest';
+--------+---------+------+-----------------------------------------------+----------------------+
| DB     | COMMAND | TIME | STATE                                         | INFO                 |
+--------+---------+------+-----------------------------------------------+----------------------+
| sbtest | Query   | 1    | wsrep: initiating pre-commit for write set (2886) | COMMIT           |
| sbtest | Query   | 1    | wsrep: initiating pre-commit for write set (2888) | COMMIT           |
| sbtest | Query   | 1    | wsrep: initiating pre-commit for write set (2884) | COMMIT           |
| sbtest | Query   | 1    | updating                                      | DELETE FROM sbtest1.. |
| sbtest | Query   | 1    | wsrep: initiating pre-commit for write set (2887) | COMMIT           |
| sbtest | Query   | 0    | wsrep: initiating pre-commit for write set (2889) | COMMIT           |
| sbtest | Query   | 1    | wsrep: initiating pre-commit for write set (2885) | COMMIT           |
| sbtest | Query   | 1    | wsrep: pre-commit/certification passed (2883)  | COMMIT              |
+--------+---------+------+-----------------------------------------------+----------------------+
8 rows in set (0.00 sec)
```

PERCONA

# PROCESSLIST: SELECT before ALTER

```
SELECT node> select DB, COMMAND, TIME, STATE, INFO from information_schema.processlist
        -> WHERE DB='sbtest';
+--------+---------+------+------------------+-----------------------------------------+
| DB     | COMMAND | TIME | STATE            | INFO                                    |
+--------+---------+------+------------------+-----------------------------------------+
| sbtest | Query   | 0    | statistics       | SELECT pad FROM sbtest2 WHERE id=5009   |
| sbtest | Query   | 0    | starting         | SELECT pad FROM sbtest3 WHERE id=4951   |
| sbtest | Query   | 0    | statistics       | SELECT pad FROM sbtest4 WHERE id=4954   |
| sbtest | Query   | 0    | System lock      | SELECT pad FROM sbtest2 WHERE id=5351   |
| sbtest | Query   | 0    | cleaning up      | SELECT pad FROM sbtest2 WHERE id=4954   |
| sbtest | Sleep   | 0    |                  | NULL                                    |
| sbtest | Query   | 0    | Sending to client| SELECT pad FROM sbtest1 WHERE id=4272   |
| sbtest | Query   | 0    | closing tables   | SELECT pad FROM sbtest4 WHERE id=4722   |
+--------+---------+------+------------------+-----------------------------------------+
8 rows in set (0.00 sec)
```

PERCONA

# ALTER

```
DDL node> use ddltest;
Reading table information for completion of table and column names
You can turn off this feature to get a quicker startup with -A

DDL node> alter table sbtest1 add key k1(c, k, pad);
Query OK, 0 rows affected (3 min 53.73 sec)
Records: 0 Duplicates: 0 Warnings: 0
```

PERCONA

# PROCESSLIST: DML during ALTER

```
DML node> select DB, COMMAND, TIME, STATE, INFO from information_schema.processlist
       -> WHERE DB in ('sbtest','ddltest');
+---------+---------+------+------------------------------------------------------+---------------------+
| DB      | COMMAND | TIME | STATE                                                |                     |
+---------+---------+------+------------------------------------------------------+---------------------+
| sbtest  | Query   | 36   | wsrep: initiating pre-commit for write set (7886)    | COMMIT              |
| sbtest  | Query   | 37   | wsrep: initiating pre-commit for write set (7882)    | COMMIT              |
| sbtest  | Query   | 27   | wsrep: initiating pre-commit for write set (7887)    | COMMIT              |
| sbtest  | Query   | 27   | wsrep: initiating pre-commit for write set (7888)    | COMMIT              |
| sbtest  | Query   | 36   | wsrep: initiating pre-commit for write set (7885)    | COMMIT              |
| sbtest  | Query   | 37   | wsrep: initiating pre-commit for write set (7883)    | COMMIT              |
| sbtest  | Query   | 37   | wsrep: initiating pre-commit for write set (7884)    | COMMIT              |
| sbtest  | Query   | 10   | wsrep: initiating pre-commit for write set (7889)    | COMMIT              |
| ddltest | Sleep   | 38   | altering table                                       | alter table sbtest1.|
+---------+---------+------+------------------------------------------------------+---------------------+
9 rows in set (0.00 sec)
```

PERCONA

# PROCESSLIST: SELECT during ALTER

```
SELECT node> select DB, COMMAND, TIME, STATE, INFO from information_schema.processlist
         -> WHERE DB in ('sbtest','ddltest');
+---------+---------+------+------------------+--------------------------------------------+
| DB      | COMMAND | TIME | STATE            |                                            |
+---------+---------+------+------------------+--------------------------------------------+
| sbtest  | Sleep   | 0    |                  | NULL                                       |
| sbtest  | Sleep   | 0    |                  | NULL                                       |
| sbtest  | Query   | 0    | Sending to client| SELECT pad FROM sbtest4 WHERE id=4989      |
| sbtest  | Sleep   | 0    |                  | NULL                                       |
| sbtest  | Query   | 0    | query end        | SELECT pad FROM sbtest2 WHERE id=4961      |
| sbtest  | Sleep   | 0    |                  | NULL                                       |
| sbtest  | Sleep   | 0    |                  | NULL                                       |
| sbtest  | Sleep   | 0    |                  | NULL                                       |
| ddltest | Sleep   | 39   | altering table   | alter table sbtest1 add key k1(c, k, pad)  |
+---------+---------+------+------------------+--------------------------------------------+
9 rows in set (0.14 sec)
```

PERCONA

# TOI Advantages

- Data always consistent
- DDL applied to all nodes at the same time
- No failure due to schema inconsistency

PERCONA

# TOI Disadvantages

- The whole cluster blocked
  - For the duration of the entire DDL operation
- Schema upgrades replicated as a statement
  - There is no guarantee that the ALTER succeed!

PERCONA

# How to Perform Upgrade with TOI

- Schedule maintenance window
- Run DDL
- Cluster won't be accessible until DDL finishes
  - SELECTs can continue
  - `wsrep_sync_wait != 1`

PERCONA

# When to Use TOI

- Quick DDL operations

PERCONA

# When to Use TOI

- Quick DDL operations
- Creating new database objects
  - `CREATE DATABASE`
  - `CREATE TABLE`

PERCONA

# When to Use TOI

- Quick DDL operations
- Creating new database objects
- Online operations which modify metadata only
  - `RENAME INDEX`
  - `RENAME TABLE`
  - `DROP INDEX`
  - `ALGORITHM=INSTANT`
  - Full list

PERCONA

# RSU

# Rolling Schema Upgrade (RSU)

- Variable `wsrep_OSU_method`
- Puts node into de-sync state
  - For the duration of DDL
- Pauses Galera provider
- Schema can get out of sync!

PERCONA

# User Responsibility

- Run DDL on the each node of the cluster
- Block read-write access that depend on DDL
  - Until all nodes are in sync
- Make sure no write is performed to the table
  - Until upgrade finishes on all nodes
- Failure makes cluster unrecoverable!

PERCONA

# RSU Workflow

- User Action
- `SET SESSION wsrep_OSU_method = 'RSU';`
- DDL


- Any other statement

- Node Operation
- Nothing
- Is `wsrep_OSU_method` set to RSU?

Yes Performs DDL
- Nothing

PERCONA

# How Node Internally Executes DDL in RSU Mode?

▼ Does node have transactions in `COMMIT` mode?

# How Node Internally Executes DDL in RSU Mode?

▼ Does node have transactions in `COMMIT` mode?
Yes Wait for 5 milliseconds

PERCONA

# How Node Internally Executes DDL in RSU Mode?

▼ Does node have transactions in `COMMIT` mode?

Yes Wait for 5 milliseconds

▼ Still transactions in the `COMMIT` mode exist?

# How Node Internally Executes DDL in RSU Mode?

▼ Does node have transactions in `COMMIT` mode?
Yes Wait for 5 milliseconds
▼ Still transactions in the `COMMIT` mode exist?
Yes Abort DDL

PERCONA

# How Node Internally Executes DDL in RSU Mode?

▼ Does node have transactions in `COMMIT` mode?
No Put node into de-sync state

PERCONA

# How Node Internally Executes DDL in RSU Mode?

▼ Does node have transactions in `COMMIT` mode?

No Put node into de-sync state

▼ Pause write-set application

PERCONA

# How Node Internally Executes DDL in RSU Mode?

▼ Does node have transactions in `COMMIT` mode?

No Put node into de-sync state

▼ Pause write-set application

▼ Execute DDL

PERCONA

# How Node Internally Executes DDL in RSU Mode?

▼ Does node have transactions in `COMMIT` mode?

No Put node into de-sync state

▼ Pause write-set application

▼ Execute DDL

▼ Bring the node back to the cluster

PERCONA

# How Node Internally Executes DDL in RSU Mode?

▼ Does node have transactions in `COMMIT` mode?

No Put node into de-sync state

▼ Pause write-set application

▼ Execute DDL

▼ Bring the node back to the cluster

● Synchronize

P E R C O N A

# RSU: Locking

- Not avoidable
- Updates to all objects on the node in RSU mode must finish before the operation
- Failure aborts DDL

PERCONA

# RSU Advantages

- Cluster remains functional
- Schedule long-running `ALTER`
  - In the best time possible

PERCONA

# RSU Disadvantages

- No checks for data and schema consistency
  - This is your responsibility!

PERCONA

# RSU Disadvantages

- No checks for data and schema consistency
- All writes must be stopped on the affected node
  - Otherwise DDL fails with an error

**PERCONA**

# RSU Disadvantages

- No checks for data and schema consistency
- All writes must be stopped on the affected node
- gcache should be big enough to hold changes
  - Made while DDL was running
  - Failure will cause SST when node re-joins cluster
  - All schema changes will be lost

PERCONA

# RSU Disadvantages

- No checks for data and schema consistency
- All writes must be stopped on the affected node
- gcache should be big enough to hold changes
- Any error can make cluster dysfunctional

PERCONA

# RSU Disadvantages

- No checks for data and schema consistency
- All writes must be stopped on the affected node
- gcache should be big enough to hold changes
- Any error can make cluster dysfunctional
- Affected table must be offline
  - Until the schema upgrade is done on all nodes
  - Unless this is schema-compatible change

PERCONA

# How to Use RSU

- Make sure gcache is big enough
  - Must hold all updates while DDL is in progress

PERCONA

# How to Use RSU

- Make sure gcache is big enough
  - Must hold all updates while DDL is in progress
- Block all writes to the table/schema

PERCONA

# How to Use RSU

↻ Choose an "upgrading node"

PERCONA

# How to Use RSU

↺ Choose an "upgrading node"
↺ Block all write requests to this node

PERCONA

# How to Use RSU

↻ Choose an "upgrading node"
↻ Block all write requests to this node
↻ `SET SESSION wsrep_OSU_method = 'RSU';`

PERCONA

# How to Use RSU

- ↻ Choose an "upgrading node"
- ↻ Block all write requests to this node
- ↻ `SET SESSION wsrep_OSU_method = 'RSU';`
- ↻ Perform DDL in the same session

PERCONA

# How to Use RSU

↻ Choose an "upgrading node"
↻ Block all write requests to this node
↻ `SET SESSION wsrep_OSU_method = 'RSU';`
↻ Perform DDL in the same session
↻ `SET SESSION wsrep_OSU_method = 'TOI';`

**PERCONA**

# How to Use RSU

- ↻ Choose an "upgrading node"
- ↻ Block all write requests to this node
- ↻ `SET SESSION wsrep_OSU_method = 'RSU';`
- ↻ Perform DDL in the same session
- ↻ `SET SESSION wsrep_OSU_method = 'TOI';`
- ↻ Re-enable writes

PERCONA

# How to Use RSU

- ↻ Choose an "upgrading node"
- ↻ Block all write requests to this node
- ↻ `SET SESSION wsrep_OSU_method = 'RSU';`
- ↻ Perform DDL in the same session
- ↻ `SET SESSION wsrep_OSU_method = 'TOI';`
- ↻ Re-enable writes
- ↻ Repeat for other nodes

PERCONA

# pt-online-schema-change (pt-osc)

# pt-online-schema-change (pt-osc)

- A tool, performing non-blocking upgrades
  - With TOI

PERCONA

# pt-online-schema-change (pt-osc)

- A tool, performing non-blocking upgrades
- Creates a copy of table with altered definition

PERCONA

# pt-online-schema-change (pt-osc)

- A tool, performing non-blocking upgrades
- Creates a copy of table with altered definition
- Creates triggers which will copy modified rows

PERCONA

# pt-online-schema-change (pt-osc)

- A tool, performing non-blocking upgrades
- Creates a copy of table with altered definition
- Creates triggers which will copy modified rows
- Starts copying data in chunks
  - Absolutely under control
  - Can be paused or stopped

–max-flow-ctl

**PERCONA**

# pt-online-schema-change (pt-osc)

- A tool, performing non-blocking upgrades
- Creates a copy of table with altered definition
- Creates triggers which will copy modified rows
- Starts copying data in chunks
  - All rows already in the table are copied in chunks
  - Newly modified rows are copied using triggers

PERCONA

# pt-online-schema-change (pt-osc)

- A tool, performing non-blocking upgrades
- Creates a copy of table with altered definition
- Creates triggers which will copy modified rows
- Starts copying data in chunks
- Once copy is complete, drops the table

**PERCONA**

# pt-online-schema-change (pt-osc)

- A tool, performing non-blocking upgrades
- Creates a copy of table with altered definition
- Creates triggers which will copy modified rows
- Starts copying data in chunks
- Once copy is complete, drops the table
- Renames the copy into the original table name

PERCONA

# pt-osc Advantages

- DDL is safe and non-blocking

PERCONA

# pt-osc Disadvantages

- Works only with InnoDB tables
- Increases IO load even for inplace operations
- Conflicts with already existing triggers
  - Unless you use MariaDB $>=$ 10.2.3
- Foreign keys updates are not effectively safe

PERCONA

# How to Use pt-osc

- Study pt-osc options
  - `--max-flow-ctl`
- Set appropriate limits
- Make sure `wsrep_OSU_method` is TOI
- Run pt-osc

PERCONA

# Which Method to Use?

▼ Will DDL be fast?
- `CREATE DATABASE`
- `CREATE TABLE`
- `DROP INDEX`
- Any `ALTER` on small tables
- Other

**PERCONA**

# Which Method to Use?

▼ Will DDL be fast?

Yes Use TOI

PERCONA

# Which Method to Use?

▼ Will DDL be fast?

Yes Use TOI

No Evaluate if you can use pt-osc
- Operation on the InnoDB table
- Table has no triggers or MariaDB $>= 10.2.3$
- Table is not referenced by a foreign key
- You can tolerate increased IO

PERCONA

# Which Method to Use?

  ▼ Will DDL be fast?
Yes Use TOI
 No Evaluate if you can use pt-osc
Yes Use pt-osc

PERCONA

# Which Method to Use?

▼ Will DDL be fast?
Yes Use TOI
No Evaluate if you can use pt-osc
Yes Use pt-osc
No Use RSU
   • Stop all write traffic on the node
   • Stop all write traffic to the modified table
   • Make sure to upgrade on all nodes

PERCONA

# Conclusion

- Use TOI whenever possible
- Then use pt-osc
- RSU is a last resort

PERCONA

# More information

Galera Cluster

MariaDB Galera Cluster

pt-online-schema-change

PERCONA

# Thank you!

www.slideshare.net/SvetaSmirnova

twitter.com/svetsmirnova

github.com/svetasmirnova

PERCONA

DATABASE PERFORMANCE MATTERS