

# MySQL Servers Working as a Team - Replication or Galera Cluster

Open Source Data Center Conference, April 26<sup>th</sup> - 28<sup>th</sup>, Berlin

**Jörg Brüche**

Senior Support Engineer, FromDual GmbH

[joerg.bruehe@fromdual.com](mailto:joerg.bruehe@fromdual.com)



CC-BY-SA



## Support



## Consulting



## remote-DBA



## Training



- **Development distributed SQL-DBMS**  
Porting mainframe -> Unix,  
Interface to archiver tools (ADSM, NetWorker)
- **MySQL Build Team**  
Release builds incl. tests, packaging, scripts, ...
- **DBA**  
MySQL running a web platform  
(master-master-replication)
- **Support-Engineer (FromDual)**  
Support + Remote-DBA for MySQL / MariaDB / Percona  
both with and without Galera Cluster

# Contents

MySQL Server: Architecture

Binlog

Replication

Galera Cluster

Comparison

Examples / When (not) to Choose Which

# General Remarks



- **Concepts rather than details:  
"the forest, not the trees"**
- **MySQL 5.6 (established GA version)**
- **Also valid for Percona and MariaDB**
  
- **Not applicable to "embedded" MySQL**
- **Not considered: NDB = "MySQL Cluster"**

## ➔ **MySQL Server: Architecture**

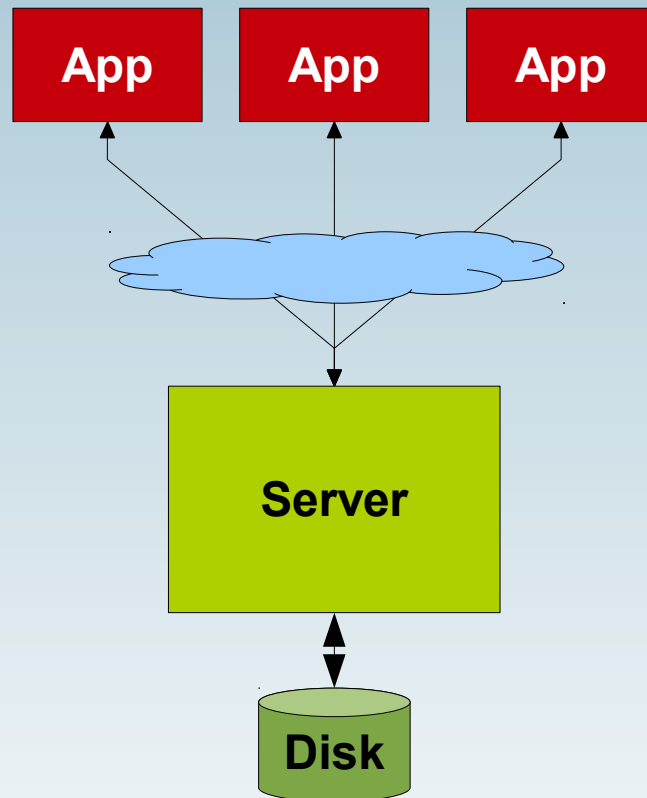
Binlog

Replication

Galera Cluster

Comparison

Examples / When (not) to Choose Which



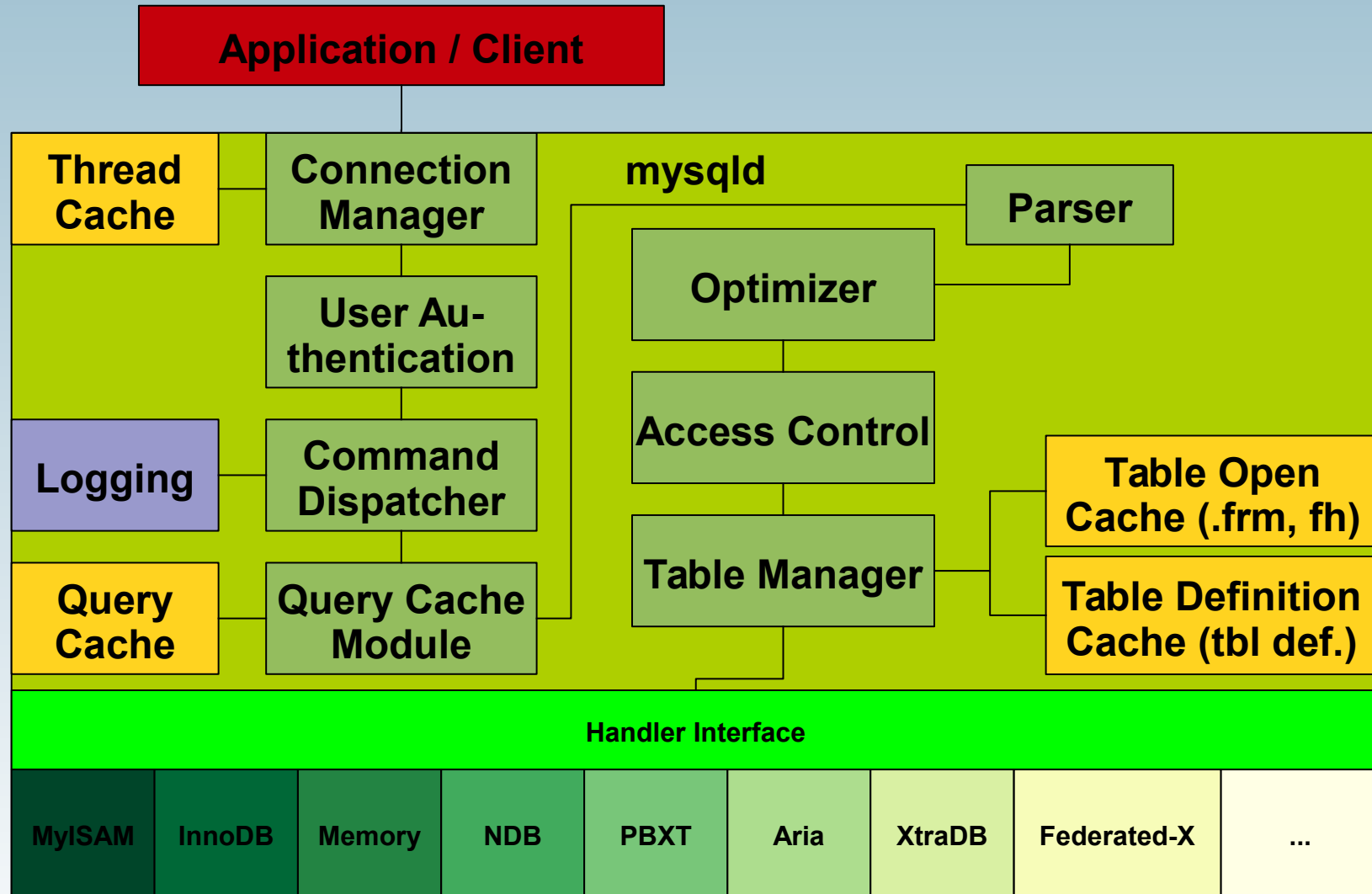
**Client (application)  
local or remote**

**Socket, LAN or internet**

**Server is separate process,  
multi-threaded:  
1 thread per user session**

**Disk / SSD, local or SAN**

# Inside the Server





# MySQL Server: Architecture

## ➔ **Binlog**

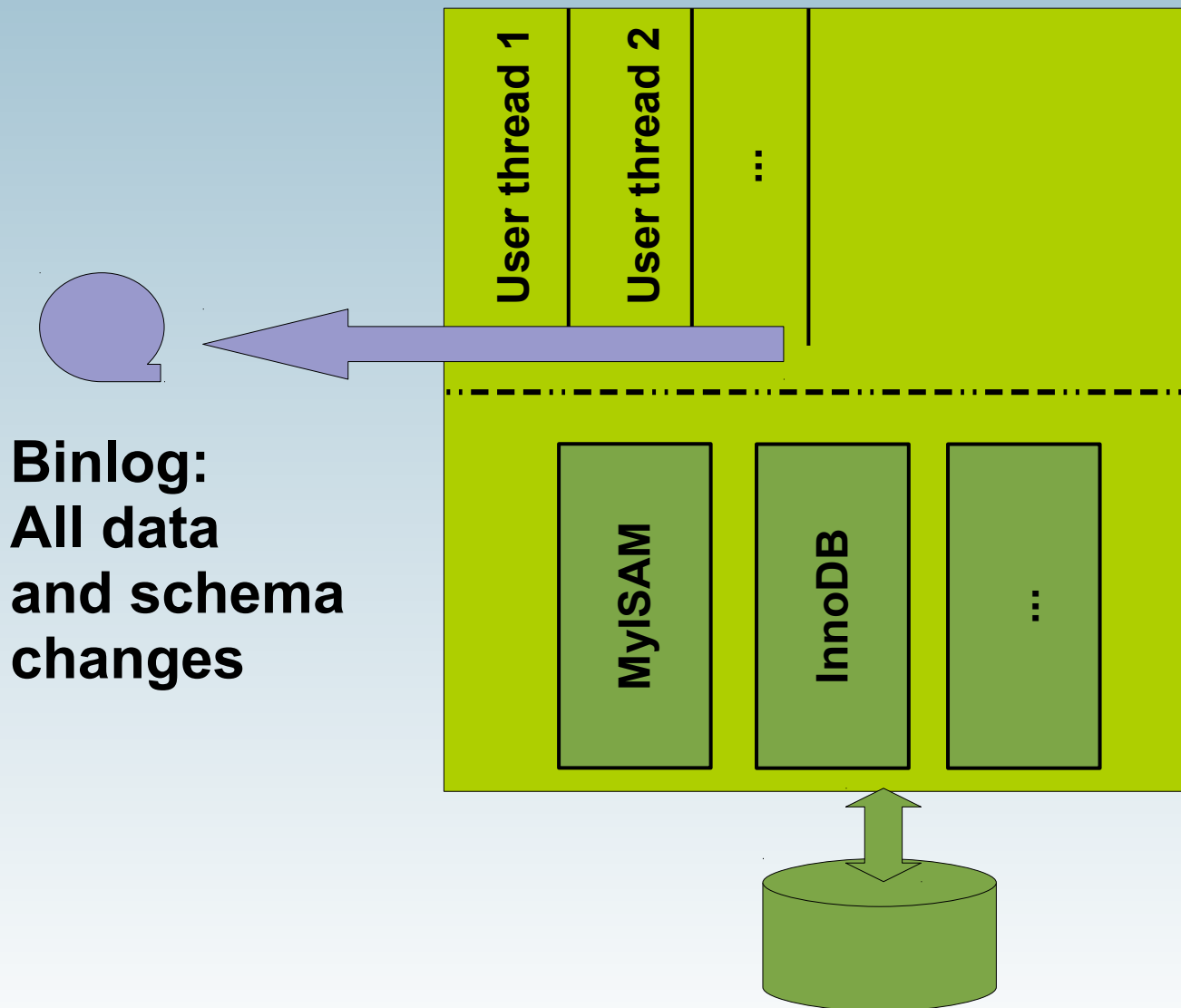
Replication

Galera Cluster

Comparison

Examples / When (not) to Choose Which

# Layers + Binlog



**Binlog:**  
All data  
and schema  
changes

## SQL layer:

- Parser
- Optimizer
- Privileges
- Query Cache
- ...

## Handler Interface

## File layer:

- Table Handler
- InnoDB:
  - Row Access
  - Row Locks
  - Recovery
- ...

- **All data changes executed**
- **All schema changes executed**
- **Timestamps**
- **Essential for Point-in-Time-Recovery "PITR"**
- **Independent of table handler**
- **Formats "statement", "row", and "mixed"**
- **Segments of configurable size**
- **Numbered sequentially**

# MySQL Server: Architecture

Binlog

➔ **Replication**

Galera Cluster

Comparison

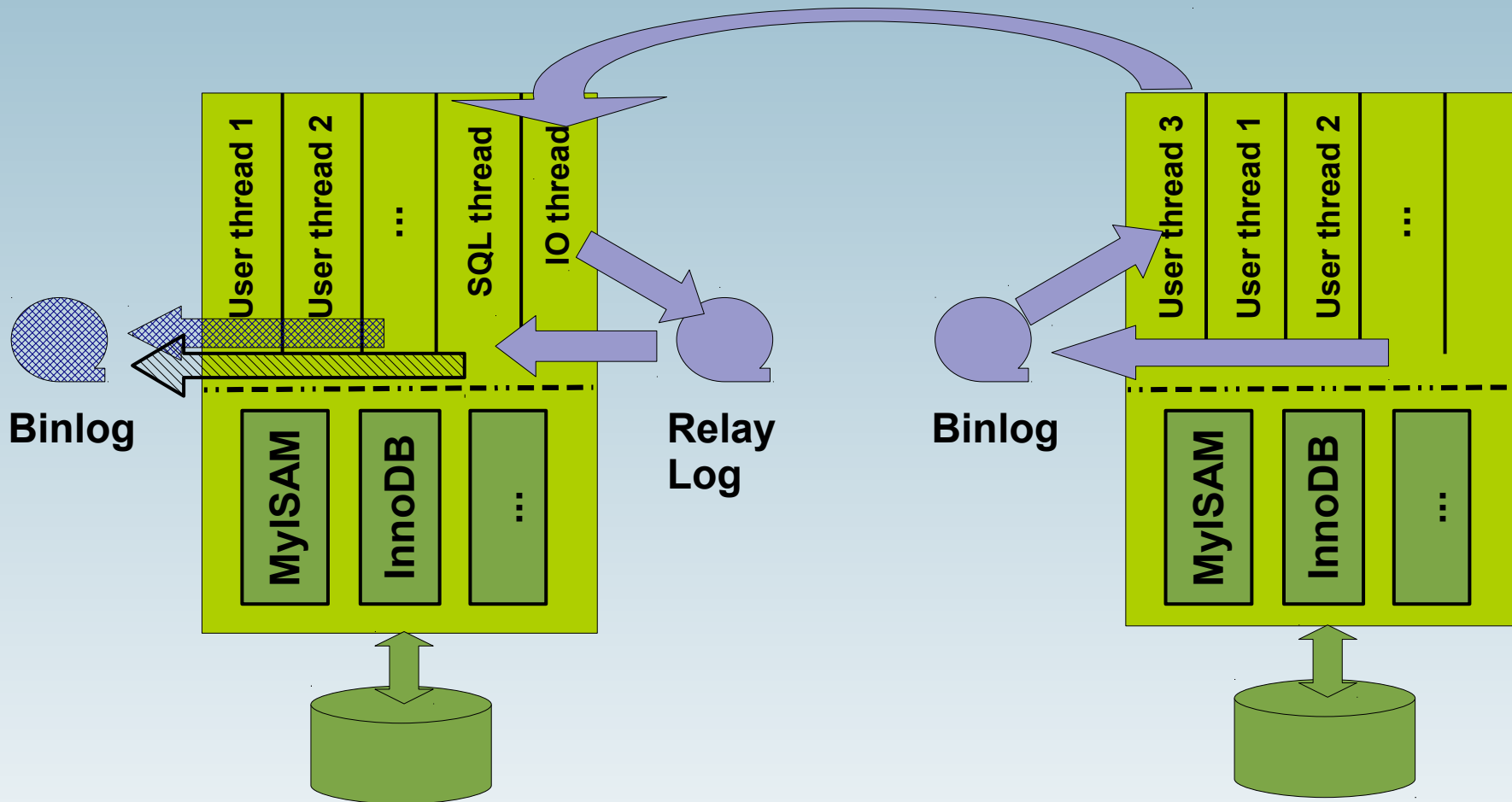
Examples / When (not) to Choose Which

# MySQL Replication



- Applications communicate with "Master"
- "Master" logs all changes
- "Slave" has identical initial state
- Slave fetches all changes from master and applies them locally
- Replication is running asynchronous
- Slave stops replication on difference

# Slave fetches binlog



**Slave:**  
“log-bin = FILE”, or else no binlog  
“log\_slave\_updates = 1” for forwarding

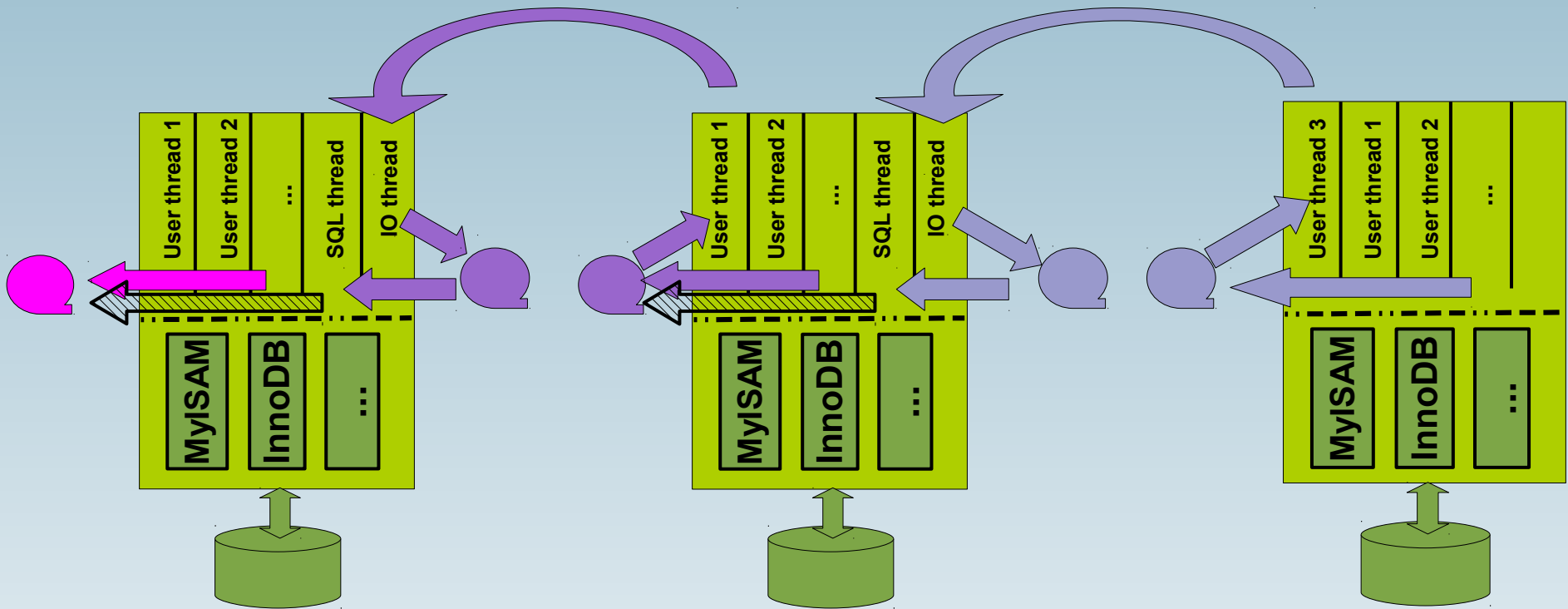
**Master:**  
“log-bin = FILE”, or else no binlog  
(no master function)

# Typical Usage



- **”High Availability“**
- **Geographic redundancy**
- **Support higher read load  
(= ”read scale-out“)**
- **Read-only instance(s)  
e.g. for backup or reports**
- **Intentional delay is possible**
- **Filtering (by DB or table) is possible**

# Replication Cascade



- Recommended: "read-only = 1" on slave  
"log\_slave\_updates = 1"
- Multiple slaves per master are possible



# Entries in Binlog

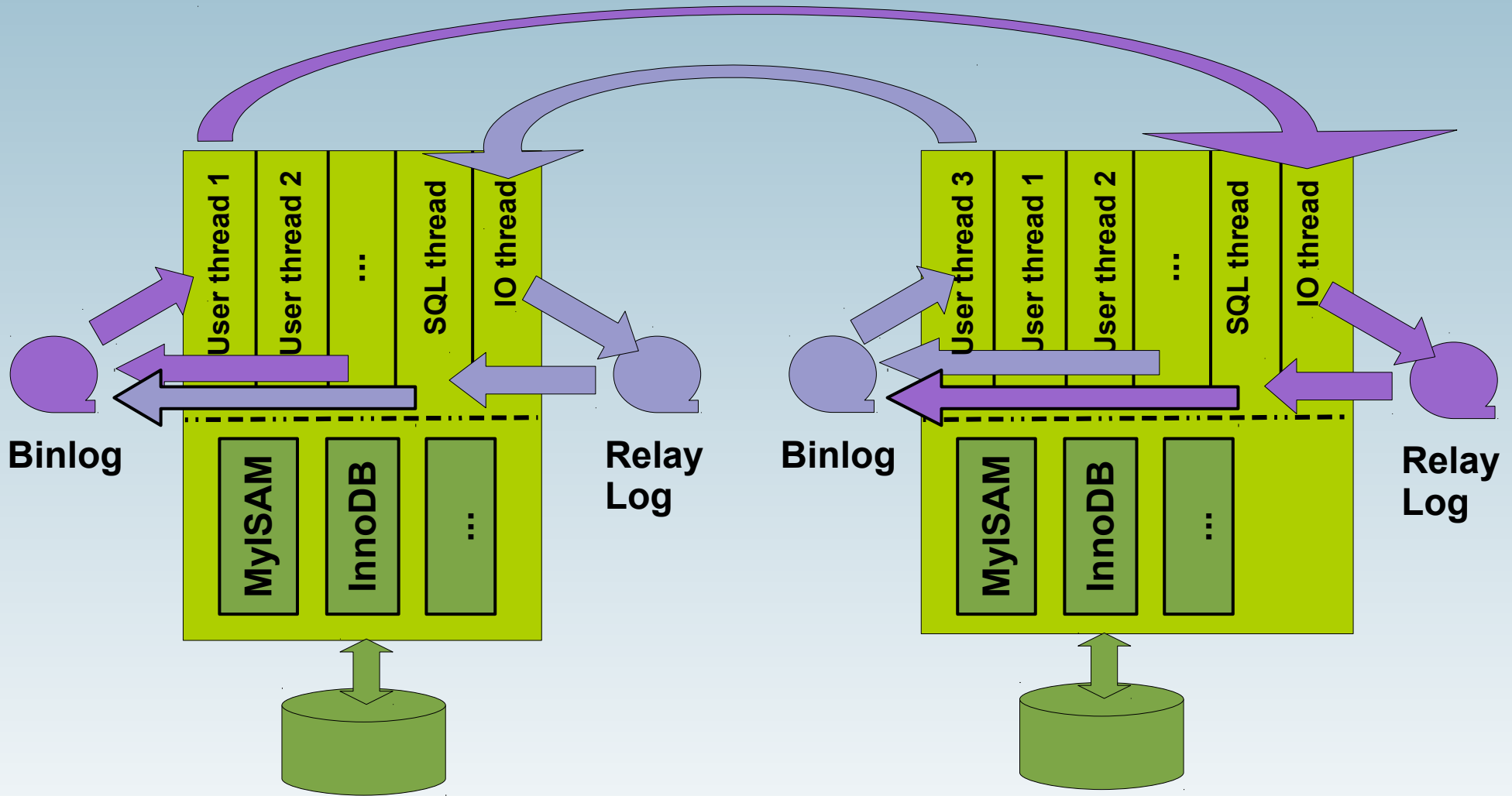
## Original:

- Identification by file name and position
- Replication: "change master to ..." specifying host, port, user, password, file, position
- See also: "mysqldump --master-data"

## From MySQL 5.6 also:

- GTID = "Global Transaction ID"
- Replication: "change master to ..." specifying host, port, user, password, "auto\_position = 1"

# Master-Master-Replication



- **Overlapping changes are fatal!**

- **Master-Master is controversial, be careful!**
- **Replication increases read throughput, but not/barely write throughput**
- **Replication causes file IO und network load**
- **Format "row" is more efficient, but less readable**
- **Multi-threaded replication since MySQL 5.6, multi-master ("multi-source") coming in MySQL 5.7**
- **Big installation: booking.com**
- **Recommended: [datacharmer.blogspot.de](http://datacharmer.blogspot.de) (Giuseppe Maxia, August 2015)**

MySQL Server: Architecture

Binlog

Replication

➔ **Galera Cluster**

Comparison

Examples / When (not) to Choose Which

# Replication Weaknesses



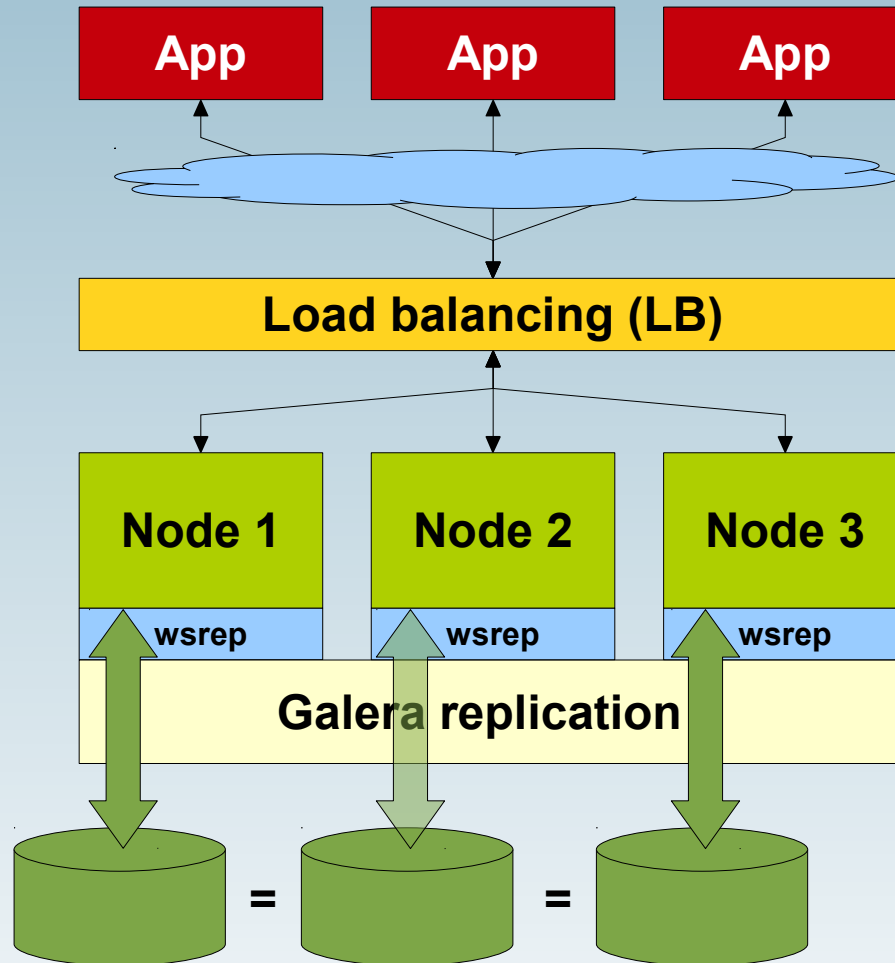
- **Asynchronous**
- **Asymmetrical**
- **Only one write node**
- **Parallel writes may cause breakage**
- **HA needs failover after node crash**
- **Each node is SPOF for its slaves,  
breakdown requires structure change**
- **Dynamic changes are complicated**

# Better Alternative



- **Synchronous transfer**
- **Symmetrical cluster**
- **Write accesses on all nodes**
- **Distributed conflict analysis and handling**
- **HA by continuity after node outage**
- **Dynamic entry / exit of nodes supported**

# Galera Cluster



Including outage detection and redirection for HA

“Working Set Replication”

Dedicated network preferred

Locale disks,  
each holding all data

“shared nothing” architecture

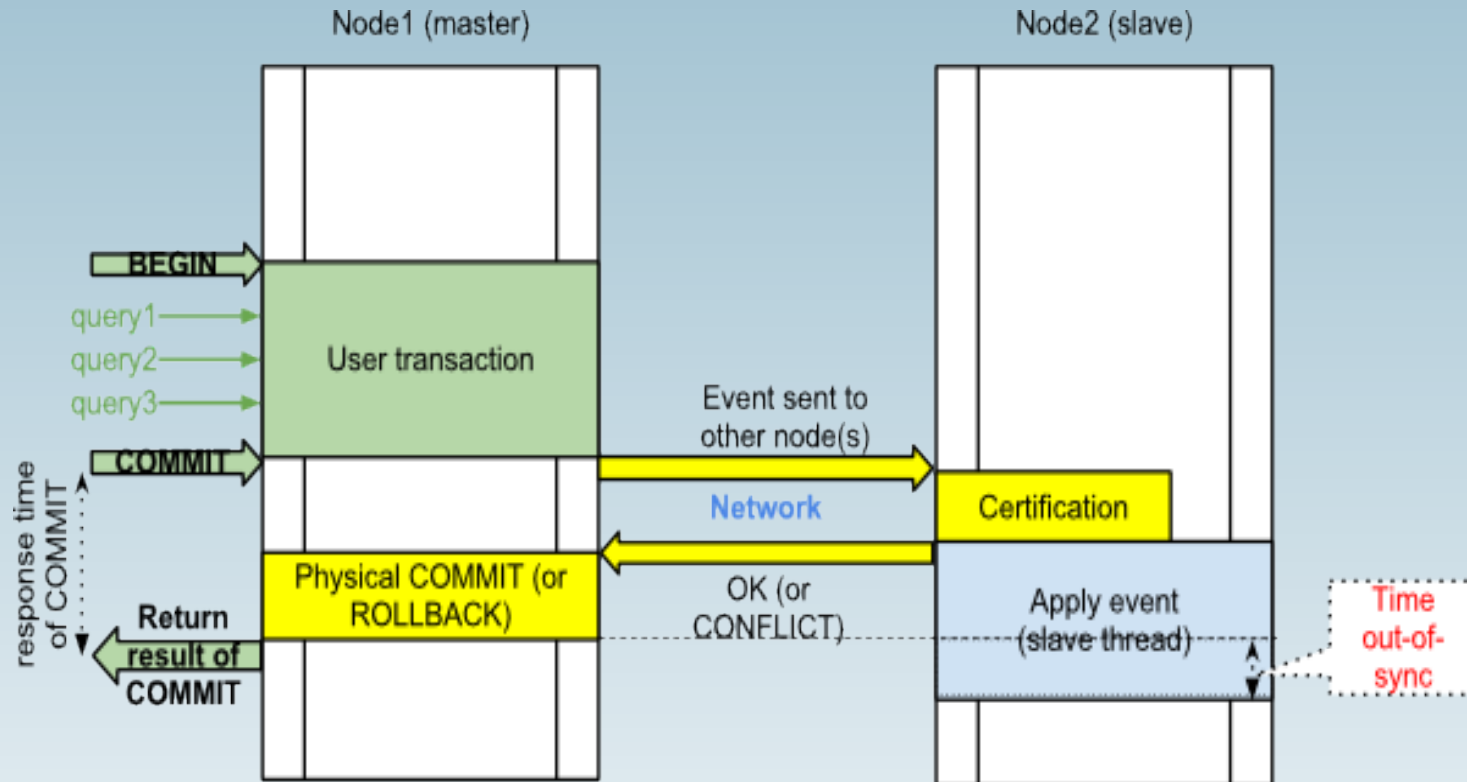
# Galera Properties (1)



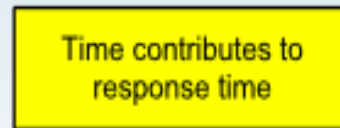
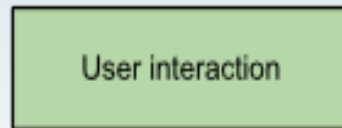
- + Based on InnoDB (due to transactions and rollback)**
- + Also transfers user definitions, privileges, ...**
- + Quasi-synchronous transfer on commit, check for conflicts, efficient**
- + Symmetrical, HA without server failover, quorum**
- + No loss of transactions**
- + Brings read scale-out, also some write increase**
- + Dynamical entry / exit possible, synchronisation is automated**



# Order of Events



Legend:



Graph by  
 Vadim Tkachenko  
 (Percona):

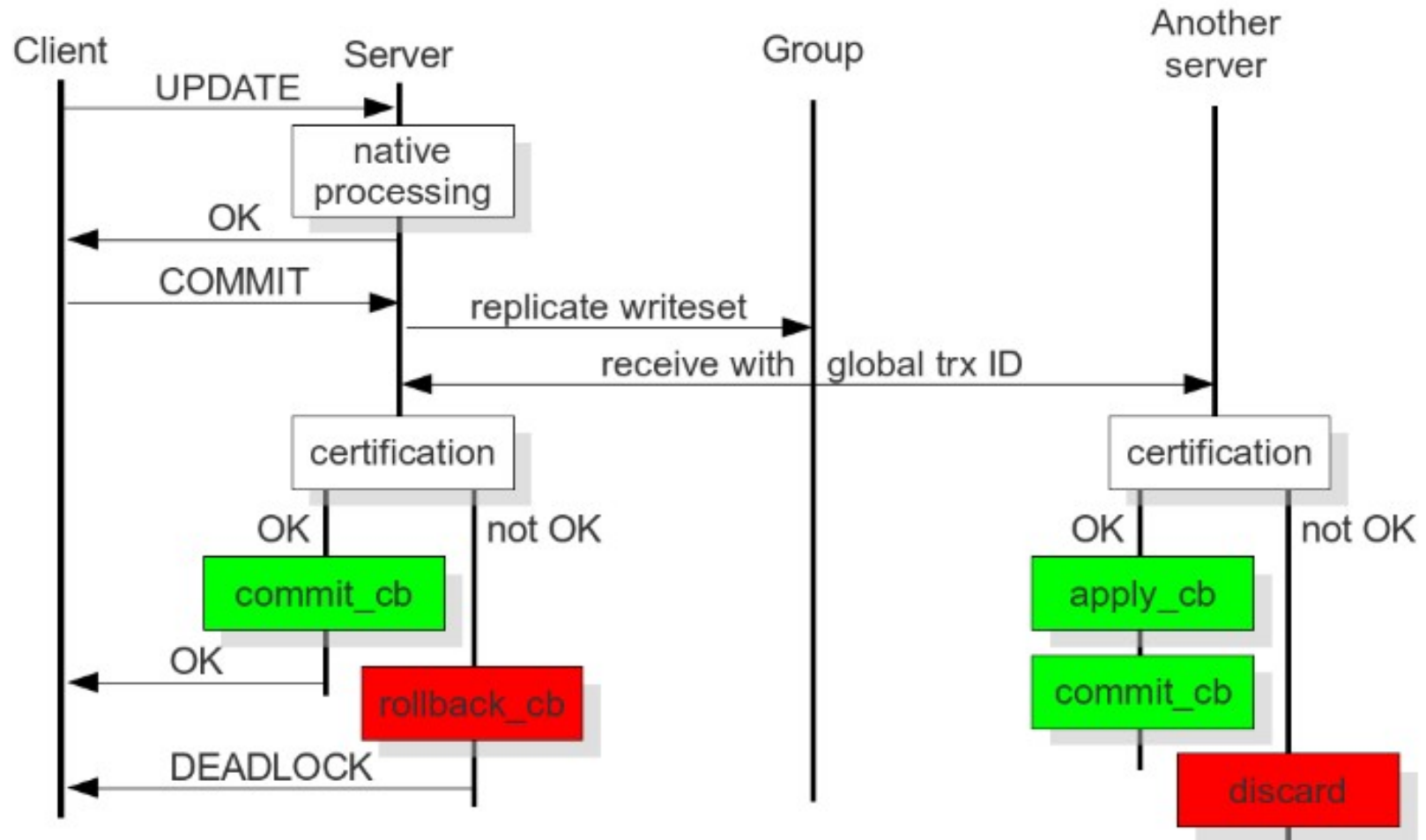
<http://www.mysqlperformanceblog.com/2012/01/19/percona-xtradb-cluster-feature-2-multi-master-replication/>

# Galera Properties (2)



- **MySQL sources need patching**  
(Codership offers binaries, ditto MariaDB and Percona)
- **Beware of hot spots (rows)**
- **Conflict detection is late, full rollback**  
(Check postponed till commit)
- **Minimum size is 3 nodes**
- **Synchronisation time for large DB**  
(mysqldump -> xtrabackup or rsync)

# Certification at Commit



<http://galeracluster.com/documentation-webpages/certificationbasedreplication.html>

# MySQL Server: Architecture

## Binlog

## Replication

## Galera Cluster

➔ **Comparison**

## Examples / When (not) to Choose Which

- **Alternatives: Replication or Galera Cluster**
- **Redundancy of machine and storage**
- **HA**
- **Scale-out, esp. for read load**
- **Instances for reports, analysis, backup**
- **Data available locally (branch offices, ...)**

# Comparison (1)

| Replication                      | Galera                       |
|----------------------------------|------------------------------|
| <b>Standard</b>                  | <b>Add-on product</b>        |
| <b>All handlers</b>              | <b>InnoDB only</b>           |
| <b>Upwards compatible</b>        | <b>Same versions</b>         |
| <b>Minimum 2 nodes</b>           | <b>Minimum 3 nodes</b>       |
| <b>HA by failover</b>            | <b>HA without changes</b>    |
| <b><i>Communication:</i></b>     |                              |
| <b>Hierarchical, chain</b>       | <b>Symmetrical, parallel</b> |
| <b>Asynchronous</b>              | <b>Quasi-synchronous</b>     |
| <b>Delay is configurable</b>     | <b>Immediate</b>             |
| <b>Filtering is configurable</b> | <b>Complete</b>              |

# Comparison (2)

| Replication                         | Galera                    |
|-------------------------------------|---------------------------|
| <b>Read scale-Out</b>               | <b>Read scale-Out</b>     |
| <b>Write unchanged</b>              | <b>Write increased</b>    |
| <b><i>1 Master:</i></b>             |                           |
| <b>1* write</b>                     | <b>1* write</b>           |
| <b><i>Local conflict :</i></b>      |                           |
| <b>Error on statement</b>           | <b>Error on statement</b> |
| <b><i>n Master:</i></b>             |                           |
| <b>n* write</b>                     | <b>n* write</b>           |
| <b><i>Distributed conflict:</i></b> |                           |
| <b>Replication breakdown</b>        | <b>Rollback on commit</b> |

# Comparison (3)

| Replication                           | Galera                                 |
|---------------------------------------|--|
| <b><i>Short interruption:</i></b>     |  |
| Replication resumes                   | IST (incremental transfer)             |
| <b><i>Long interruption:</i></b>      |  |
| Replication resumes                   | SST (full transfer)                    |
| <b><i>Structure change:</i></b>       |  |
| Manual / separate Tool                | Automatic / dynamic                    |
| <b><i>Initial setup:</i></b>          |  |
| Snapshot,<br>master remains available | Full transfer,<br>donor may be blocked |



# CAP Theorem



**”For a distributed computer system, it is impossible to simultaneously provide all three of the following guarantees:**

- **C = Consistency** (identical data throughout)
- **A = Availability** (system is operational)
- **P = Partition Tolerance** (network outage)“

**Eric Brewer, 1998 ff.**

[https://en.wikipedia.org/wiki/CAP\\_theorem](https://en.wikipedia.org/wiki/CAP_theorem)

# Prospect: MySQL 5.7



- **MySQL 5.7 is GA (5.7.9, 2015-Oct-21)**
- **Replication like in MySQL 5.6, added: multi-source replication (one slave reading from several masters)**
- **Codership is working on adding Galera Cluster to MySQL 5.7**
- **Oracle is working on "Group replication", currently available as "labs release" (= "not fit for production")**

# MySQL Server: Architecture

## Binlog

## Replication

## Galera Cluster

## Comparison

➔ **Examples / When (not) to Choose Which**

## Galera Cluster:

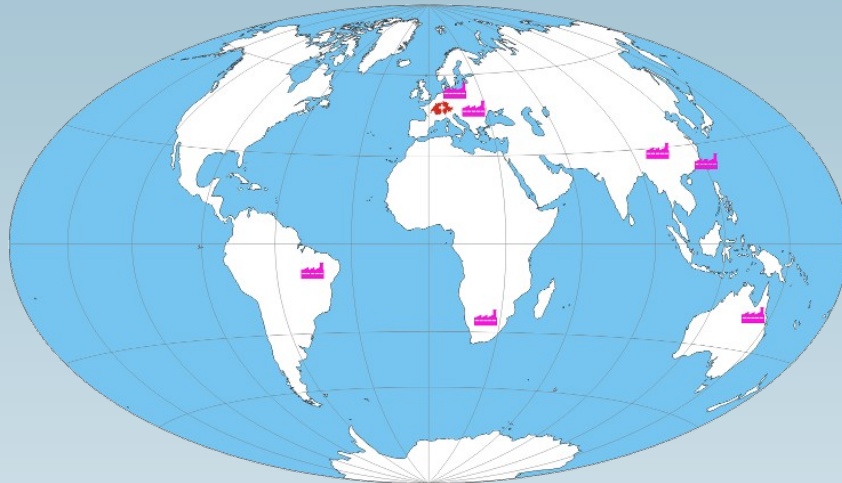
- **Isolated node has no quorum  
=> will not serve applications**
- **Quorum is at risk!**
- **Active nodes write "gcache" to files,  
storage period?**
- **Switch to SST threatens**

## Replication:

- **Master writes log segments to files**
- **IO-Thread asks to read from binlog position / GTID, retries periodical until successful**
- **Avoid "purge log"!**

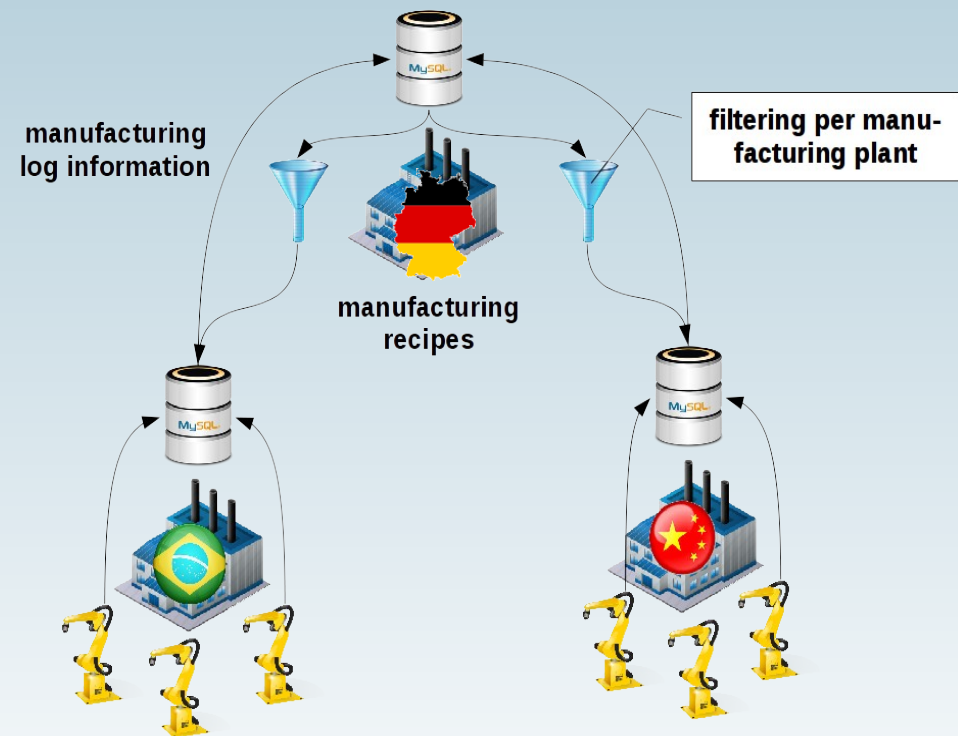
**Replication is more tolerant  
than Galera Cluster !**

# Global Production



## Solution: Replication with filtering

**Requirement:**  
Head office (D) and  
factories (BR, CN, ...)  
with selective transfer



# Parallel Writes + Conflict

## Galera:

- **Retry of autocommit statements configurable**
- **Transaction conflict causes rollback**  
**=> Application repeats complete transaction**

## Replication:

- **Slave detects conflict, no contact to application**  
**=> Replication stops**

**Replication needs admin action on conflict !**

# Hot Spot



- **Replication: frequent aborts**
  - **Galera: frequent rollbacks**
- => Agree on a single write node !**



# High Availability

## Replication:

- **Failover manual (reaction time) or automated (correct?)**
- **Slave lag, selection of new master**

## Galera:

- **Symmetrical, no change of roles**
- **Virtually synchronous replication (no lag)**

**=> Advantage Galera**



## Questions ?

## Discussion?

- **FromDual provides neutral and independent:**
  - **Consulting**
  - **Remote-DBA**
  - **Support for MySQL, Galera, Percona Server and MariaDB**
  - **Training**

[www.fromdual.com/presentations](http://www.fromdual.com/presentations)