

HBASE-6381 AssignmentManager should use the same logic for clean startup and failover

Purpose of the Patch

Currently, AssignmentManager has lots of ad-hoc changes to resolve individual issues. The overall logic is not easy to understand and maintain because there are lots of specific code to handle specific issue related to the failover case.

When a cluster starts up clean, things are simple and everyone is happy. However, during the failover mode,

- we need to restore the region states lost when a master dies.
- we need to respond those ZK events triggered by existing region servers in assigning/closing regions in transition.
- we need to handle some user region assignments started by the ServerShutdownHandler in case some region servers die in this wrong time.
- we need to handle tables in the middle of enabling/disabling.

As you can see, it is really busy. This patch is to clean it up, try to share some logic, and handle things one by one in order. For example,

- we use the ServerShutdownHandler to handle those dead region servers in the failover mode,
- we hold the ServerShutdownHandler before the region transition states are restored from the transition data in ZK,
- we don't listen to ZK events before we are ready,
- we share the same bulk assigner.

High Level Changes

I split this issue into four sub-tasks, hoped to make sure the patch is smaller. Unfortunately, it doesn't work out. So I ended up with a big patch. However, the ideas are in the four sub-tasks:

- **HBASE-6482** In AssignmentManager failover mode, use ServerShutdownHandler to handle dead regions

In AssignmentManager failover mode, a special failoverProcessedRegions map is used to manage regions in transition before the failover cleanup is completed. In failover mode, some regions may be still in transition. Once AssignmentManager starts to listen to ZK assign nodes, events will come in. However, since AssignmentManager already lost its region transition states, it needs some special handles.

For this patch, we use ServerShutdownHandler to process those regions. So that we can share some code and make the logic of AssignmentManager a little bit simpler.

More importantly, we don't listen to any ZK assignment events before we fully restore the region transition states based on the transition data in the ZK assign nodes.

- **HBASE-6483** Fully enable ServerShutdownHandler after master joins the cluster

Once ROOT and META are assigned, ServerShutdownHandler is enabled. So that we can handle meta/root region server failure before joinCluster is completed. However, we can hold ServerShutdownHandler a little bit more for the user region assignments, i.e. doesn't assign user regions (in ServerShutdownHandler) before joinCluster is ready. If so, we can avoid some region assignments racing: same regions are trying to be assigned in both joinCluster and ServerShutdownHandler.

We cannot enable ServerShutdownHandler after joinCluster is completed since META/ROOT region server could be down in the middle. In this patch, I introduced another queue in ServerManager to hold those dead region servers ServerShutdownHandler has already processed once: log splitting, ROOT/META reassignment, however, user regions are not processed yet. I'd like to enable the user region processing after the user region states are restored, i.e. the failover cleanup of AssignmentManger is completed. I don't want to wait in ServerShutdownHandler either since it holds the ServerShutdownHandler threads. If we run out of ServerShutdownHandler threads, new dead server will not be processed. So ROOT/META could not be assigned if the corresponding region server dies and all ServerShutdownHandler threads are used.

- **HBASE-6484** Make AssignmentManger#enablingTables and disablintTables local variables

Those enablingTables and disablingTables, are used only during the startup time to handle some table transistions which are happened when the master fails over. They should be some local variables.

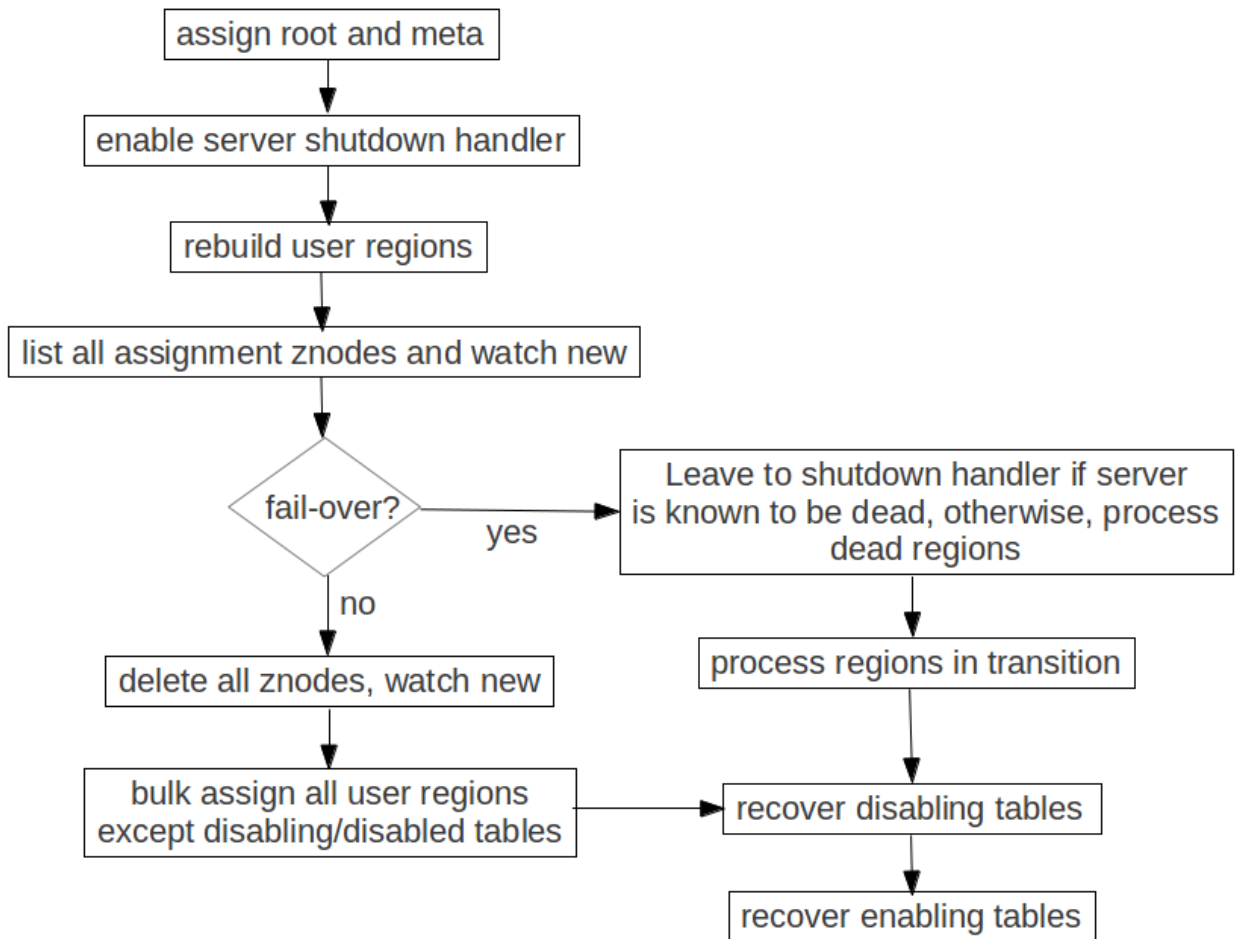
- **HBASE-6485** Share bulk assign code in AssignmentManager

AssignmentManager has several bulk assign functions: for startup bulk assigner, for ServerShutdownHandler bulk assign, etc. They can be shared.

There is a bulk assigner used by ServerShutdownHandler only, a bulk assigner used by AssignmentManger for startup bulk assigning and assigning multiple user regions, called StartupBulkAssigner. There is also a GeneralBulkAssigner which is not used at all.

In this patch, StartupBulkAssigner is merged to GeneralBulkAssigner. The bulk assign method used by ServerShutdownHandler is removed, some logic of which is folded into GeneralBulkAssigner. The GeneralBulkAssigner is enhanced to handle timeout properly.

How AssignmentManager Works Before the Patch



How AssignmentManager Works After the Patch

