

RocksDB Brownbag: Write Paths

Siyong Dong

9/11/2020

What does RocksDB do in write path?

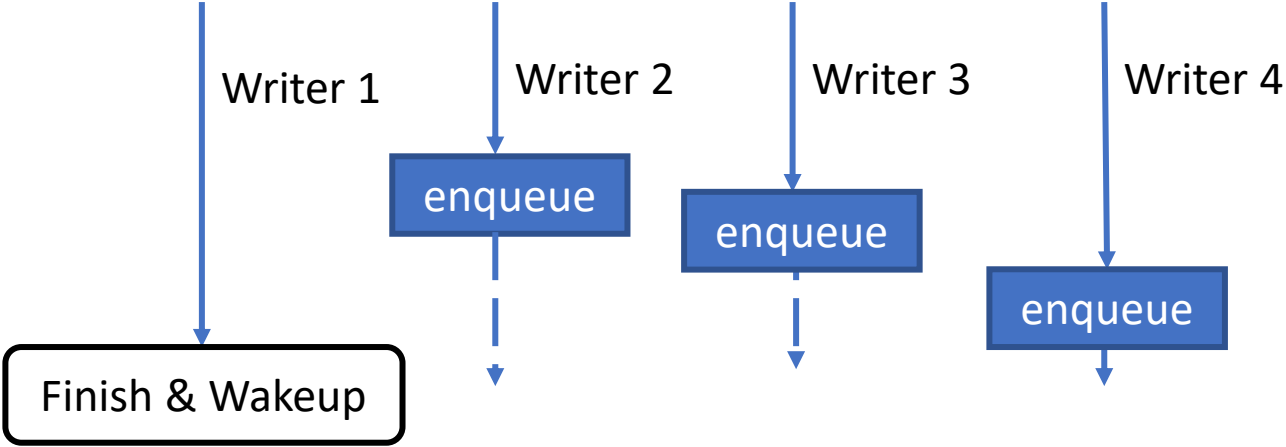
- Write to WAL
- Write to Memtable

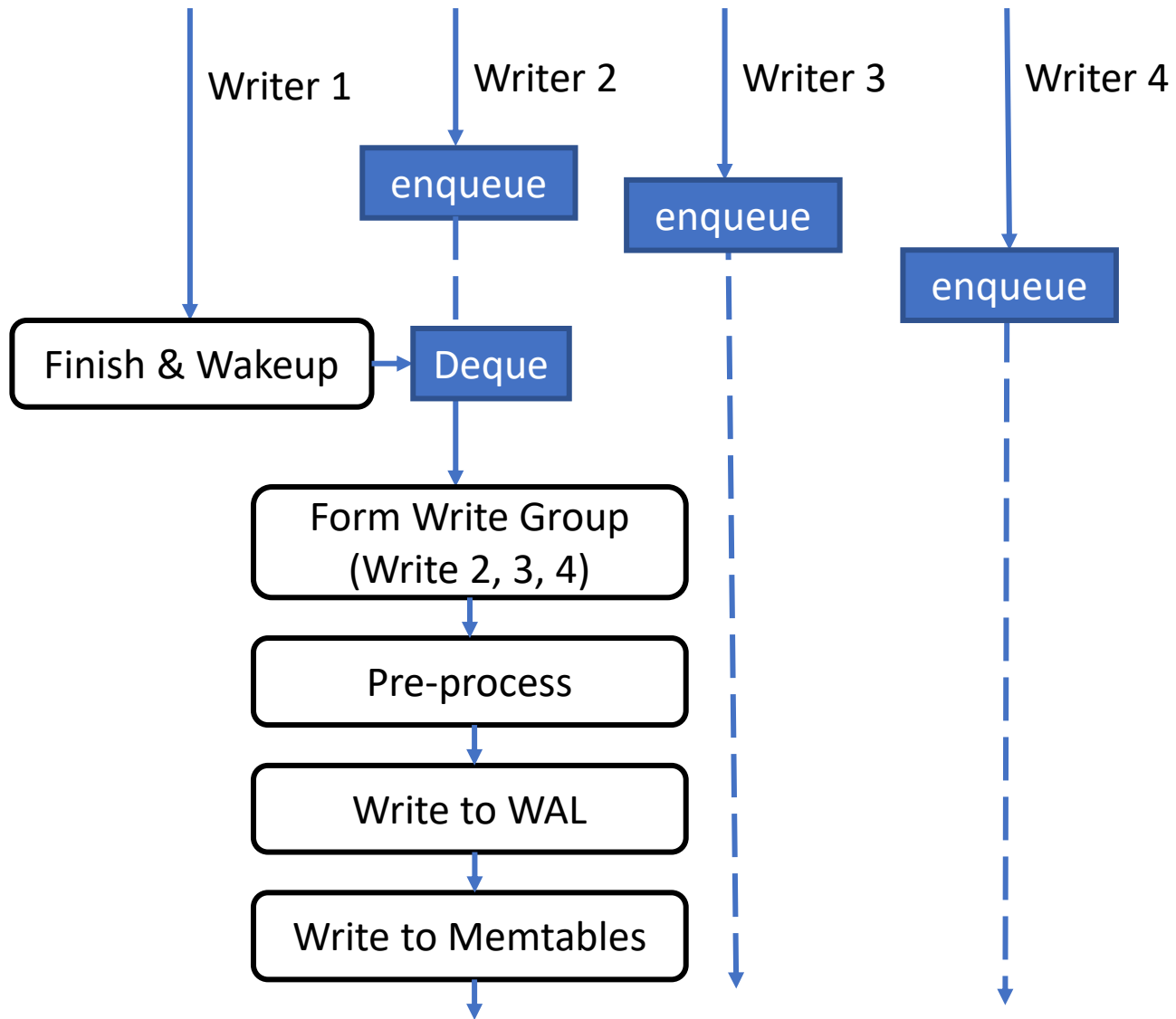
Why is write path complex?

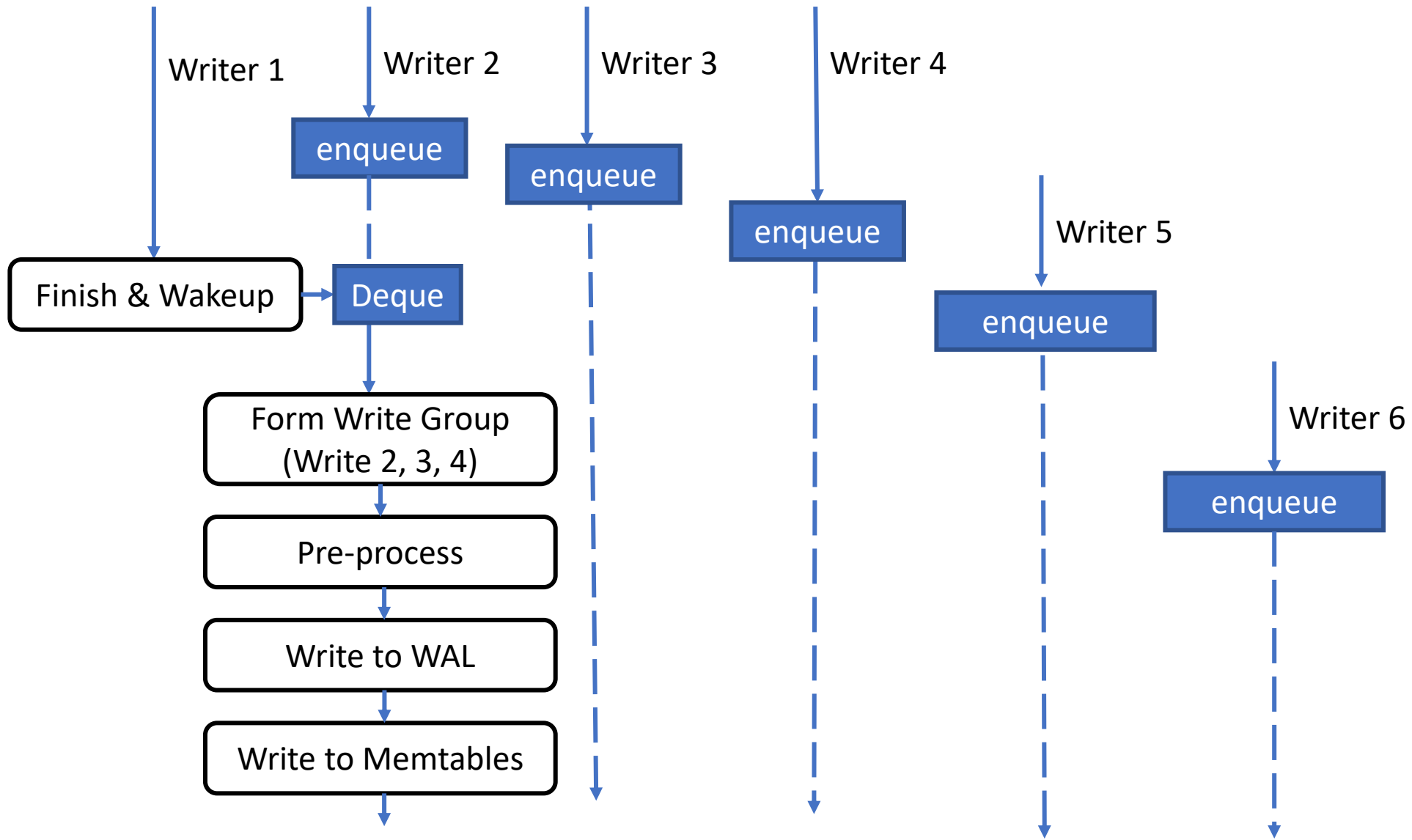
- Writing to WAL and memtable are not fully parallelizable
- Batching write (group commit)
- Pipelining
- Two-Phase-Commit (2PC)
- Unordered write

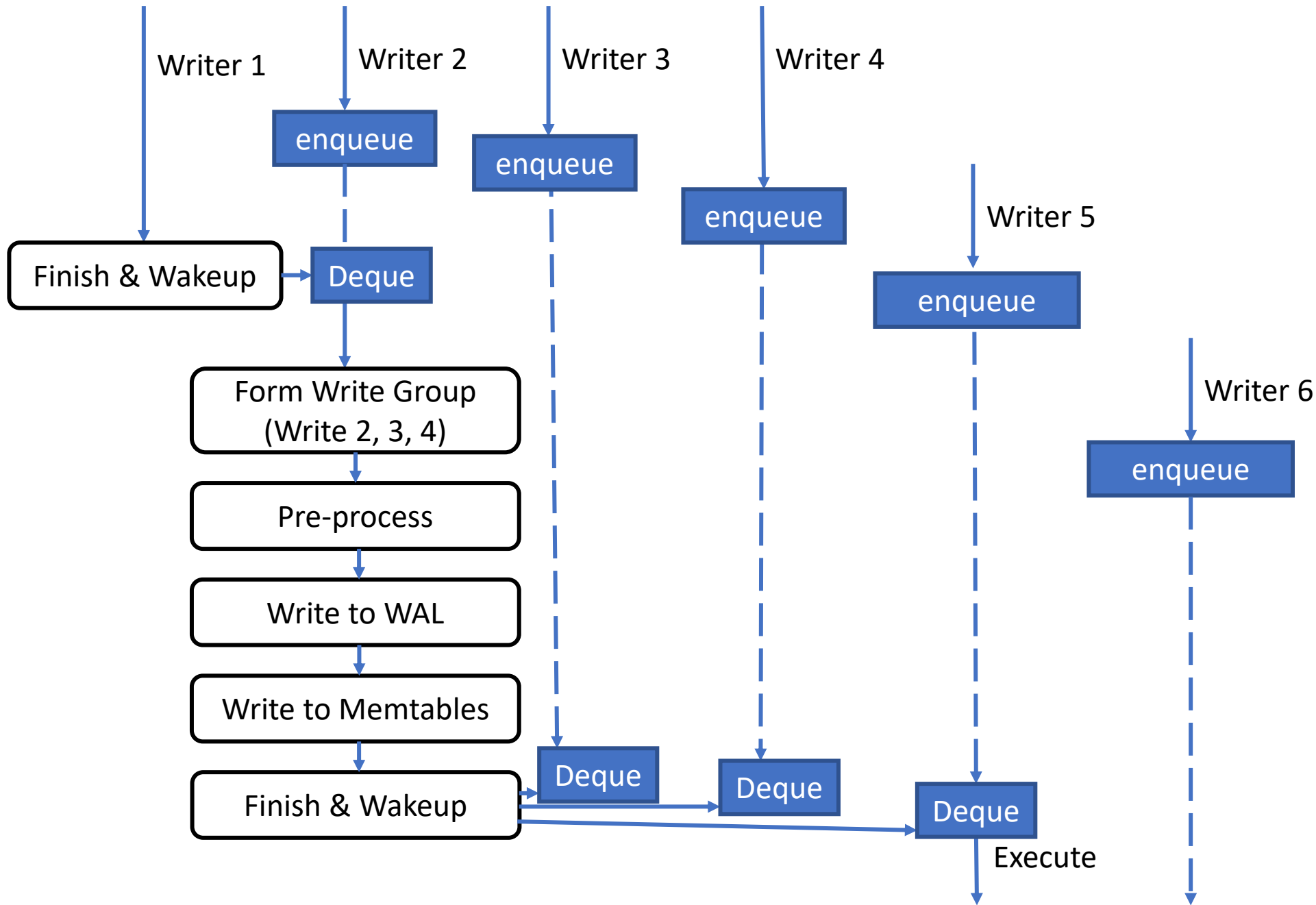
Where is it?

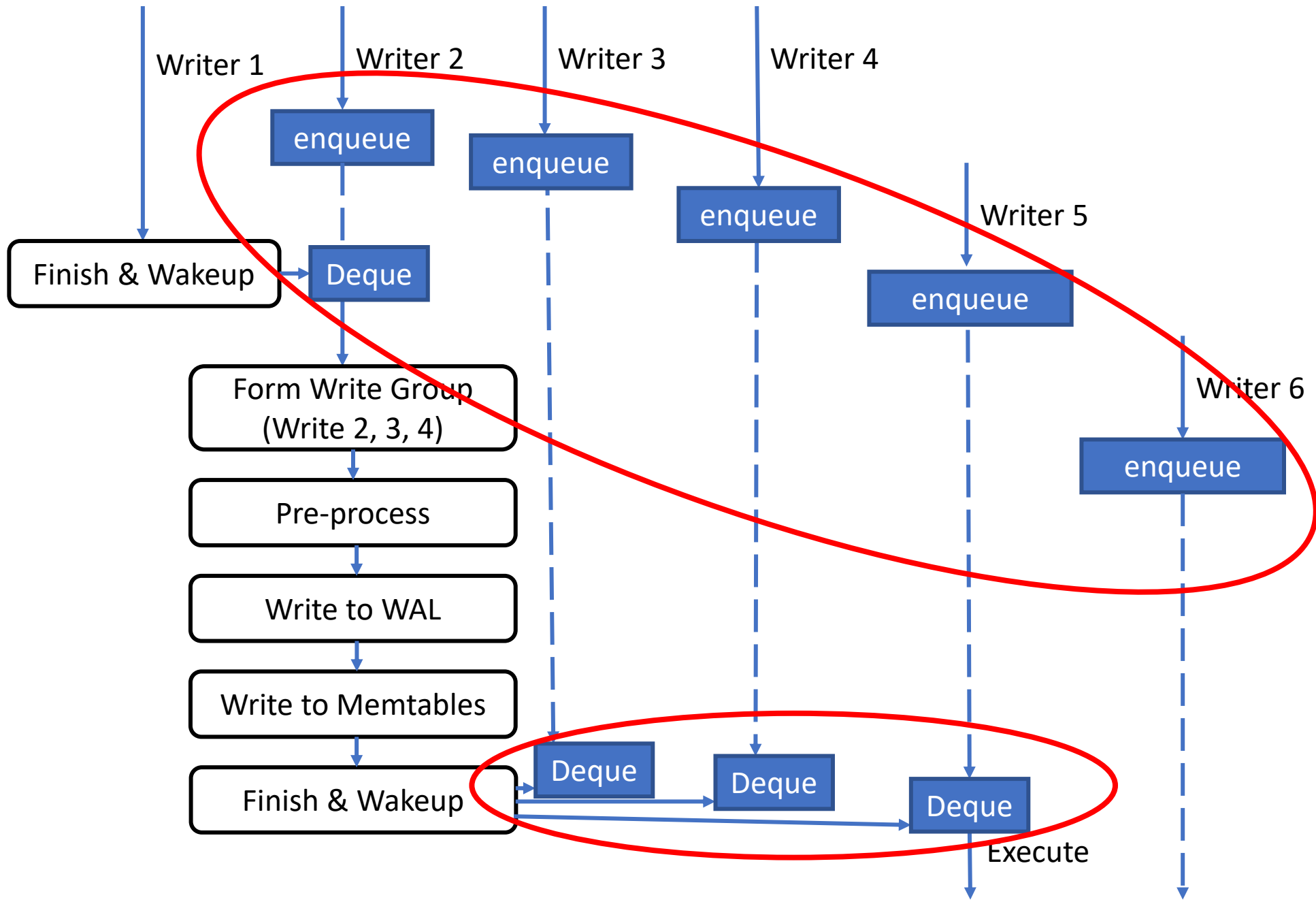
Function `DBImpl::WriteImpl()`



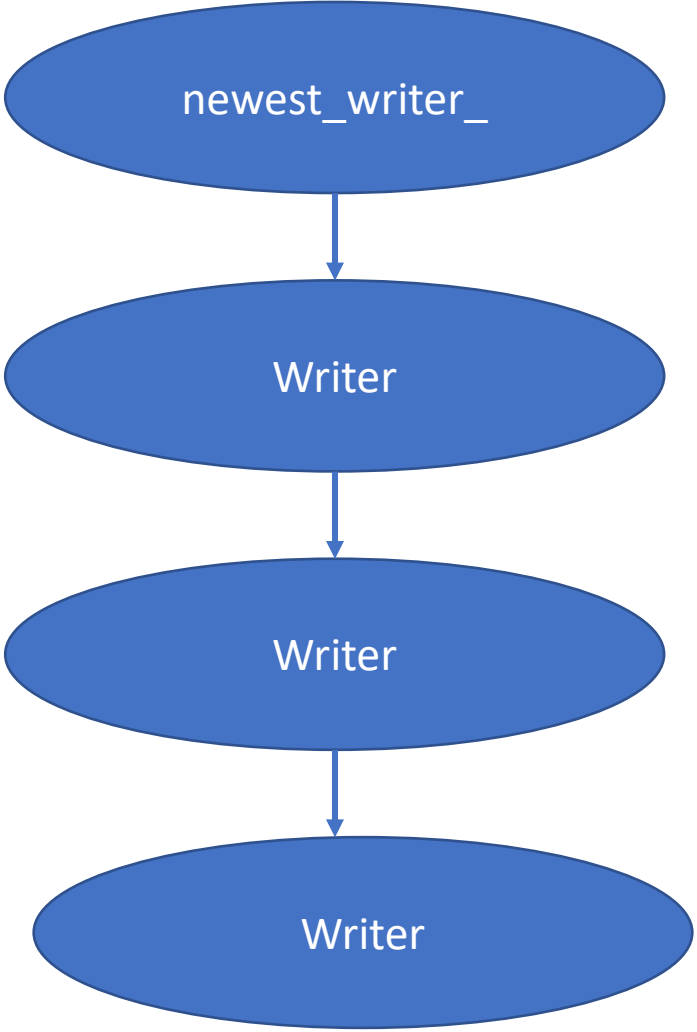




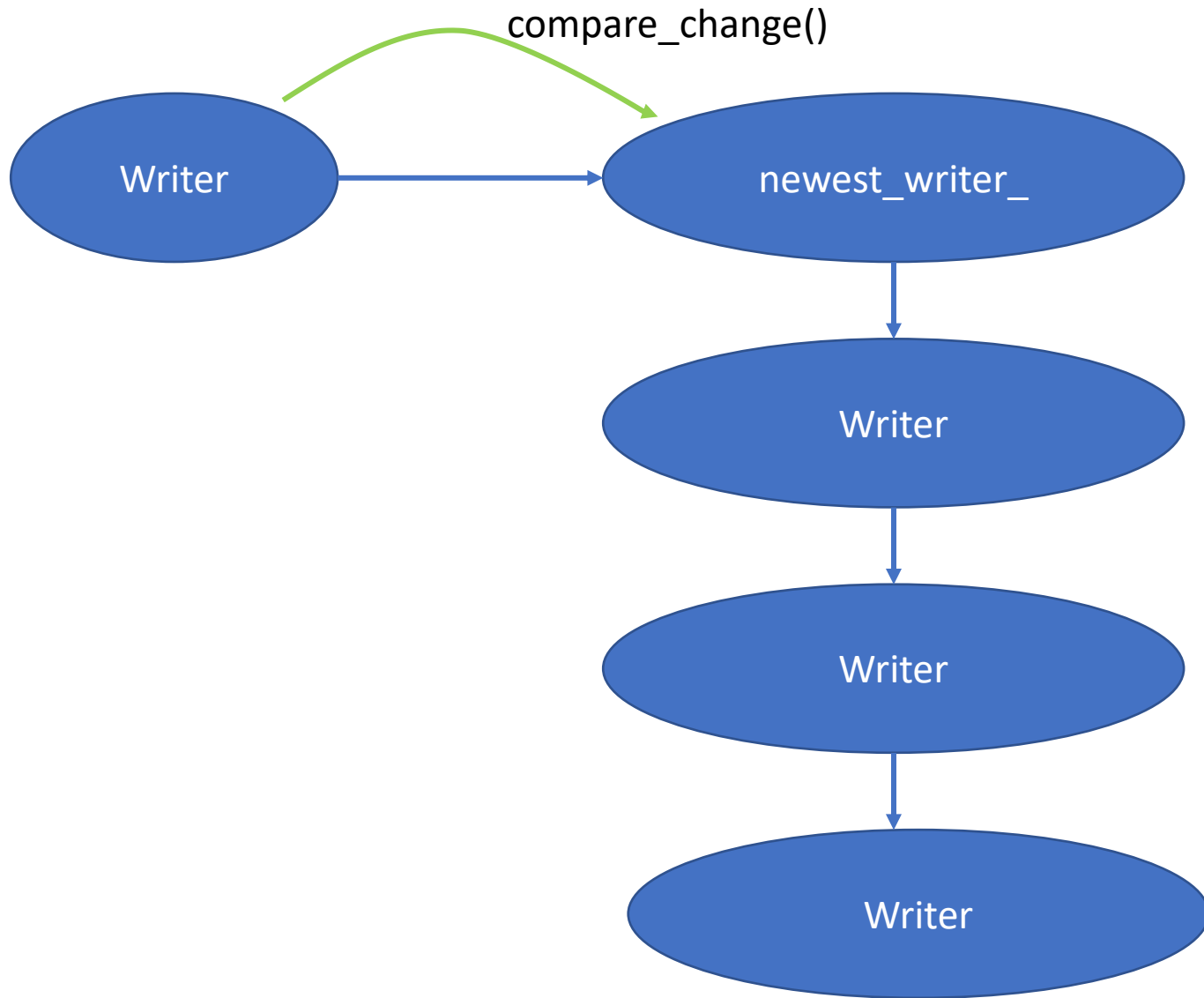




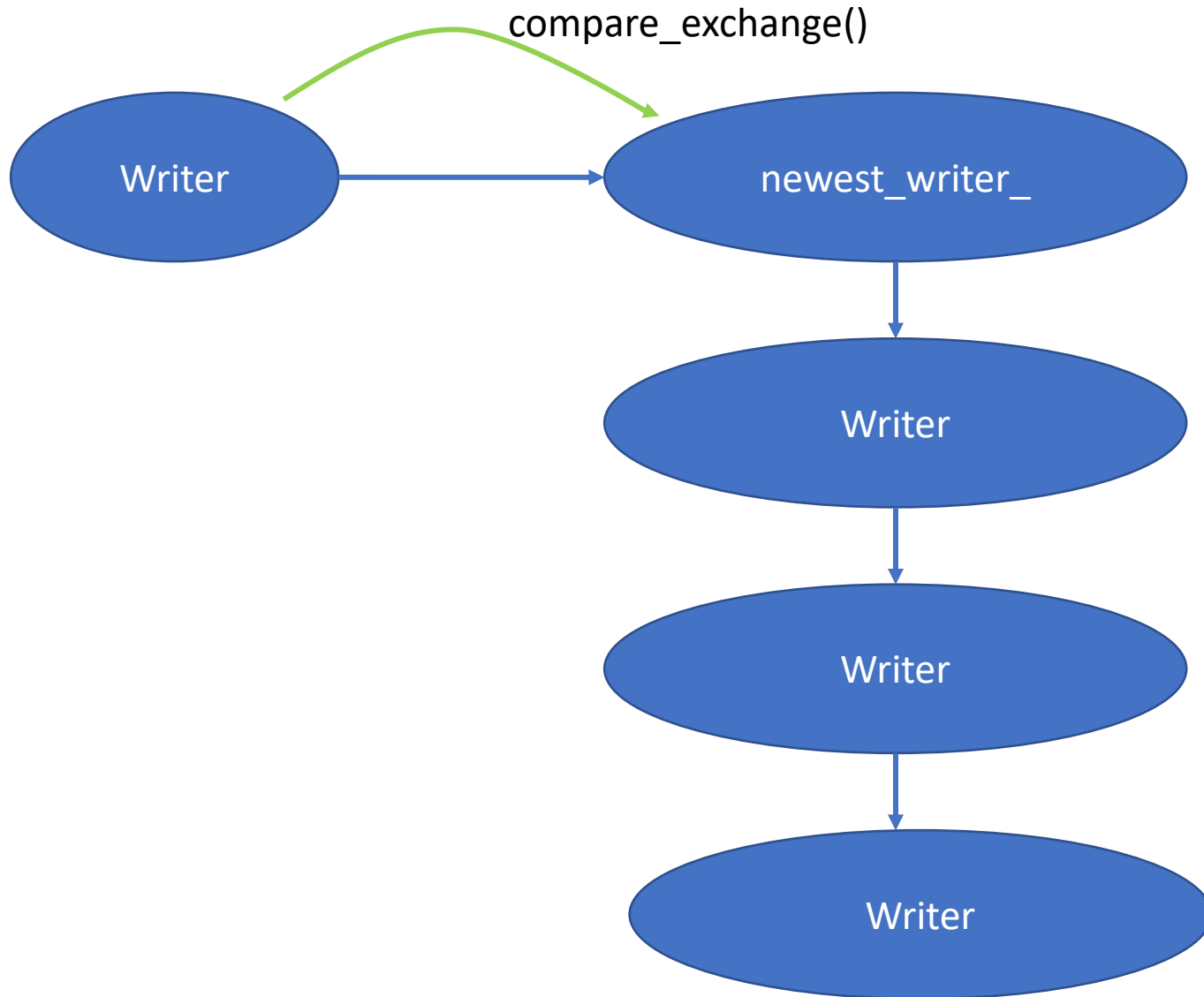
The write queue: *class* WriteThread



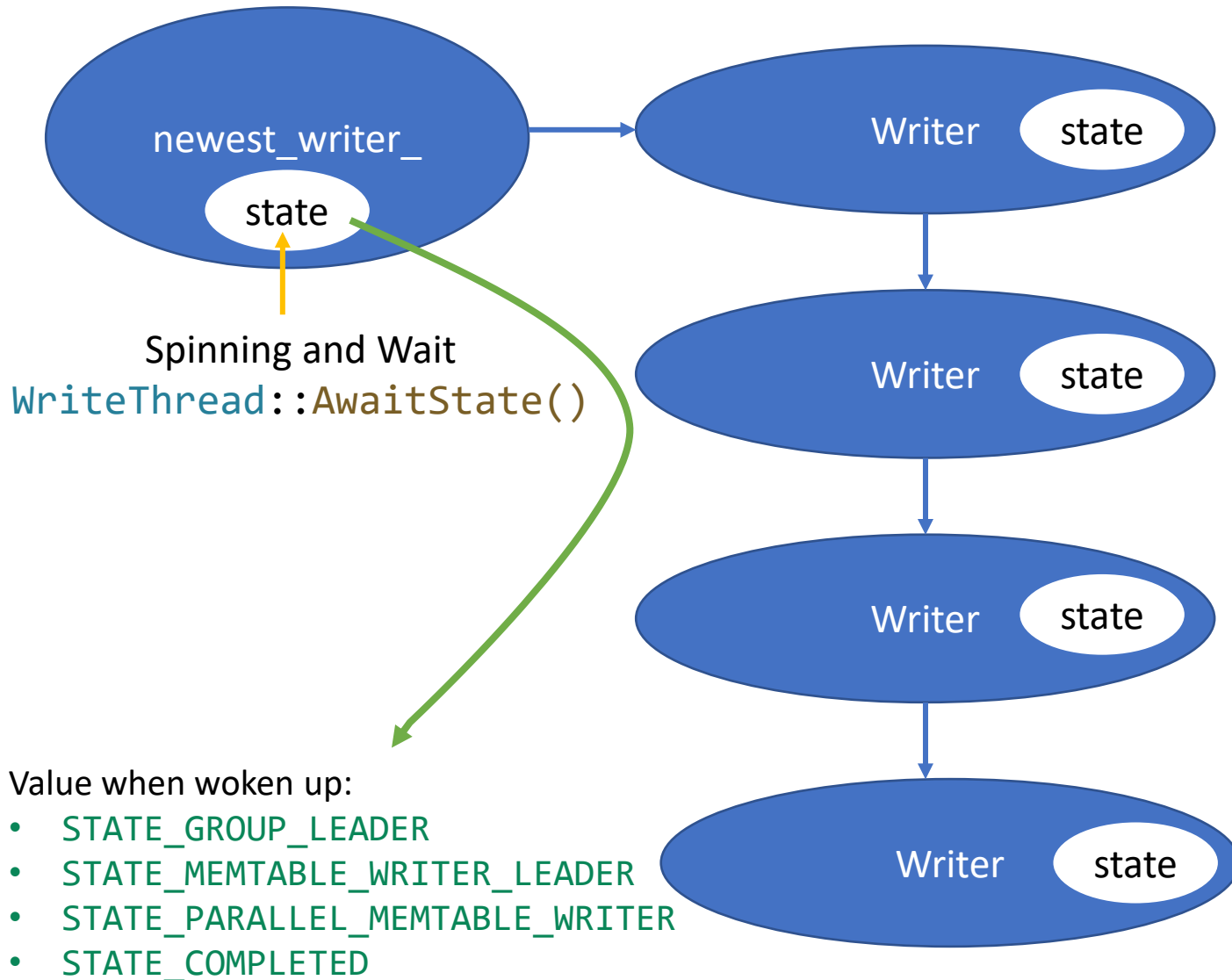
Enqueue: `WriteThread::JoinBatchGroup()`



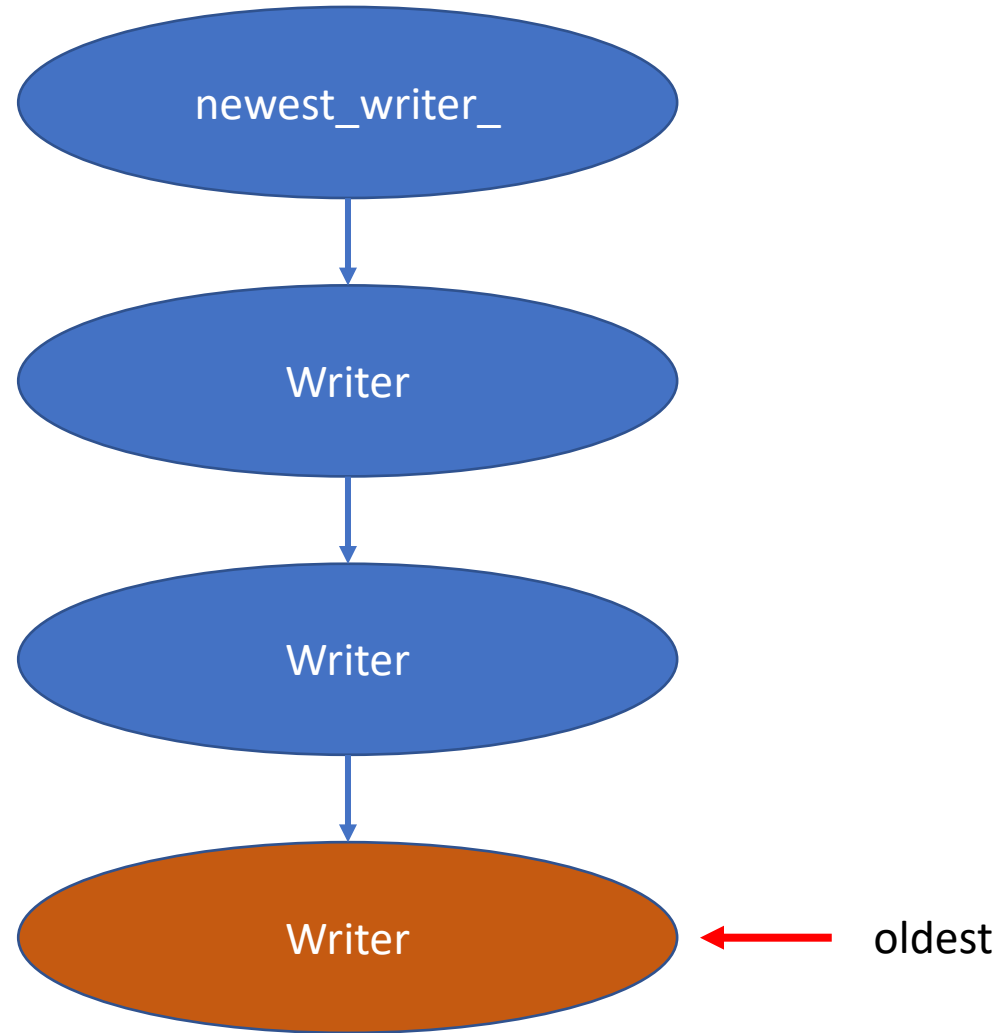
Enqueue: `WriteThread::JoinBatchGroup()`



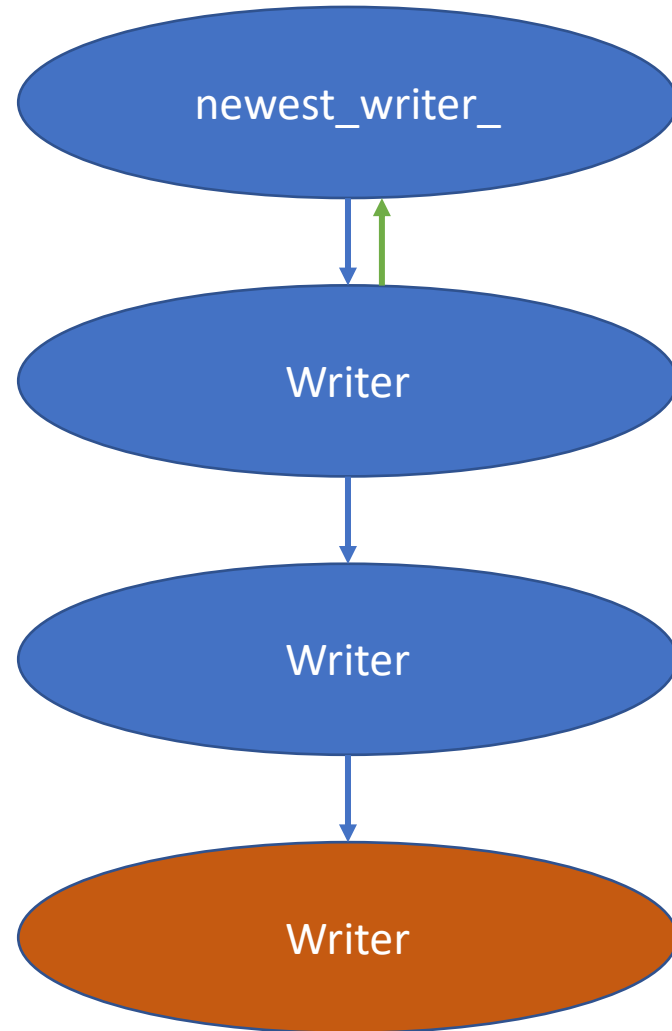
Enqueue: `WriteThread::JoinBatchGroup()`



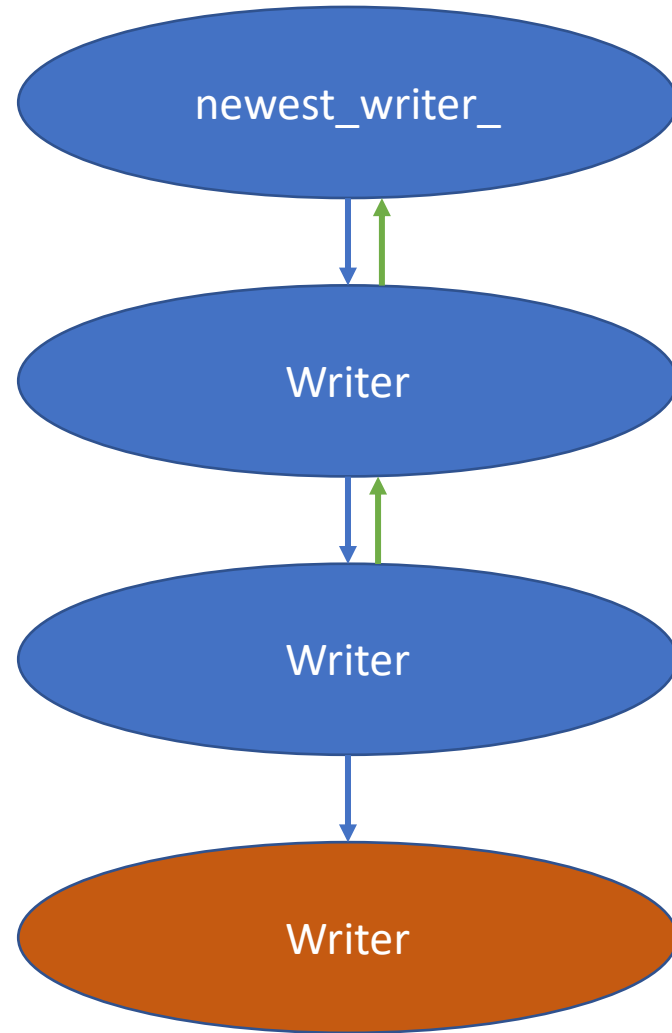
WriteThread::EnterAsBatchGroupLeader()



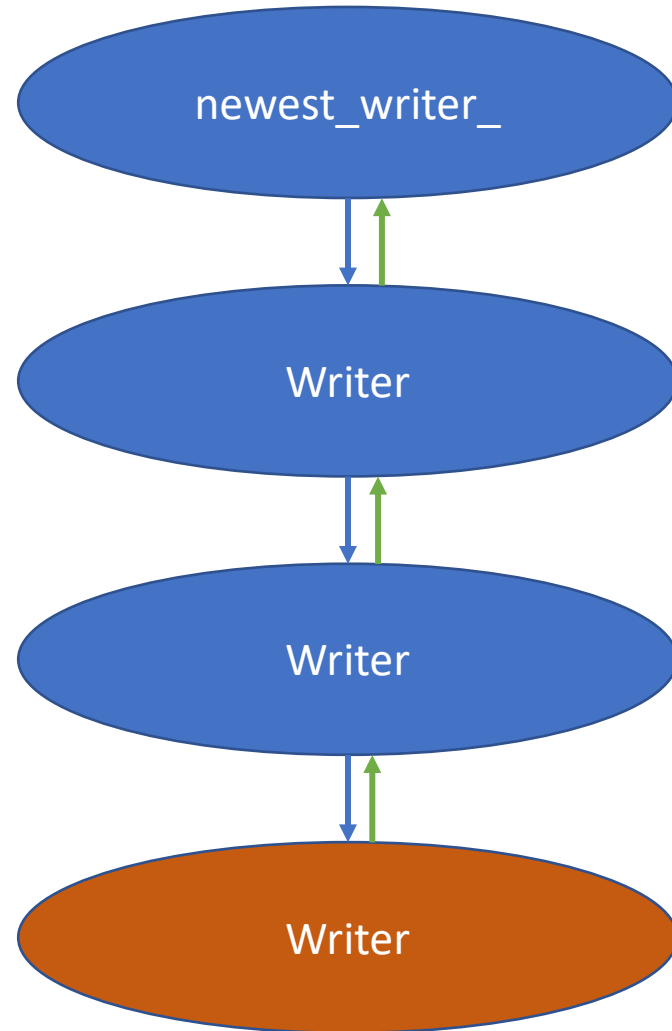
WriteThread::EnterAsBatchGroupLeader()



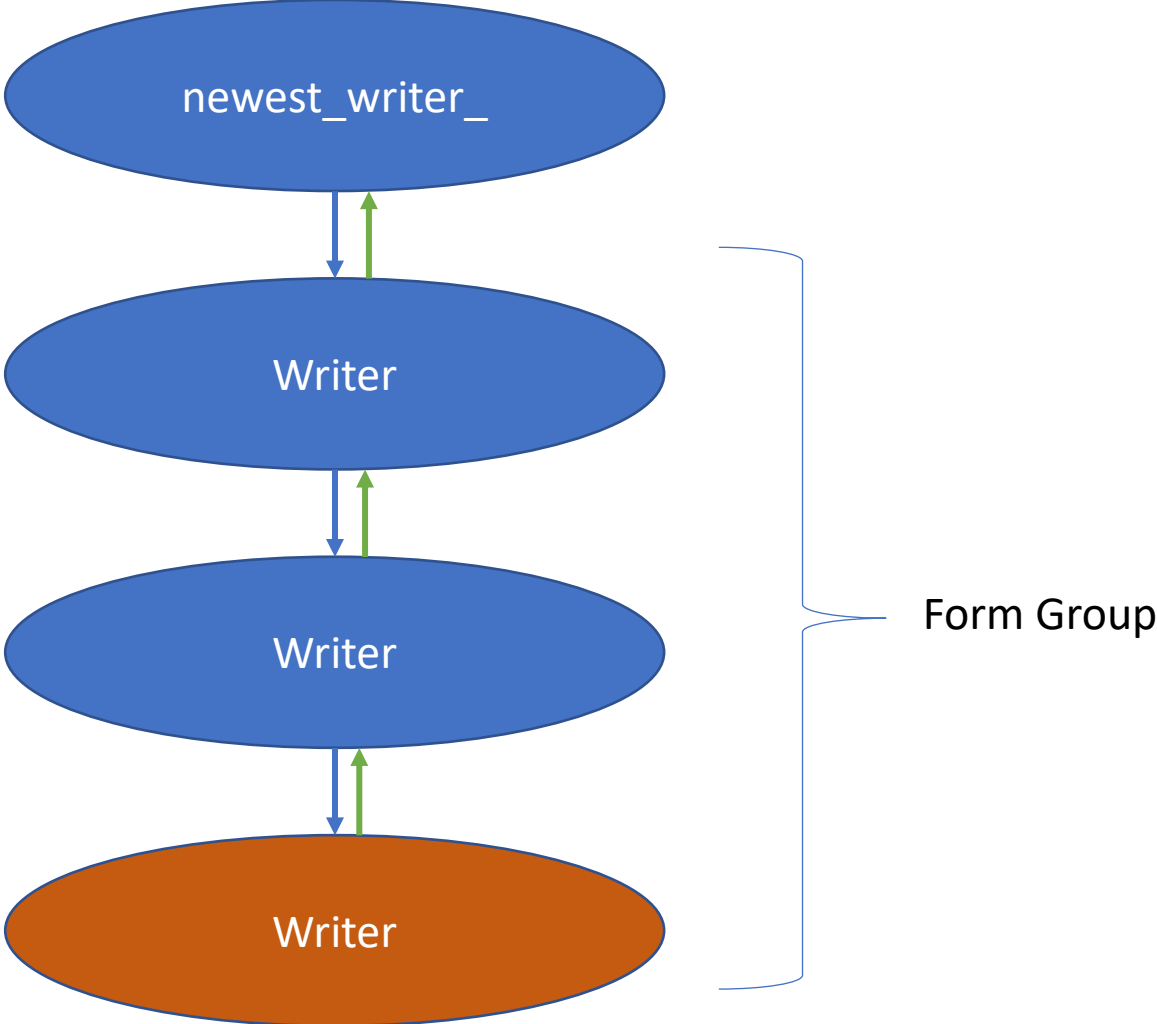
WriteThread::EnterAsBatchGroupLeader()



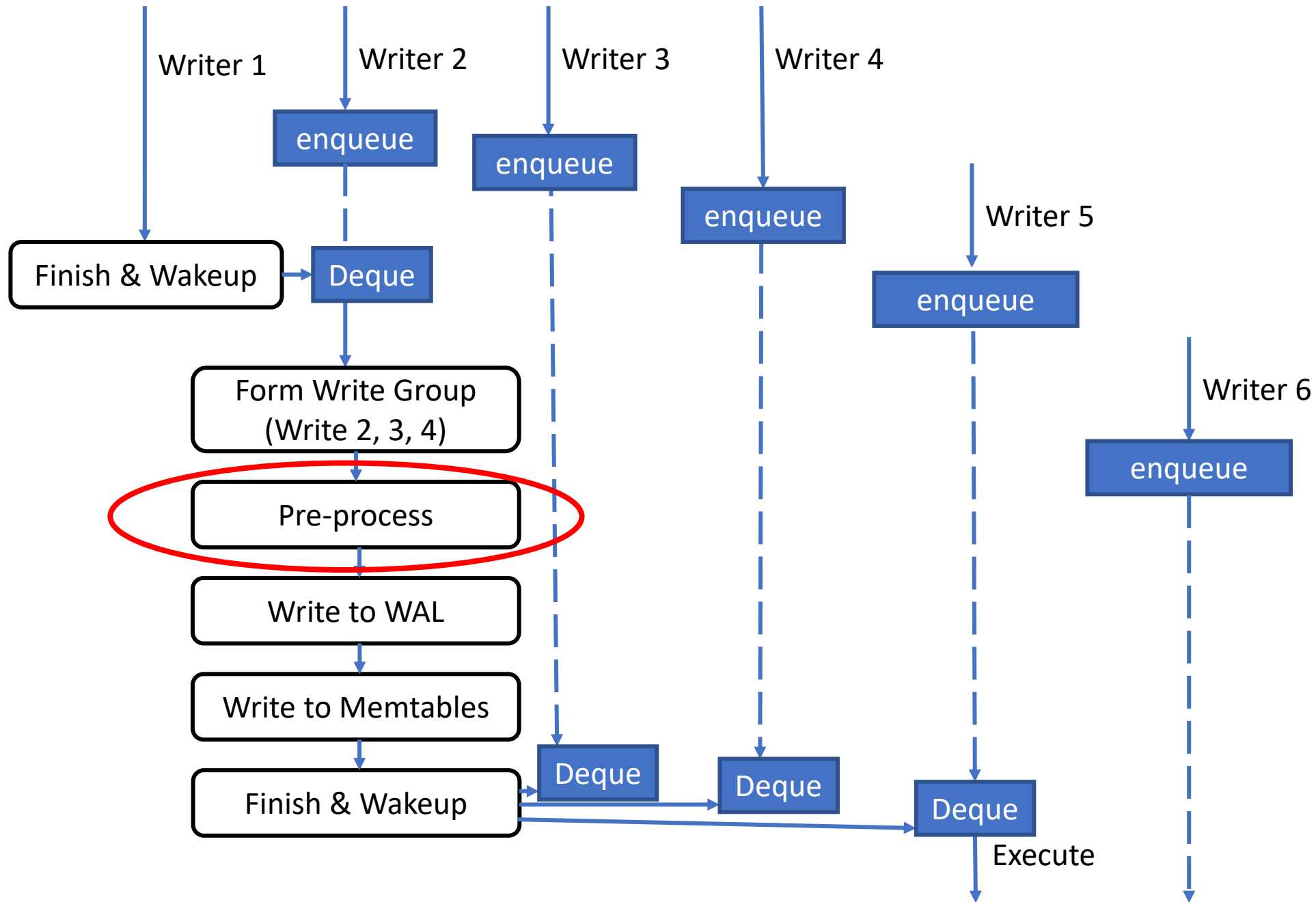
WriteThread::EnterAsBatchGroupLeader()



WriteThread::EnterAsBatchGroupLeader()



- Terminal Conditions:
- 1. Max group size
 - 2. Merge Operands

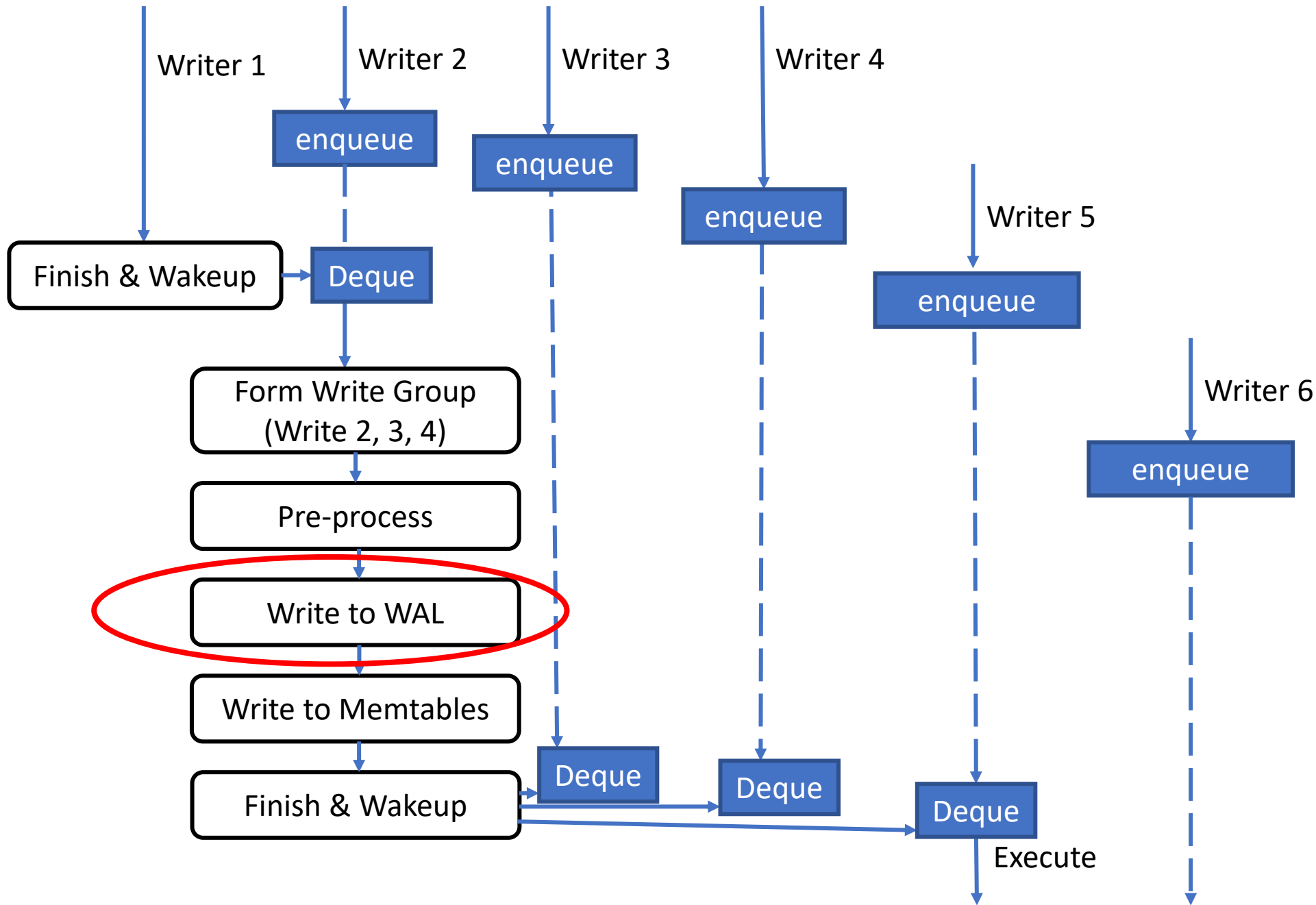


Pre-process: `DBImpl::PreprocessWrite()`

- Switch memtable when needed: `DBImpl::SwitchMemtable()`
 - `total_log_size_`
 - `write_buffer_manager_ ->ShouldFlush()`
 - Memtable full: `!flush_scheduler_.Empty()`
- Trim memtable history: `trim_history_scheduler_.Empty()`
- Throttling: checking `write_controller_`

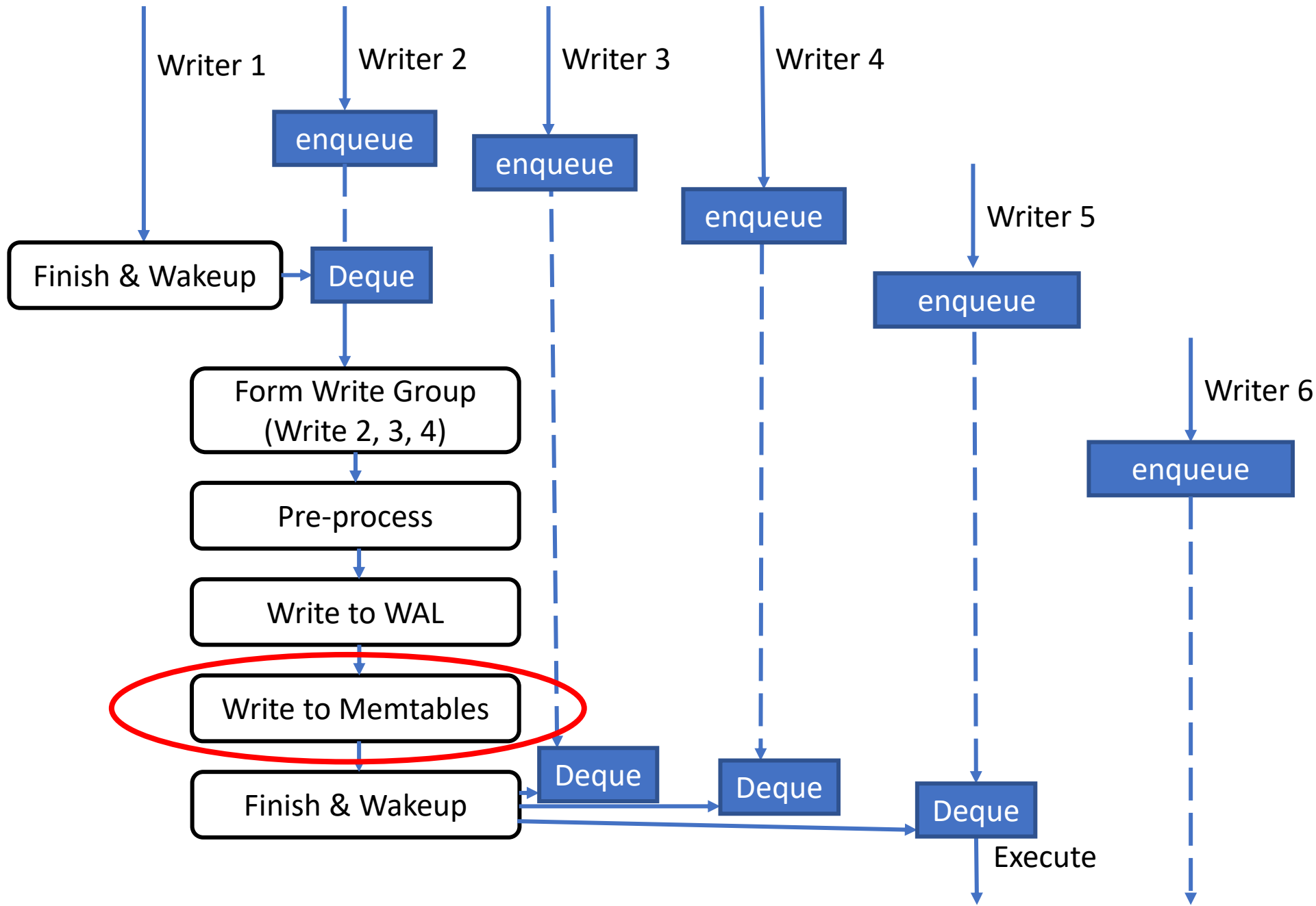
DBImpl::SwitchMemtable()

- Create new WAL file
- Create a new memtable
- Flush old WAL file buffer (to OS page cache)



Write to WAL

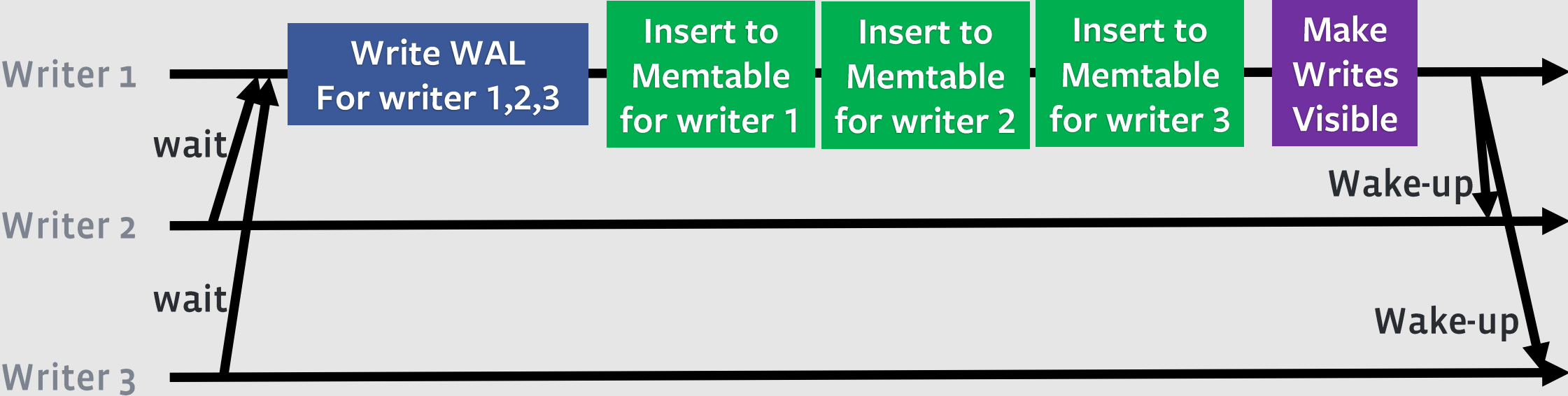
- Normal case: `DBImpl::WriteToWAL()`
- Unordered write: `DBImpl::WriteImplWALOnly()`
(write thread queuing is done there)
- `two_write_queues_`: `ConcurrentWriteToWAL()`
(synchronize through `log_write_mutex_`)



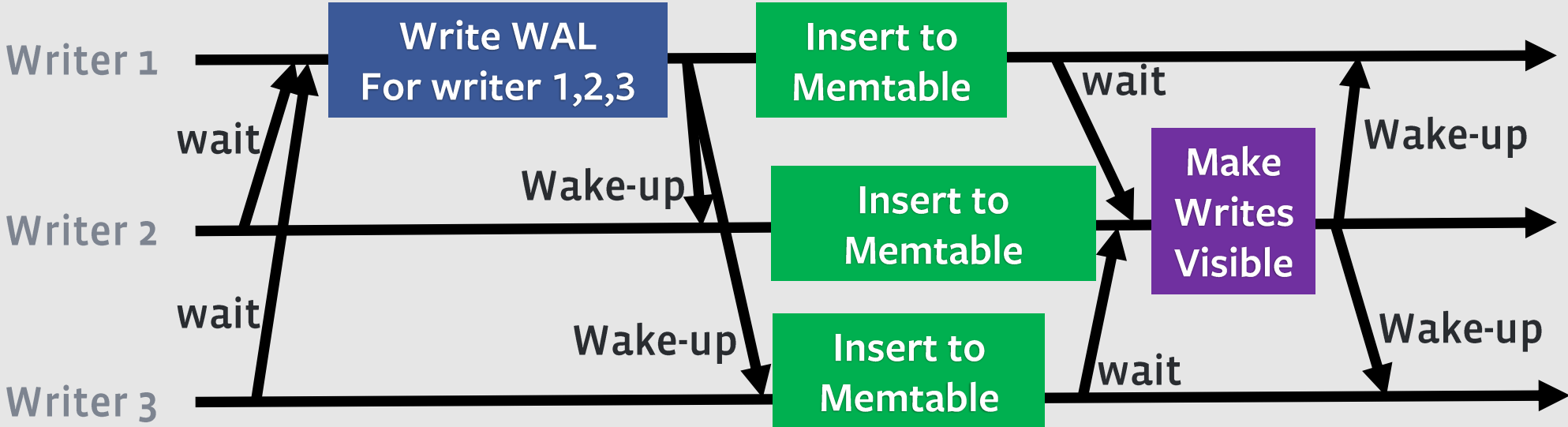
Memtable Writes: Three Modes

- Mode 1: group leader writes for all in the group
- Mode 2: concurrent memtable write
- Mode 3: unordered write: writes and return

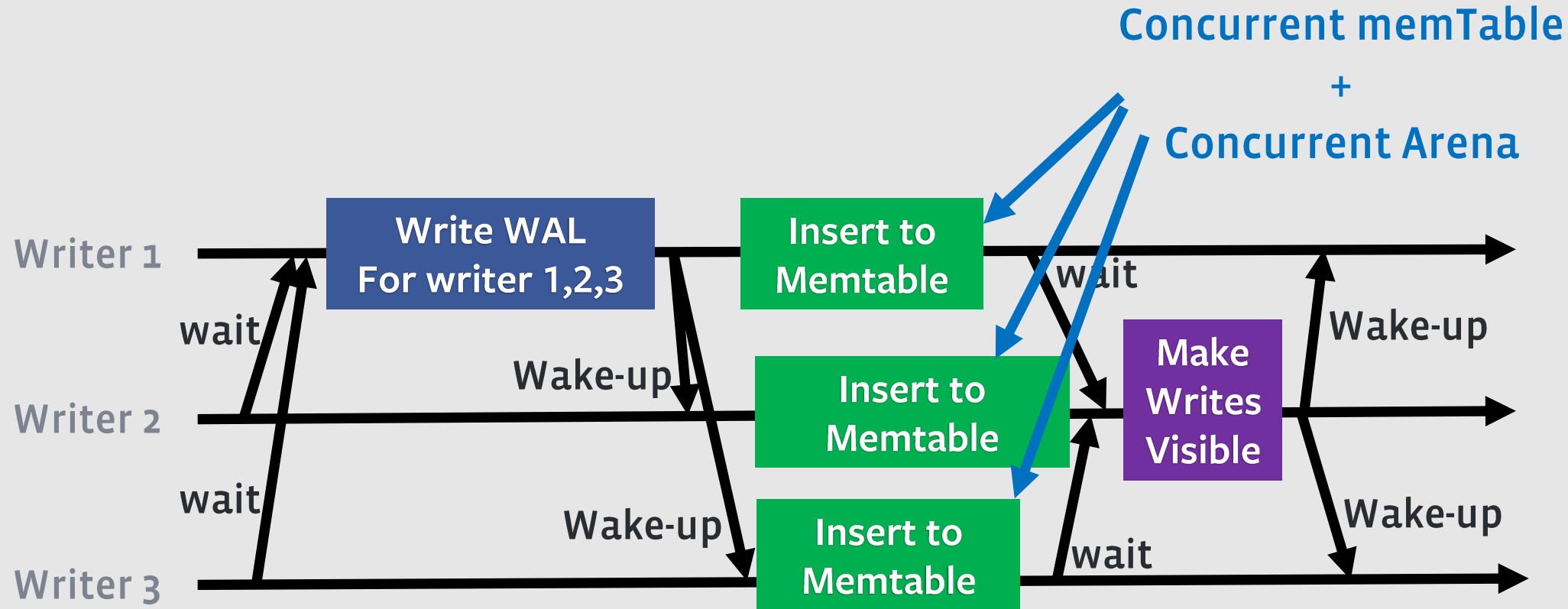
Mode 1: Group leader writes for all



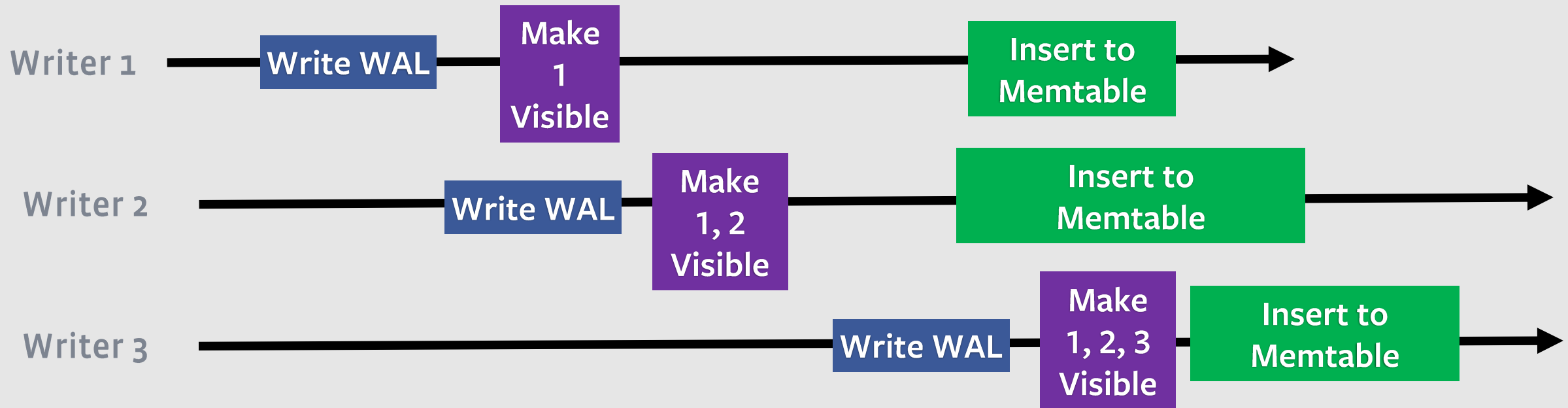
Mode 2: Concurrent memtable writes



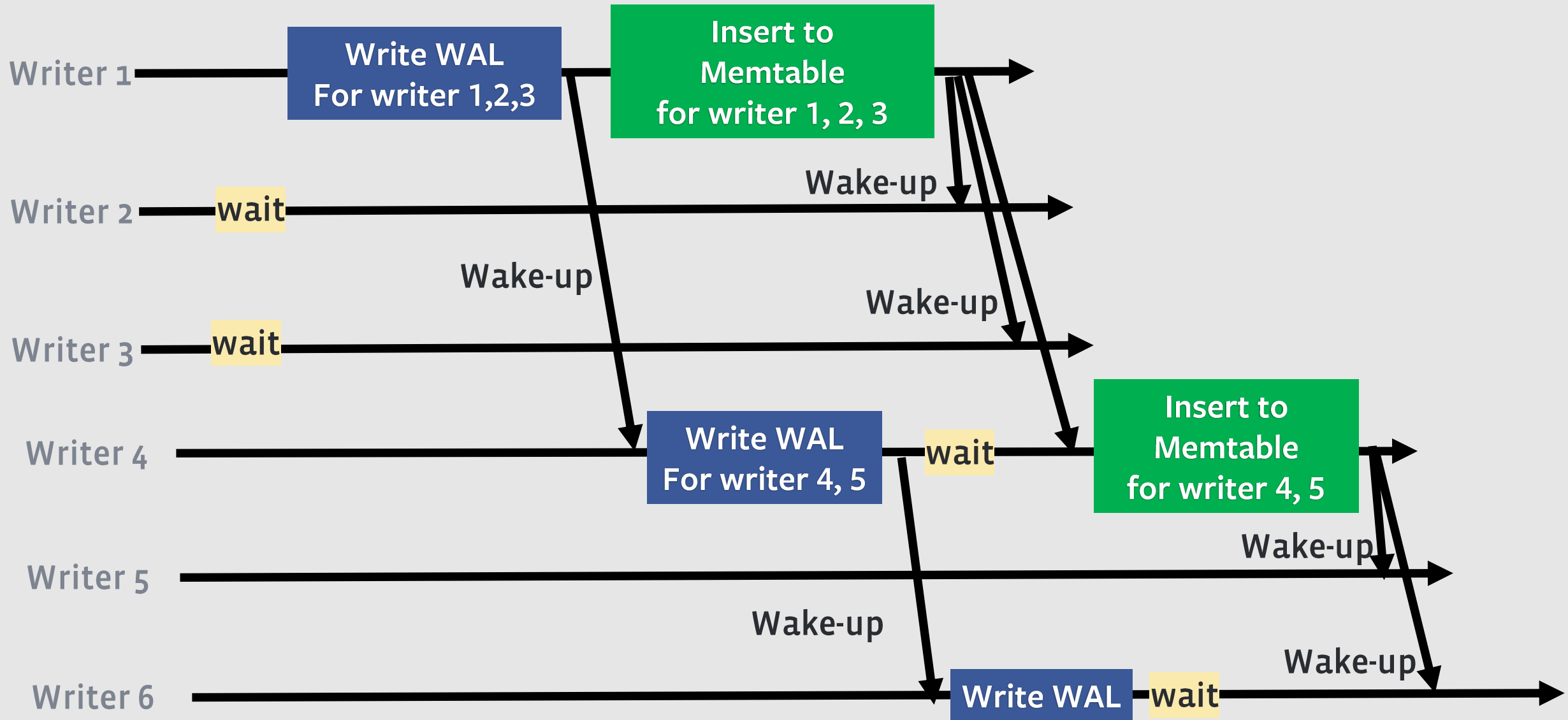
Mode 2: Concurrent memtable writes



Mode 3: Unordered Write



Pipelined Write (With Mode 1)



WriteBatchInternal::InsertInto()

- Go through every entry in the write batch:
 - Call `MemTable::Add()`
 - `MemTableInserter::CheckMemtableFull()`
 - Check memtable full
 - Check memtable history trimming condition

Two-Phase-Commit

	Write Committed	Write Prepared	Write UnPrepared
Before Prepare			<ol style="list-style-type: none">1. Prepare Entry to WAL2. Write to memtable3. Add uncommitted seqNum to tracking data structure <pre>WriteImpl(disable_memtable = false, add_prepared_callback)</pre>
Prepare	Prepare Entry to WAL	<ol style="list-style-type: none">1. Prepare Entry to WAL2. Write to memtable3. Add uncommitted seqNum to tracking data structure <pre>WriteImpl(disable_memtable = false, add_prepared_callback)</pre>	
Commit	Write to memtable Commit Entry to WAL	Commit Entry to WAL: <pre>WriteImpl(disable_memtable = true, update_commit_map_with_aux_batch)</pre>	Commit Entry to WAL: <pre>WriteImpl(disable_memtable = true, update_commit_map_with_aux_batch)</pre>

What's write_callback used for?

Optimistic Transaction's legacy conflict resolving mechanism.

`OptimisticTransaction::CommitWithSerialValidate()`

How should we deal with it?

Recap: Write Path

