

# HDFS16942

## Issue link

<https://issues.apache.org/jira/browse/HDFS-16942>

## Symptom

When a datanode sends a FBR to the namenode, it requires a lease to send it. On a couple of busy clusters, we have seen an issue where the DN is somehow delayed in sending the FBR after requesting the lease. Then the NN rejects the FBR and logs a message to that effect, but from the Datanodes point of view, it thinks the report was successful and does not try to send another report until the 6 hour default interval has passed.

If this happens to a few DNs, there can be missing and under replicated blocks, further adding to the cluster load.

## Root cause

When the namenode receives a block report, it firstly check if the lease is valid. It will not process the block report if the lease is not valid. As is shown in the code below, if `bm.checkBlockReportLease(context, nodeReg)` is true, the namenode will process the block report. However, if `bm.checkBlockReportLease(context, nodeReg)` is false, namenode will do nothing about the block report and datanode will not informed that its block report is rejected. Therefore, there may be missing blocks in the cluster.

```
1  try {
2      if (bm.checkBlockReportLease(context, nodeReg)) {#####
3          for (int r = 0; r < reports.length; r++) {
4              final BlockListAsLongs blocks = reports[r].getBlocks();
5              //
6              // BlockManager.processReport accumulates information of prior calls
7              // for the same node and storage, so the value returned by the last
8              // call of this loop is the final updated value for noStaleStorage.
9              //
10             final int index = r;
11             noStaleStorages = bm.runBlockOp(() ->
12                 bm.processReport(nodeReg, reports[index].getStorage(),
13                     blocks, context));
14         }
15     }
16 } catch (UnregisteredNodeException une) {
17     LOG.warn("Datanode {} is attempting to report but not register yet.",
18         nodeReg);
19     return RegisterCommand.REGISTER;
20 }
```

If there is some delay between datanode gets the lease and namenode processes block report, the lease will be expired so the value of `bm.checkBlockReportLease(context, nodeReg)` is false, which triggers this issue.

## Reproduce

Hadoop version: 3.1.3

To reproduce HDFS16942, we can set the configuration value `dfs.namenode.full.block.report.lease.length.ms` as 1ms by adding the following code into file `hdfs-site.xml`.

```

1 <property>
2   <name>dfs.namenode.full.block.report.lease.length.ms</name>
3   <value>1</value>
4   <description>The number of milliseconds that the NameNode will wait before invalidating a
full block report lease.</description>
5 </property>

```

Then it is very likely that the lease is expired when the namenode processes the block report from datanode, since making and sending block report take more than 1ms. To make the reproducing process deterministic, we can add some delay before the namenode processes block report:

```

1 public boolean checkBlockReportLease(BlockReportContext context, //#####
2   final DatanodeID nodeID) throws UnregisteredNodeException {
3   if (context == null) {
4     return true;
5   }
6   DatanodeDescriptor node = datanodeManager.getDatanode(nodeID);
7   if (node == null) {
8     throw new UnregisteredNodeException(nodeID, null);
9   }
10  + // add some delay here
11  + try {
12  +   System.out.println("### Sleep 5 seconds before checking lease");
13  +   Thread.sleep(5000);
14  + } catch (InterruptedException e) {
15  +   e.printStackTrace();
16  + }
17   final long startTime = Time.monotonicNow();
18   //check lease
19   return blockReportLeaseManager.checkLease(node, startTime,
20     context.getLeaseId());
21 }

```

By make the modifications above, every time we start the cluster, there is a WARN message `BR Lease 0x{} is not valid for DN {}`, because the lease has expired..

2023-07-06 14:16:45,905 WARN org.apache.hadoop.hdfs.server.blockmanagement.BlockReportLeaseManager: BR Lease 0xa5224a02e40ce2f4 is not valid for DN e31b77b3-e600-4db7-b6e7-dc5b9be8236b, because the lease has expired.

What's more, after uploading a file named `file_2M` and restarting the cluster, there is a missing block reported by the namenode and the location of `file_2M` can't be found.

There are 1 missing blocks. The following files may be corrupted: ✕

blk\_1073741825 /file\_2M

Please check the logs or run fsck in order to identify the missing blocks. See the Hadoop FAQ for common causes and potential solutions.

File information - file\_2M ✕

[Download](#)
[Head the file \(first 32K\)](#)
[Tail the file \(last 32K\)](#)

Block information -- Block 0

Block ID: 1073741825

Block Pool ID: BP-1027407363-127.0.1.1-1688624188512

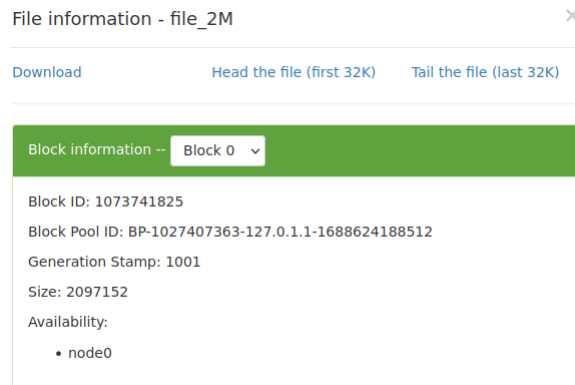
Generation Stamp: 1001

Size: 2097152

Availability:

Then we can trigger a full block report manually, which will always be accepted by namenode without expiry check. Now the missing block has been recovered.

```
1 bin/hdfs dfsadmin -triggerBlockReport localhost:9867
```



## Detailed steps to reproduce

Please make sure the version of Hadoop source code is 3.1.3

Reproducing this bug need only one server, which will run namenode and datanode at the same time.

Then install these packages if they can't be found in your servers.

```
1 sudo apt install git
2 sudo apt install maven
3 sudo apt-get install openjdk-8-jdk
4 sudo apt install curl
5 sudo apt install autoconf
6 sudo apt-get install libtool
7 sudo apt install make
8 sudo apt install g++
```

Then change the permission of scripts:

```
1 chmod -R 755 path_to_T2C/experiments/reproduce/HDFS-16942
```

Run `install_HDFS-16942.sh`. This script will apply the patch to modify source code and make the necessary configurations.

```
1 experiments/reproduce/HDFS-16942/install_HDFS-16942.sh [path_to_T2C] [path_to_HADOOP]
```

If there is a problem about the version of `protobuf` like `'protoc --version did not return a version'`, please run the script `build_protoc.sh`. After running this script, there should be a folder at the same directory as `build_protoc.sh`. Add `protoc` to PATH:

```
1 export PROTOC_HOME=/path_to_protoc/protoc/2.5.0
2 export HADOOP_PROTOC_PATH=$PROTOC_HOME/dist/bin/protoc
3 export PATH=$PROTOC_HOME/dist/bin/:$PATH
4 source ~/.bashrc
```

Then there should be no error message when running `install_HDFS-16942.sh`.

Then run `trigger_HDFS-16942.sh`.

```
1 experiments/reproduce/HDFS-16942/trigger_HDFS-16942.sh [path_to_HADOOP]
```

Now there should be a WARN message `BR lease 0x{} is not valid for DN{}, because the lease has expired.` in `logs/hadoop-username-namenode-hostname.log`. If so, the bug has been reproduced.

2023-07-20 02:00:34,766 WARN org.apache.hadoop.hdfs.server.blockmanagement.BlockReportLeaseManager: BR lease 0x9fdd38bc2a5197c7 is not valid for DN 025564f2-b79b-49dc-999d-56134dae81a1, because the lease has expired.