General *GFS* module. This module implements the general functioning of *GFS*, w/o focusing on write types. It specifies master functioning and adds some "time" functionality, but does not specify write/read-related "holes" in *PrimSecRep*.

EXTENDS *MCPrimSecRep*

CONSTANTS

$\qquad$ *Chunk*, $\qquad$ chunk *IDs*
$\qquad$ *RepFactor* $\quad$ replication factor

$gfsvars \;\triangleq\; mcvars$

Fill in some of the Holes in *PrimSecRep*:
$GFS\_ErrorTypes \;\triangleq\; \{\text{"badver"}\}$ $\;$ error types at secondaries:
$\qquad\qquad\qquad\qquad\qquad$ only one general type of error (or maybe none?)

$GFS\_TypeInvariant \;\triangleq\;$
$\qquad \wedge\; TypeInvariant$

$GFS\_Init \;\triangleq\;$
$\qquad \wedge\; MCInit$

Util functions
$chunk\_exists(c) \;\triangleq\; master.objects[c].version > 0$

$chunk\_full(c) \;\triangleq\;$ IF $\;master.objects[c].prim = NoRep\;$ THEN FALSE
$\qquad\qquad$ ELSE $\;Len(data[master.objects[c].prim][c]) = ChunkSize$
$chunk\_order \;\triangleq\; \langle FirstChunk,\, SecondChunk,\, ThirdChunk,\, FourthChunk,\, FifthChunk \rangle$
$prev\_chunk(c) \;\triangleq\;$
$\quad$ LET $idx \;\triangleq\;$ CHOOSE $n \in 1\,..\,Len(chunk\_order)\!: chunk\_order[n] = c$ IN $\;$ IF $\;(idx \text{ - } 1 = 0)\;$ THEN $0$
$\qquad$ ELSE $\;chunk\_order[idx-1]$

**MASTER** functioning

Chunk creation. Masster creates a chunk if it doesn't already exist. The new chunk is replicated at a number of replicas. These replicas must be believed as alive by the master. No lease is given out yet. Normally, chunk creation should be only when needed, but now I do it anytime, for simplicity.

$\_CreateChunk(c,\, reps) \;\triangleq\;$
$\qquad \wedge\; \forall\, r \in reps :$
$\qquad\qquad master.health[r] = \text{"alive"}\;$ master thinks that all proposed *reps* are alive

$\qquad$ Set the *reps* as secondaries. No replica has lease for this chunk yet.
$\qquad \wedge\; master' = [master$ EXCEPT $\,!.objects[c].prim = NoRep,\;$ no replica has lease yet
$\qquad\qquad\qquad\qquad\qquad\qquad\quad !.objects[c].sec = reps,$
$\qquad\qquad\qquad\qquad\qquad\qquad\quad !.objects[c].version = 1]$

$\qquad \wedge\;$ UNCHANGED $\langle data,\, cache,\, stat,\, channel,\, resps \rangle$
$CreateChunk \;\triangleq\;$

1

$\exists\, c \in Chunk :$
  $\land\, \neg chunk\_exists(c)$ new chunk does not exist

  $\land\, \lor\, prev\_chunk(c) = 0$
  $\phantom{\land}\, \lor\, \land\, prev\_chunk(c) \neq 0$
  $\phantom{\land\, \lor\, \land}\, \land\, chunk\_exists(prev\_chunk(c))$ \ * previous chunk already exists
  $\phantom{\land\, \lor\, \land}\, \land\, chunk\_full(prev\_chunk(c))$ \ * and is full

  $\land\, \exists\, reps \in \text{SUBSET } (Rep) :$ Find a set of replica to store the new chunk
  $\phantom{\land}\, \land\, card(reps) = RepFactor$
  $\phantom{\land}\, \land\, \_CreateChunk(c,\, reps)$

The primary releases a new lease for a chunk that doesn't have a primary. This release only happens after an older lease expires (so this action is logically after *LeaseExpiry*.

$\_ReleaseLease(c,\, r) \triangleq$
  $\land\, master.objects[c].prim = NoRep$ No lease given for chunk $c$
  $\land\, r \in master.objects[c].sec$ replica $r$ is one of secondaries
  $\land\, master.health[r] = \text{"alive"}$ replica $r$ is believed to be alive by master (redundant)

  Set secondary $r$ to be the primary of chunk $c$ (so it will no longer be a *sec* for that chunk)
  $\land\, master' = [master \text{ EXCEPT } !.objects[c].sec\ \ = @ \setminus \{r\},$
  $\phantom{\land\, master' = [master \text{ EXCEPT } }\, !.objects[c].prim = r,$
  $\phantom{\land\, master' = [master \text{ EXCEPT } }\, !.objects[c].version = @ + 1]$ get to a new version

  $\land\, \text{UNCHANGED } \langle data,\, stat,\, channel,\, cache,\, resps\rangle$
$ReleaseLease \triangleq$
  $\exists\, c \in Chunk :$
  $\phantom{\exists}\, \land\, chunk\_exists(c)$ the chunk exists
  $\phantom{\exists}\, \land\, master.objects[c].sec \neq \{\}$ there're replicas from where master can choose a new primary
  $\phantom{\exists}\, \land\, \exists\, r \in Rep :$
  $\phantom{\exists\, \land\, \exists}\, \_ReleaseLease(c,\, r)$

Master detects a replica is dead. This might be the reality or not. Events such as network partitions can lead to the primary thinking a replcia is dead when it's alive.

$\_DetectDeath(r) \triangleq$
  $\land\, master.health[r] = \text{"alive"}$ master thinks replica $r$ used to be alive

  Update all the chunks that the replica stored.
  If $r$ was a primary in any of these, force lease expiry
  before the master can release the lease to someone else.
  $\land\, master' =$
  $\phantom{\land}\, [objects \mapsto [c \in Chunk \mapsto$
  $\phantom{\land\, [objects \mapsto [}\, \text{IF } master.objects[c].prim = r \text{ THEN}$
  $\phantom{\land\, [objects \mapsto [\text{IF }}\, master.objects[c]$ DO NOT revoke lease of primary, wait for it to expire
  $\phantom{\land\, [objects \mapsto [}\, \text{ELSE IF } r \in master.objects[c].sec \text{ THEN}$
  $\phantom{\land\, [objects \mapsto [\text{ELSE }}\, [master.objects[c] \text{ EXCEPT } !.sec = @ \setminus \{r\}]$ invalidate a secondary
  $\phantom{\land\, [objects \mapsto [}\, \text{ELSE } master.objects[c]],$
  $\phantom{\land\, [}\, health \mapsto [master.health \text{ EXCEPT } ![r] = \text{"dead"}]]$

2

$\land$ UNCHANGED $\langle data,\ cache,\ stat,\ channel,\ resps \rangle$

$DetectDeath \triangleq$
    $\land \exists\, r \in Rep : \_DetectDeath(r)$

$GFS\_MasterActions \triangleq$
    $\lor CreateChunk \lor ReleaseLease \lor DetectDeath \lor MasterActions$

---

$\_LeaseExpiration(c,\ r) \triangleq$   lease of replica $r$ for chunk $c$ expires
    $\land master.objects[c].prim = r$   master knows $r$ used to have the lease

    $\land master' = [master \text{ EXCEPT } !.objects[c].prim = NoRep,$
                                    $!.objects[c].sec = \text{IF } master.health[r] = \text{"alive"}$
                                                           $\text{THEN } @ \cup \{r\} \text{ ELSE } @]$
    $\land cache' = [r1 \in Rep \mapsto [cache[r1] \text{ EXCEPT}$
                        $![c] = \text{IF } r = @.prim \text{ THEN}$
                                    $[prim \mapsto NoRep,$
                                     $sec \mapsto \text{IF } stat[r].phase = \text{"alive" THEN } @.sec \cup \{r\}$
                                                   $\text{ELSE } @.sec,$
                                     $version \mapsto @.version]$
                                 $\text{ELSE } @]]$
    $\land$ UNCHANGED $\langle data,\ stat,\ channel,\ resps \rangle$

$LeaseExpiration \triangleq$
    $\land \exists\, c \in Chunk,\ r \in Rep : \_LeaseExpiration(c,\ r)$

$GFS\_TimeActions \triangleq$
    $\lor LeaseExpiration \lor TimeActions$

---

$GFS\_ReplyFinishedWrite(r,\ c,\ w) \triangleq$
    LET $CommittedEverywhere \triangleq$   the write has been committed to all replicas
                                          believed alive at the time the write was initiated
            $\forall\, r1 \in Rep : \lor resps[r][c][r1] = \text{"ok"}$   $r1$ responded $ok$
                                 $\lor resps[r][c][r1] = \text{"n/a"}$   OR $r$ wasn't expecting a reply from $r1$
    IN
        IF $CommittedEverywhere$ THEN
                $Reply(r,\ \text{"wrFinished"},\ [ack \mapsto \text{"ok"},\ object \mapsto c,\ w \mapsto w])$
        ELSE

$$Reply(r, \text{"wrFinished"}, [ack \mapsto \text{"error"}, object \mapsto c, w \mapsto w])$$

---

**SPECIFICATION**

$GFS\_Next \triangleq$
$\quad\quad \vee\ GFS\_MasterActions$
$\quad\quad \vee\ GFS\_TimeActions$
$\quad\quad \vee\ ReplicaActions$
$\quad\quad \vee\ ClientActions$

$GFS\_Spec \triangleq\ \wedge\ GFS\_Init \wedge \Box[GFS\_Next]_{gfsvars}$

---

**Invariants**

DEBUG: Primary is unique. If a replica believes it's prim for a chunk, the master agrees with this.

$PrimIsFreshAndUniqueI \triangleq$
$\quad \wedge\ \forall\, c \in Chunk :$

$\quad\quad$ If any replica believes it's primary, then it's really primary
$\quad\quad \wedge\ \forall\, r \in Rep : (cache[r][c].prim = r \Rightarrow master.objects[c].prim = r)$

$\quad\quad$ the primary knows it's primary
$\quad\quad \wedge\ \vee\ master.objects[c].prim = NoRep \setminus *$ either no primary
$\quad\quad\quad\ \vee\ stat[master.objects[c].prim].phase = \text{"dead"}$
$\quad\quad\quad\ \vee\ master.health[master.objects[c].prim] = \text{"dead"}$
$\quad\quad\quad\ \vee\ cache[master.objects[c].prim][c].prim = master.objects[c].prim$
$\quad\quad\quad\quad\quad \setminus *$ OR the prim knows it's prim

DEBUG: No replica that master believes as dead is a secondary. Note that the same does NOT hold for primaries!! There can be a primary that's known dead to master but that the master still has registered as prim until its lease expires.

$MasterCorrectnessI \triangleq$
$\quad \wedge\ \forall\, r \in Rep :$
$\quad\quad (master.health[r] = \text{"dead"} \Rightarrow \neg(\exists\, c \in Chunk : r \in master.objects[c].sec))$

DEBUG: At no time is the primary a secondary.

$PrimIsNotSecI \triangleq$ a primary is not a secondary for the same chunk
$\quad \wedge\ \forall\, c \in Chunk :$
$\quad\quad \vee\ \neg chunk\_exists(c)$
$\quad\quad \vee\ master.objects[c].prim \notin master.objects[c].sec$

$GFS\_AllInvariants \triangleq$
$\quad \wedge\ GFS\_TypeInvariant$
$\quad \wedge\ PrimIsFreshAndUniqueI$
$\quad \wedge\ MasterCorrectnessI$
$\quad \wedge\ PrimIsNotSecI$

---

THEOREM $GFS\_Spec \Rightarrow \Box GFS\_AllInvariants$

4