

Semi-sync Group Ack

# Problem

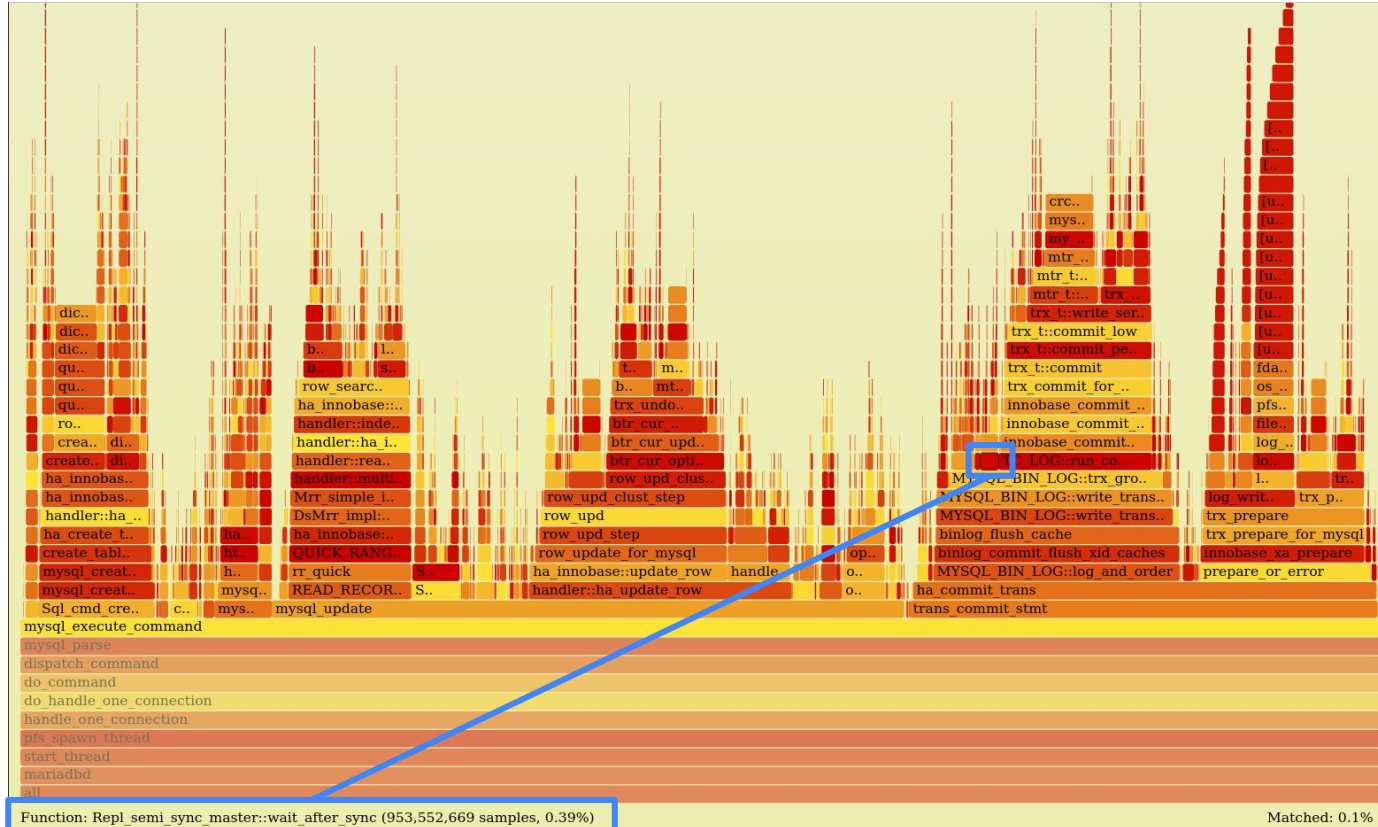
All transactions wait on one cond var

With `rpl_semi_sync_master_wait_point = AFTER_COMMIT`

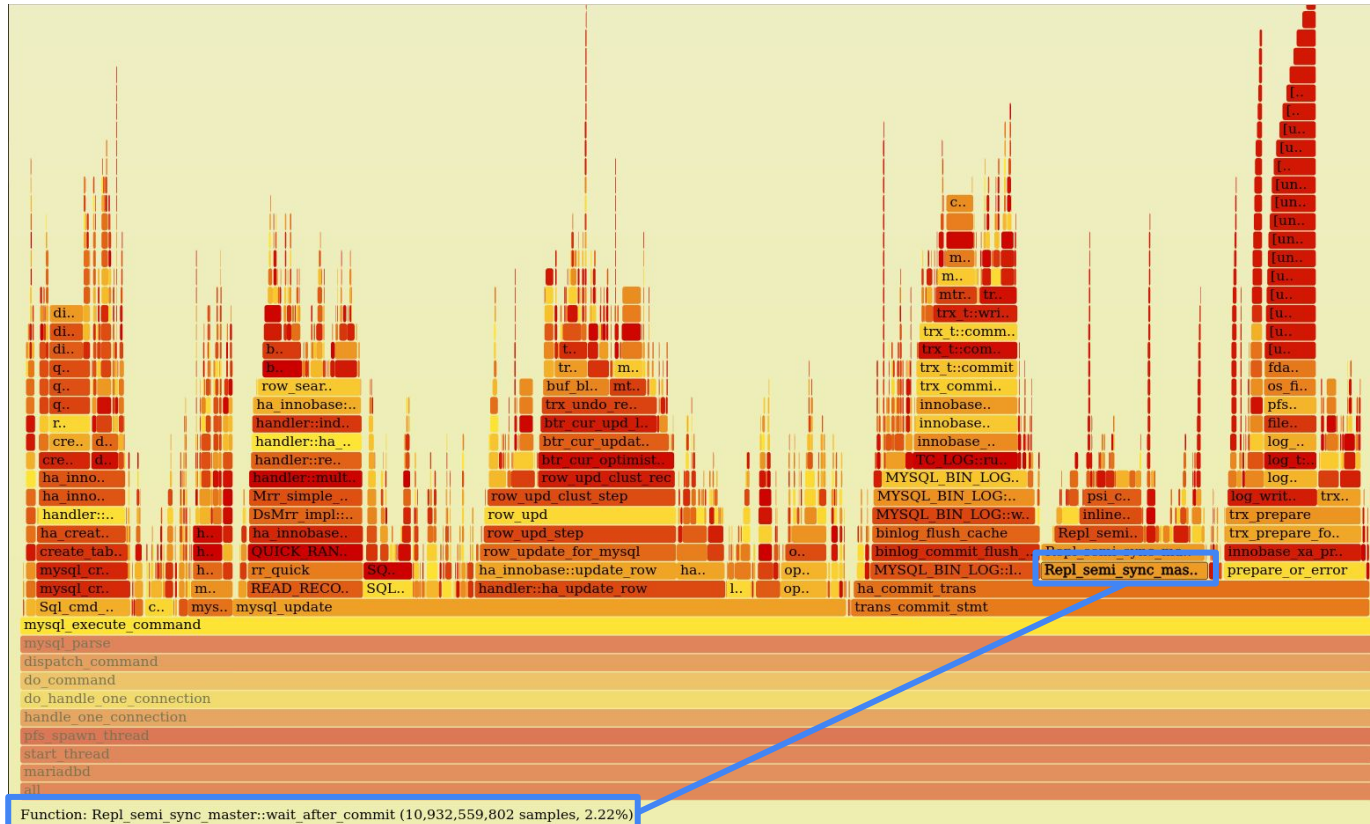
All concurrent transactions are signalled

# Flamegraphs

# repl\_semi\_sync\_master\_wait\_point = AFTER\_SYNC

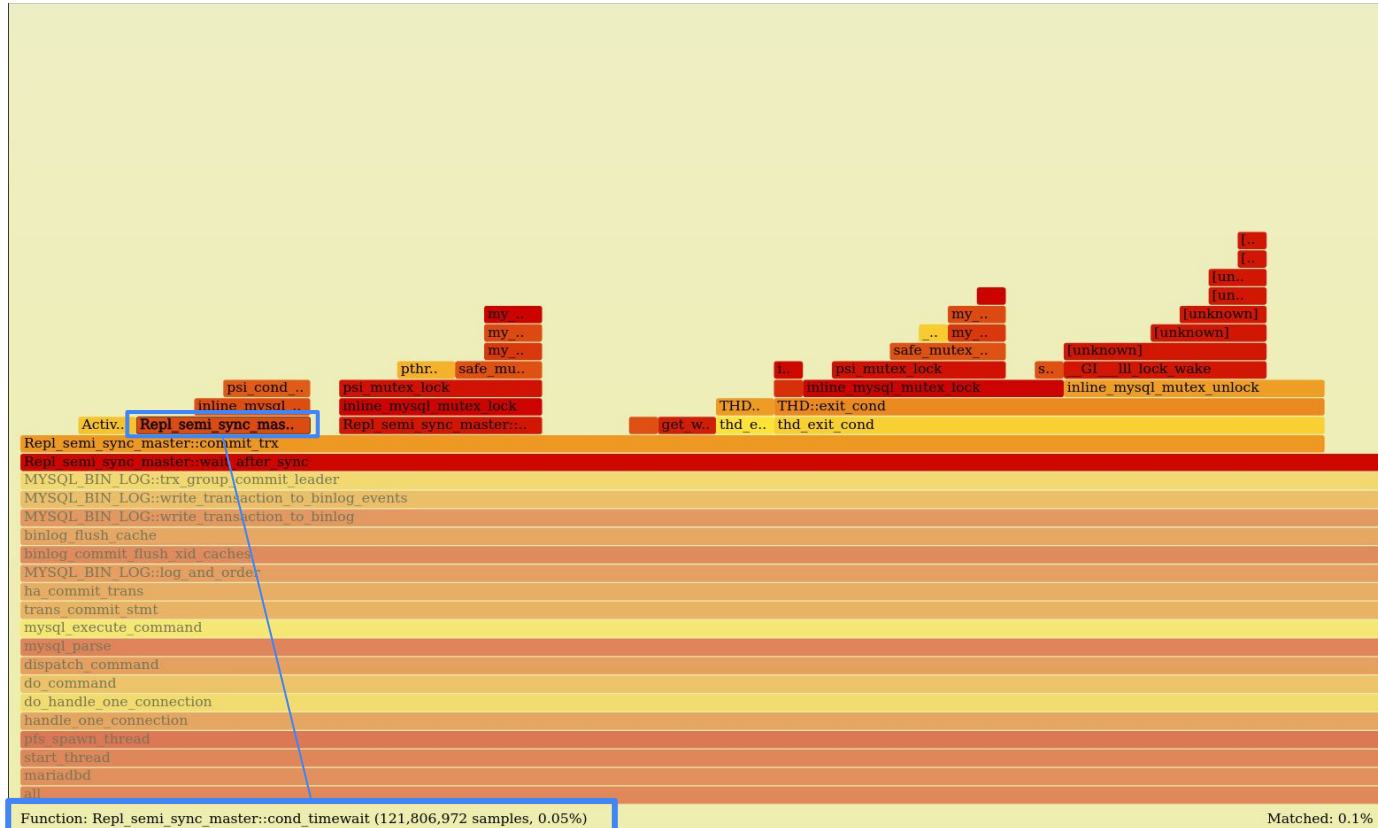


# rpl\_semi\_sync\_master\_wait\_point = AFTER\_COMMIT



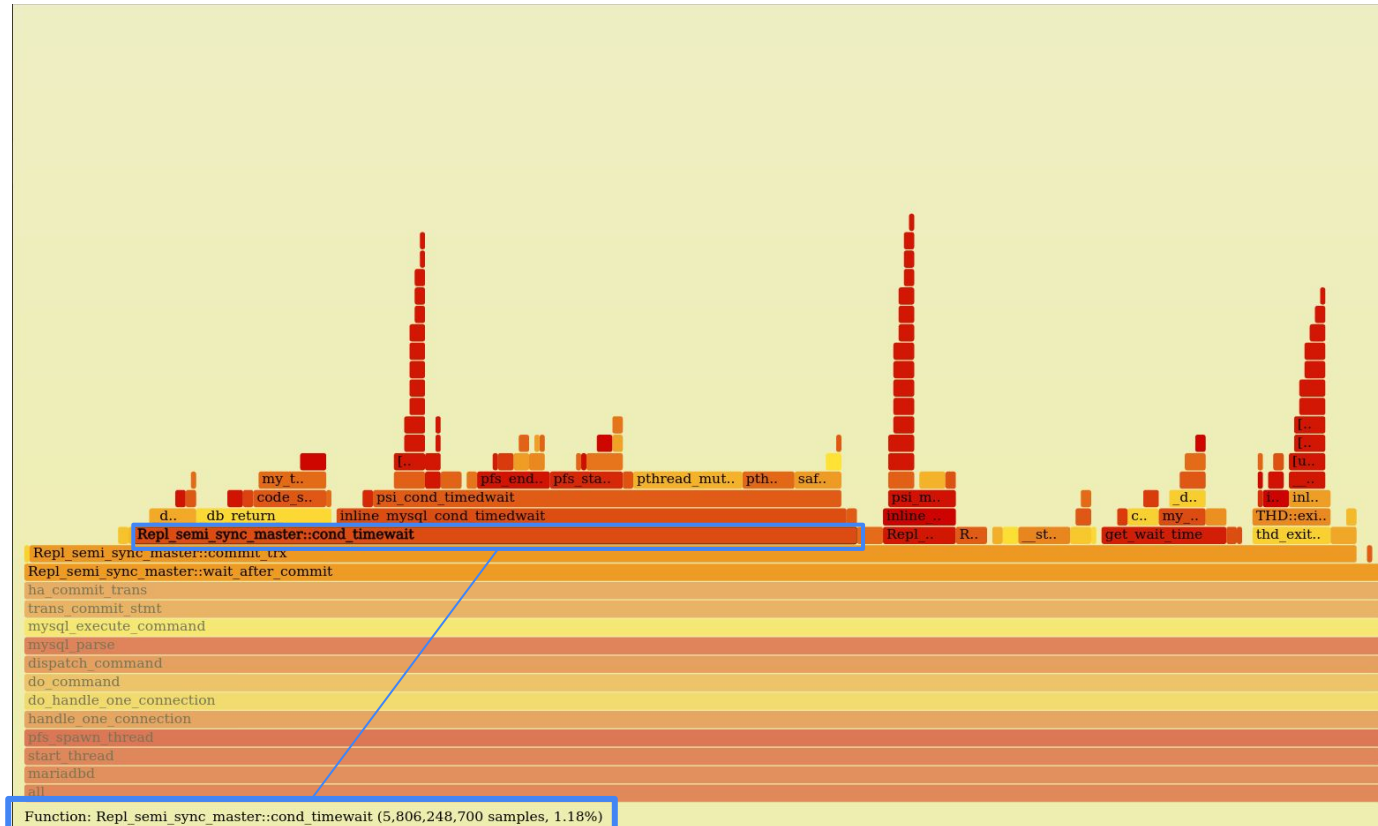
Zooming into `commit_trx()`

# AFTER\_SYNC time in cond\_timewait



AFTER\_SYNC  
100k samples

# AFTER\_COMMIT time in cond\_timewait



AFTER\_SYNC  
100k samples

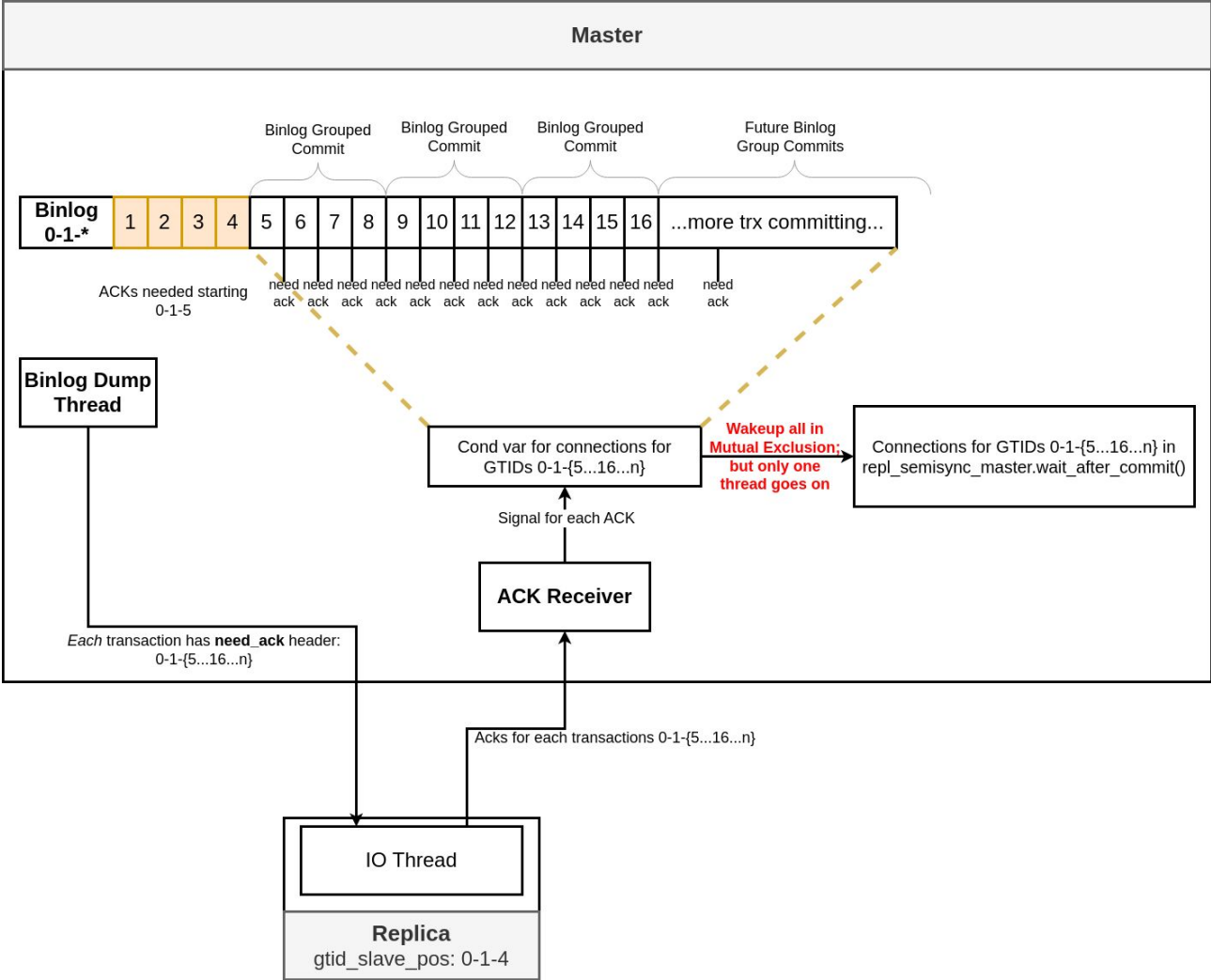
AFTER\_COMMIT  
5.8B samples



# AFTER\_COMMIT Complications

Multiple binlog group commits can happen ahead of current waiting trx

# Example



# Proposed Solution

**At binlogging time**, transactions register themselves with `repl_semisync_master` in binlog order (but don't yet wait)

..then **after** storage engine commit..

**repl\_semisync\_master** waits on just the oldest registered transactions (can be individual or groups)

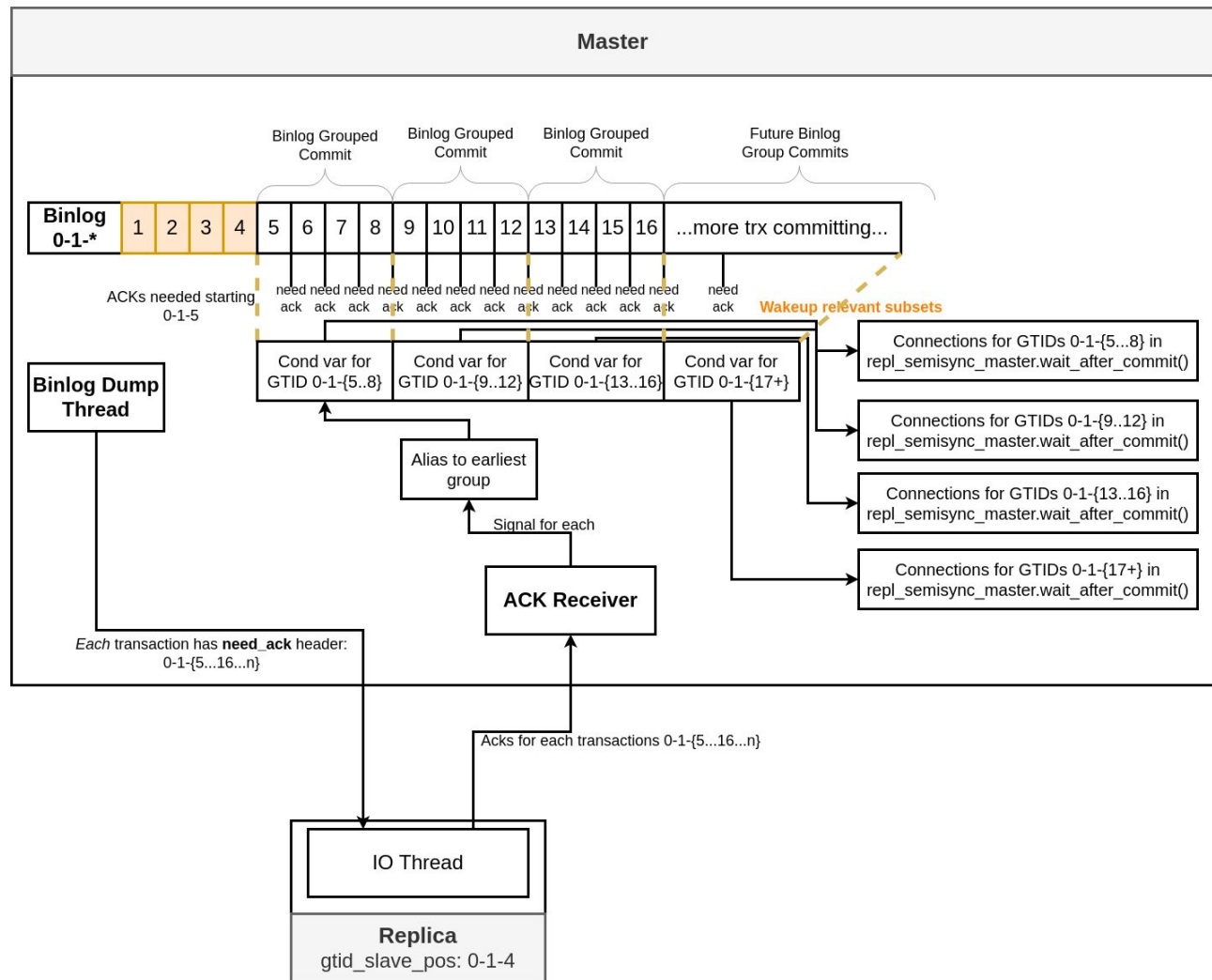
# Proposed Solution

Many options here:

1. 1 Cond Var per Registered transactions
  - a. No context switching problem, most memory and overhead (more acks from slave)
  - b. Can be GA
2. 1 Cond Var for a group of transactions
  - a. Limited context switching (only within group), slightly less memory footprint
  - b. Can align with binlog group commit boundary
  - c. Can be GA
3. 1 Cond Var for a group of transactions + group ack
  - a. No context switching, least overhead
  - b. New feature I imagine

Option 2 Would Look Something Like...

# Option 2



Option 1 would just be one cond\_var per  
binlog trx

# Conclusion After Discussion...

1. In 10.6, do option 1: One condition variable per connection thread
  - a. Can reuse THD::COND\_wakeup\_ready
  - b. MDEV-33551
2. in 11.X (new feature), implement group ack
  - a. MDEV-33491

[Link](#) to Zulip discussion