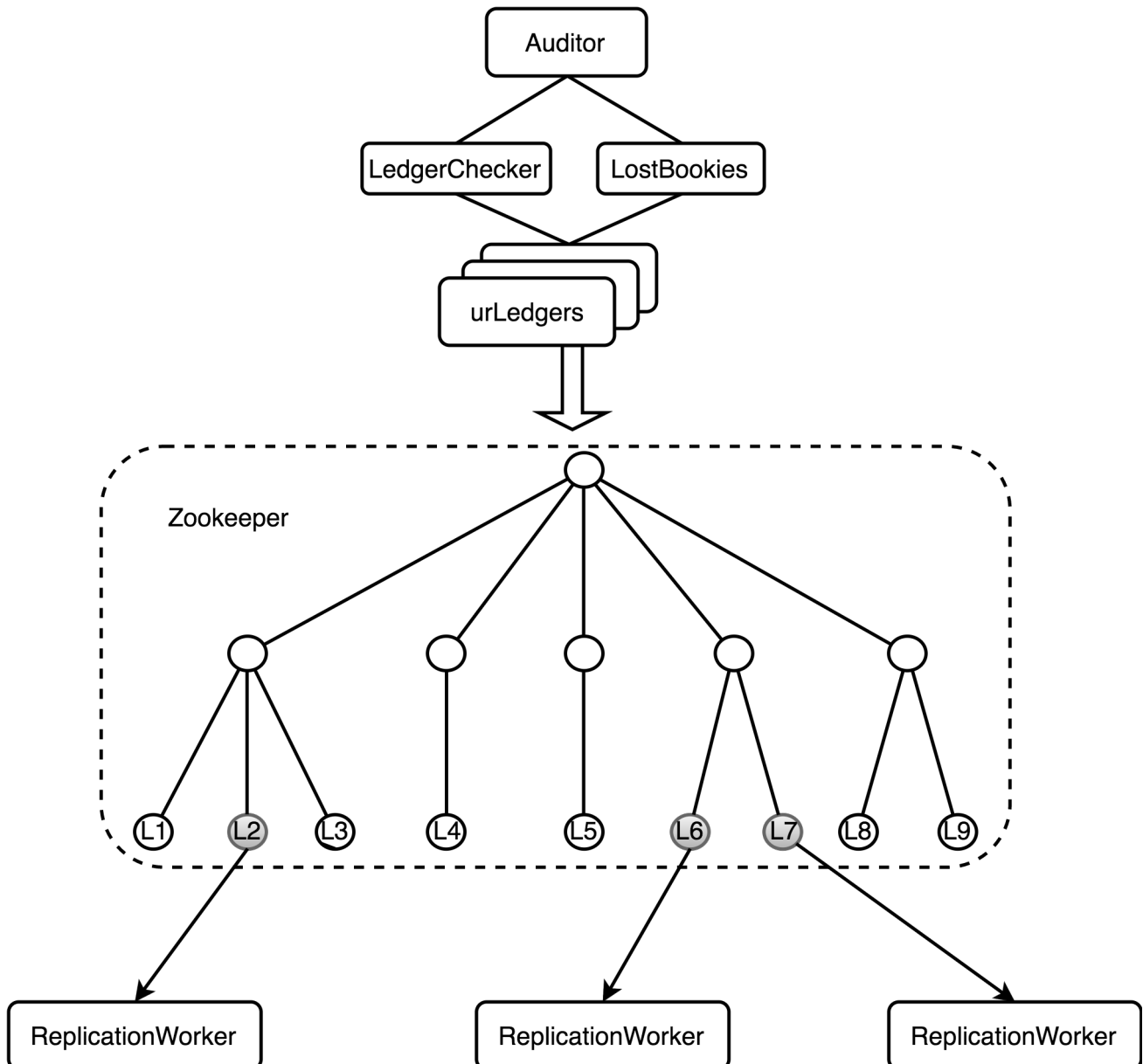


BP-8 - Queue based auto rereplicator

Current Design

Here's how the current Auditor and ReplicationWorker collaborate with each other.

1. Periodically, Auditor publishes underreplicated ledgers based on lost bookies and bad ledgers from LedgerChecker results.
2. Auditor publishes the underreplicated ledgers to Zookeeper, organized in hierarchical way, similar as how we store the bookkeeper ledgers in Zookeeper store.
3. ReplicationWorker poll an underreplicated ledger from Zookeeper by scanning the tree and acquiring the lock
4. ReplicationWorker check and rereplicate underreplicated ledgers
5. If replication succeed, ReplicationWorker mark the ledger as replicated, otherwise ReplicationWorker release the lock
6. Replication poll the next ledger to rereplicate



Problem

Auditor is not able to get replication result from ReplicationWorker. If there're ledgers failed to be replicated, Auditor will keep publishing these ledgers, and ReplicationWorker will keep replicating these ledgers over and over. This will generate lot of unnecessary traffic and generate duplicated data which might fill up the bookie disks.

Another minor existing issue is GC issue. Each time when ReplicationWorker tries to get the next underreplicated ledger, it needs to register a new watcher for each getChildren so that if there's currently no ledger to replicate, it'll wait for the notification from the watcher of a new ledger to replicate. However, this will sometimes cause memory leak. For example, if there always exists some underreplicated ledgers failed to replicate, and there's no new underreplicated ledgers published, the watcher will not be fired, so we will end up creating more and more watcher objects.

Proposed Design