

EXTENDS *Integers, Sequences*

Spec for a Write Ahead Logger (*WAL*)

WAL Logger appends data to a file called write extent (*WE*) which has max capacity of 1 *GB*.

Once write extent is full, we move it to a set of read files called read extents (*RE*) and create new *WE*.

Each extent/file has data in range from start to end *LSN* (logical sequence numbers).

extent: [start, end) where $start \leq end$.

means that extent contains data from [start, end) //end not included

$start \triangleq end$ means no data in file.

Data is appened at the end of log called tail. So, data grows from head towards tail.

Other operations performed by logger are truncation of log at head or tail,

Close and recovery of log at Open.

State machine valid steps:

1. There is only 1 *Append* request at a time.
2. After close or crash, only recovery runs.
3. Only 1 *TruncateHead* request is initiated at a time.
4. Only 1 *TruncateTail* request is initiated at a time.
5. *TruncateHead* can be initiated with *Append*.
6. No other request is served till *TruncateTail* finishes.

Todo:

- 1) *MetadataFile* corruption is single point of failure.

Todo: Create *metadataFile.mdlog.tmp* file first,

then delete *metadataFile.mdlog* and rename *.mdlog.tmp* to *.mdlog*

Handle crash after each step in *Recovery*.

Variables are divided into 2 categories:

1. Variables representing on disk data structures

MAK: This reads a lot like as if E_{-} are history variables (<https://arxiv.org/pdf/1703.05121.pdf>). It usually makes the spec mo

2. Variables representing expected values, prefixed by E_{-}

This means that during recovery, we can't use E_{-} * variables.

E_{-} * variables are used in Invariants to check that disk state is correct.

VARIABLES

MAK: Like with *Vladimir's* spec of *Floyd's* algorithm, you should see if you really need this variable.

PrevState, For specifying state machine

WE, $WE \Rightarrow$ current file to which logger is appending.

REs, $REs \Rightarrow$ sequence of read only files, which became read only after they were full.

MetadataFile, metadata file storing metadata information of the logger.

$E_{-}LowLSN$, lowest valid *LSN* of the logger

Data is stored from [$E_{-}LowLSN$, $E_{-}HighLSN$)

$E_{-}HighLSN$, next valid *LSN* of the logger. *i.e.* not including $E_{-}HighLSN$

$E_{-}THIP$, *TruncateHeadInProgress* : Is *TruncateHead* in progress ?

Truncate head means removing old data before offset *Head*

E_TTIP , $TruncateTailInProgress$: Is $TruncateTail$ in progress ?

Truncate tail means removing current data from end of file.

This is sometimes needed if we want to reset the file or other false progress cases.

Note: New data is appened at tail of the log/file.

E_NWEIP , WE is full and Is New WE creation In Progress ?

New_WE , New_WE file while E_NWEIP is TRUE

$TornWrite$, Did last crash caused torn write ?

MAK: It is better to use state or action constraints (<https://tla.msr-inria.inria.fr/tlatoolbox/doc/model/spec-optio>)

$MaxNum$ Variable to restrict TLC Model Checker (MC) to $MaxNum$ steps.

$TypeOK \triangleq$

$\wedge WE \in [id : 1 .. MaxNum, start : 1 .. MaxNum, end : 1 .. MaxNum, version : 1 .. MaxNum]$

$\wedge WE.start \leq WE.end$

$\wedge New_WE \in [exist : \{TRUE, FALSE\}, id : 1 .. MaxNum, start : 1 .. MaxNum, end : 1 .. MaxNum, version : 1 .. MaxNum]$

$\wedge New_WE.start \leq New_WE.end$

$\wedge REs \in Seq([id : 1 .. MaxNum, start : 1 .. MaxNum, end : 1 .. MaxNum, version : 1 .. MaxNum])$

$\wedge E_LowLSN \in 1 .. MaxNum$

$\wedge E_HighLSN \in 1 .. MaxNum$

$\wedge PrevState \in \{ "start", "append", "WE_full_New_WE", "New_WE_in_MDT", "crash", "recovery", "close", "truncate_head_p1", "truncate_head", "truncate_tail_p1", "truncate_tail" \}$

$MetadataFile \Rightarrow$ Stores $headLSN$, $tailLSN$, $tailVersionNum$, $fileNames$

When a new file is created after last WE fills up, it's entry is added in $metadataFile$

When Truncation happens, head and tail are updated in $metadataFile$

Recovery uses $metadataFile$ for knowing list of valid files in log

$headLSN$ corresponds to current E_LowLSN .

$lastTailLSN$ corresponds to last tail truncation.

$lastTailVersion$ is needed to know if we crashed during truncating tail or new data has been added after last tail truncation.

This is needed because we don't update metadata file on every write to log WE file.

$\wedge MetadataFile \in [headLSN : 1 .. MaxNum, lastTailLSN : 1 .. MaxNum, lastTailVersion : 1 .. MaxNum, cleanShutdown : \{TRUE, FALSE\}, fileIds : Seq(1 .. MaxNum)]$

$\wedge TornWrite \in \text{BOOLEAN}$

$\wedge E_THIP \in \{TRUE, FALSE\}$

$\wedge E_TTIP \in \{TRUE, FALSE\}$

$\wedge E_NWEIP \in \{TRUE, FALSE\}$

MAK: This says that $MaxNum$ is always 7. Why isn't this a constant?

$\wedge MaxNum = 7$

Initial state of the system.

$Init \triangleq$

$\wedge REs = \langle \rangle$

$\wedge WE = [id \mapsto 1, start \mapsto 1, end \mapsto 1, version \mapsto 1]$

$\wedge New_WE = [exist \mapsto FALSE, id \mapsto 1, start \mapsto 1, end \mapsto 1, version \mapsto 1]$

$\wedge E_LowLSN = 1$
 $\wedge E_HighLSN = 1$
 $\wedge PrevState = \text{"start"}$
 $\wedge MetadataFile = [headLSN \mapsto 1, lastTailLSN \mapsto 1, lastTailVersion \mapsto 1,$
 $cleanShutdown \mapsto \text{FALSE}, fileIds \mapsto \langle 1 \rangle]$
 $\wedge TornWrite = \text{FALSE}$
 $\wedge E_THIP = \text{FALSE}$
 $\wedge E_TTIP = \text{FALSE}$
 $\wedge E_NWEIP = \text{FALSE}$
 $\wedge MaxNum = 7$

Helper functions – begin
 Don't use E_* variables in helper functions.
 $GetFileIds(files) \triangleq$
 $[i \in 1 \dots Len(files) \mapsto files[i].id]$
 $GetMetadataFiles \triangleq$
 $LET PresentInMetadataFiles(r) \triangleq$
 $LET SameId(r2Id) \triangleq r.id = r2Id$
 $IN Len(SelectSeq(MetadataFile.fileIds, SameId)) > 0$
 $IN SelectSeq(Append(REs, WE), PresentInMetadataFiles)$
 $GetValidFiles(files, lowLSN, highLSN) \triangleq$
 $LET ValidFile(f) \triangleq \wedge f.start < f.end$
 $\wedge \neg(\vee (f.end \leq lowLSN) \vee (f.start > highLSN))$
 $IN SelectSeq(files, ValidFile)$
 Helper functions – end

Append keeps appending to WE increasing end LSN .
 No writes allowed while we do Truncate tail.
 Writes allowed while we do Truncate head.
 $AppendToFile \triangleq$
 Append to file is always allowed except close or crash.
 After crash/close, we first do recovery.
 $\wedge PrevState \notin \{\text{"crash"}, \text{"close"}\}$
 $\wedge E_NWEIP = \text{FALSE}$ WE is not full
 No writes allowed while $truncate_tail$ is in progress.
 $\wedge E_TTIP = \text{FALSE}$
 $\wedge E_HighLSN < MaxNum - 1$ Stop TLC model checker to generate more cases.
 $\wedge WE' = [WE \text{ EXCEPT } !.end = WE.end + 1,$
 Every write needs a metadata header with version number
 Next write after $TruncateTail$ will append to file with new version number. Thanks TLA+
 $!.version = MetadataFile.lastTailVersion]$
 $\wedge E_HighLSN' = E_HighLSN + 1$ Ack to customer that write succeeded
 $\wedge PrevState' = \text{"append"}$
 $\wedge \text{UNCHANGED } \langle E_LowLSN, MaxNum, REs, MetadataFile, TornWrite, E_THIP, E_TTIP, E_NWEIP,$

Action : Write extent is full - Create new Write Extent/file

1. Create a new *WE* file with right header and append data

Next Action: *NewWriteExtentAddToMetadataFile* :

Add to metadata file and move *WE* to RE in memory and $WE' = new_WE$

These 2 steps are divided in separate Actions to simulate crash and concurrency with different actions.

If we crash before updating *Metadata* file, we ignore this write and *New_WE* file is deleted on recovery.

Todo: Allow *WriteExtentFullNewWE* to run after recovery. Add a field in file : full : {TRUE, FALSE}

$WriteExtentFullNewWE \triangleq$

$\wedge \vee PrevState = \text{"append"}$

$\vee \wedge E_THIP = TRUE$ Truncate *Head* is allowed concurrent with writes.

$\wedge PrevState \neq \text{"crash"}$

$\wedge WE.id < MaxNum - 1$ Stop *MC* after these steps

No writes allowed while *truncate_tail* is in progress.

$\wedge E_TTIP = FALSE$

$\wedge E_NWEIP = FALSE$ Only *WE* full workflow at a time

$\wedge E_NWEIP' = TRUE$ Stop appends to *WE*

Create new *WE*

$\wedge New_WE' = [exist \mapsto TRUE, id \mapsto WE.id + 1, start \mapsto WE.end, end \mapsto WE.end + 1,$

Next write after *TruncateTail* will append to file with new version number. Thanks TLA+
 $version \mapsto MetadataFile.lastTailVersion]$

$\wedge PrevState' = \text{"WE_full_New_WE"}$

$\wedge UNCHANGED \langle E_LowLSN, E_HighLSN, MaxNum, REs, WE, MetadataFile, TornWrite, E_THIP, E_TTIP \rangle$

Add new write extent file to *MetadataFile* and open for new appends

$NewWriteExtentAddToMetadataFile \triangleq$

$\wedge E_NWEIP = TRUE$

$\wedge PrevState \neq \text{"crash"}$ only recovery runs after crash

$\wedge E_HighLSN < MaxNum - 1$ Stop *TLC* model checker (*MC*) to generate more cases to finish *MC*.

No writes allowed while *truncate_tail* is in progress.

$\wedge E_TTIP = FALSE$ How to assert that *E_TTIP* is false ?

First change *MetadataFile* on disk:

It might seem that, We could have easily done $!fileIds = Append(MetadataFile.fileIds, New_WE.id)$

but If you do that , invariant *AllMetadataFilesPresentOnDisk* will fail.

i.e. You will see that we can have file entries in *MetadataFile* which don't exist on disk.

This can be because of concurrency with *TH*. Thanks to TLA+.

Don't use *E_LowLSN*, *E_HighLSN* for *MetadataFile*:

Idea is not use Low/High *LSN* for changing disk data structures. They are maintained paralelly and used for assertion in Invariants.

$\wedge LET validFiles \triangleq GetValidFiles(Append(REs, WE), MetadataFile.headLSN, WE.end)$

$\wedge MetadataFile' = [MetadataFile EXCEPT !.fileIds = Append(GetFileIds(validFiles), New_WE.id)]$

Todo: Split changing RE in separate action to see what concurrency can do.

In-memory data structure change

Because of concurrency in *TH*, it is possible to get Truncation till last *WE* while we are adding new *WE*

$\wedge REs' = GetValidFiles(validFiles, E_LowLSN, E_HighLSN)$
 $\wedge WE' = [id \mapsto New_WE.id,$
 $\quad start \mapsto New_WE.start,$
 $\quad end \mapsto New_WE.end,$
 $\quad version \mapsto New_WE.version]$
 Reset other fields
 $\wedge New_WE' = [New_WE \text{ EXCEPT } !.exist = FALSE]$
 $\wedge E_HighLSN' = E_HighLSN + 1$ Ack to customer that write succeeded
 $\wedge PrevState' = \text{"New_WE_in_MDT"}$
 $\wedge E_NWEIP' = FALSE$ allow appends to WE now
 $\wedge UNCHANGED \langle E_LowLSN, MaxNum, TornWrite, E_THIP, E_TTIP \rangle$

Crash: torn write : last write ignored

We can't have torn write in case of New_WE , as only after write is successful, we update metadata and ack to caller.

$CrashWhileAppend \triangleq$

$\wedge PrevState = \text{"append"}$
 $\wedge PrevState' = \text{"crash"}$
 $\wedge E_HighLSN' = E_HighLSN - 1$ Simulate : we crashed before acking to customer
 $\wedge TornWrite' = TRUE$
 don't change metadata file as we can't do it during crash
 Invariant : *CleanShutdownOnlyAfterClose* makes sure that we have clean bit set only after close
 $\wedge MetadataFile' = [MetadataFile \text{ EXCEPT } !.cleanShutdown = FALSE]$
 $\wedge UNCHANGED \langle E_LowLSN, MaxNum, REs, WE, E_THIP, E_TTIP, E_NWEIP, New_WE, MetadataFile \rangle$

Normal crash that does not cause data loss

$CrashNoDataLoss \triangleq$

$\wedge PrevState \notin \{\text{"crash"}, \text{"close"}\}$ we can't crash after close as we aren't running
 $\wedge PrevState' = \text{"crash"}$
 don't change metadata file as we can't do it during crash
 Invariant : *CleanShutdownOnlyAfterClose* makes sure that we have clean bit as TRUE only after close
 $\wedge MetadataFile' = [MetadataFile \text{ EXCEPT } !.cleanShutdown = FALSE]$
 $\wedge UNCHANGED \langle E_LowLSN, MaxNum, E_HighLSN, REs, WE, TornWrite, E_THIP, E_TTIP, E_NWEIP \rangle$

Close the log file.

Waits for all operations to finish on log and sets the clean shutdown bit to true.

$Close \triangleq$

$\wedge PrevState \notin \{\text{"crash"}, \text{"close"}\}$
 Close waits for workflows to finish:
 $New_WE, truncate_head, truncate_tail$
 $\wedge E_NWEIP = FALSE$
 $\wedge E_TTIP = FALSE$
 $\wedge E_THIP = FALSE$
 $\wedge PrevState' = \text{"close"}$
 $\wedge MetadataFile' = [MetadataFile \text{ EXCEPT } !.cleanShutdown = TRUE,$
 $\quad !.lastTailLSN = WE.end]$

\wedge UNCHANGED $\langle E_LowLSN, MaxNum, E_HighLSN, REs, WE, TornWrite, E_THIP, E_TTIP, E_NWE \rangle$

Action: *Recovery* : It happens on Open

After crash, we can't look at value of E_* variables to know the state of the system before close/crash.

After recovery, we set cleanshutdown bit in metadatafile to false, which can only be set to TRUE during close.

We don't do anything if clean shutdown bit is set in the metadata file - fast open case.

Recovery \triangleq

$\wedge \vee PrevState = \text{"crash"}$

$\vee PrevState = \text{"close"}$

\wedge IF *MetadataFile*.cleanShutdown

THEN $\wedge REs' = REs$

$\wedge WE' = WE$

$\wedge MetadataFile' = [MetadataFile \text{ EXCEPT } !.cleanShutdown = FALSE]$

Otherwise we crashed.

Read metadata file and check *WE/REs* file to rebuild the state.

ELSE LET *allFiles* \triangleq *GetMetadataFiles*

lowLSN \triangleq *MetadataFile*.headLSN

last file in metadatafile is supposed to be *WE*

Todo : add a flag in metadata file for which file is *WE*.

lastWE \triangleq LET *lastWEId* \triangleq *MetadataFile*.fileIds[Len(*MetadataFile*.fileIds)]

SameId(*r*) \triangleq *r*.id = *lastWEId*

IN *Head*(*SelectSeq*(*allFiles*, *SameId*))

highLSN: Thanks TLA+

case : Crash during *append* – *TornWrite* : Last *LSN* in write file.

case : Crash during *TruncateTail phase1*

LastTailLSN : if version of *WE* is < *metadataFile*'s version

We can't append while *TruncateTail* is going on,

so we can't have both cases occurring at same time together.

highLSN \triangleq LET *lastValidWrite* \triangleq *lastWE*.end

IN IF *TornWrite*

THEN *lastValidWrite* – 1

ELSE IF *lastWE*.version < *MetadataFile*.lastTailVersion

THEN *MetadataFile*.lastTailLSN

ELSE *lastValidWrite*

goodExtents \triangleq *GetValidFiles*(*allFiles*, *lowLSN*, *highLSN*)

cleanState \triangleq Len(*goodExtents*) = 0

IN IF *cleanState* if we don't have any data

THEN $\wedge REs' = \langle \rangle$

$\wedge WE' = [id \mapsto 1, start \mapsto lowLSN, end \mapsto highLSN,$
 $version \mapsto MetadataFile.lastTailVersion]$

Not setting clean bit to false, as it is expected to be false because we crashed.

$\wedge MetadataFile' = [MetadataFile \text{ EXCEPT } !.fileIds = \langle 1 \rangle]$

ELSE $\wedge REs' = SubSeq(goodExtents, 1, Len(goodExtents) - 1)$

$\wedge WE' = [goodExtents[Len(goodExtents)]]$

EXCEPT $!.end = highLSN,$
 $!.version = MetadataFile.lastTailVersion]$
 $\wedge MetadataFile' = [MetadataFile \text{ EXCEPT } !.fileIds = GetFileIds(goodExtents)]$
 Reset variables correctly so that appends can work.
 $\wedge PrevState' = \text{"recovery"}$
 $\wedge E_THIP' = \text{FALSE}$
 $\wedge E_TTIP' = \text{FALSE}$
 $\wedge E_NWEIP' = \text{FALSE}$
 Delete *New_WE* file if it exists – crash happened before updating *MetadataFile*
 $\wedge New_WE' = [New_WE \text{ EXCEPT } !.exist = \text{FALSE}]$
 $\wedge TornWrite' = \text{FALSE}$
 $\wedge \text{UNCHANGED } \langle E_LowLSN, MaxNum, E_HighLSN \rangle$

TruncateHead (TH):
 Remove old data from the log.
 ASSUMPTIONS:
 1. There is only 1 *TH* call at a time - but because of 2 phases, there can be multiple *TH* in *phase2*.
 2. 1 *TH* and 1 append can happen concurrently.
 3. No *TT* when *TH* starts.
 We broke truncate head in 2 phases to simulate a crash in between 2 stages.
 Other states like appends can happen in between 2 phases.
Phase1 : Update metadata file first.
TruncateHeadP1 \triangleq
 $\wedge PrevState \notin \{\text{"crash"}, \text{"close"}\}$
truncate_head waits for *new_WE* workflow to finish.
 Todo: This is possibly bad as even starting the *TruncateHead* is waiting.
 It is not very bad because *New_WE* workflow should finish fast and is also rare.
 $\wedge E_NWEIP = \text{FALSE}$
 $\wedge E_TTIP = \text{FALSE}$ No truncate tail in progress.
 $\wedge E_LowLSN < E_HighLSN$
 $\wedge PrevState' = \text{"truncate_head_p1"}$
 $\wedge E_LowLSN' = E_LowLSN + 1$
WE is never removed from *MetadataFile* in case of *TH*
 as we need at least 1 file in logger at all time.
 $\wedge \text{LET } newREs \triangleq \text{LET } nonTruncatedRE(re) \triangleq re.end > E_LowLSN'$
 $\quad \text{IN } SelectSeq(REs, nonTruncatedRE)$
 $\quad \text{IN } MetadataFile' = [MetadataFile \text{ EXCEPT } !.headLSN = E_LowLSN',$
 $\quad \quad \quad !.lastTailLSN = E_HighLSN,$
 $\quad \quad \quad !.fileIds = GetFileIds(Append(newREs, WE))]$
 $\wedge E_THIP' = \text{TRUE}$
 $\wedge \text{UNCHANGED } \langle E_HighLSN, MaxNum, REs, WE, TornWrite, E_TTIP, E_NWEIP, New_WE \rangle$

Delete/Zero out RE files in 2nd phase of *TruncateHead*
TruncateHeadP2 \triangleq

$\wedge \vee PrevState = \text{"truncate_head_p1"}$
 $\vee \wedge E_THIP = \text{TRUE}$
 $\wedge PrevState \neq \text{"crash"}$ After crash, only recovery runs
 $\wedge PrevState' = \text{"truncate_head"}$
 $\wedge REs' = \text{LET } nonTruncatedRE(re) \triangleq re.end > E_LowLSN$
 $\quad \text{IN } SelectSeq(REs, nonTruncatedRE)$
 $\wedge E_THIP' = \text{FALSE}$
 $\wedge \text{UNCHANGED } \langle E_LowLSN, E_HighLSN, MaxNum, WE, MetadataFile, TornWrite, E_TTIP, E_NWEIP \rangle$

TruncateTailP1 :

Remove data from tail of the log.

ASSUMPTIONS:

1. No *TruncateTail* (*TT*) while append is called.
2. No *Append/TruncateHead* (*TH*) while *TT* is going on.

Phase1 : Update metadata file first.

We broke truncate tail in 2 phases to simulate a crash in between 2 stages.

Update metadata file first:

If we crash after updating metadata file, we can truncate tail of *WE* on recovery.

Other valid states like appends can't run between 2 phases.

$TruncateTailP1 \triangleq$
 $\wedge PrevState \notin \{\text{"crash"}, \text{"close"}\}$
 $\wedge E_NWEIP = \text{FALSE}$
 $\wedge E_THIP = \text{FALSE}$
 $\wedge E_TTIP = \text{FALSE}$ Only one truncate tail allowed at a time.
 $\wedge E_LowLSN < E_HighLSN$
 $\wedge MetadataFile.lastTailVersion < MaxNum - 1$ Restrict *MC* to finite states.
 $\wedge PrevState' = \text{"truncate_tail_p1"}$
 $\wedge E_HighLSN' = E_HighLSN - 1$
 $\wedge E_TTIP' = \text{TRUE}$
 $\wedge \text{In } TruncateTail - \text{update } tailLsn, \text{version, } fileIds \text{ in } MetadataFile$
 $\wedge \text{LET } validExtents \triangleq \text{IF } WE.start < WE.end \text{ } WE \text{ has data}$
 $\quad \text{THEN } Append(REs, [WE \text{ EXCEPT } !.end = WE.end - 1])$
 $\quad \text{ELSE LET } lastRE \triangleq REs[Len(REs)]$
 $\quad \quad \text{IN } Append(SubSeq(REs, 1, Len(REs) - 1), [lastRE \text{ EXCEPT } !.end = lastRE.end])$
 $\quad \text{IN } MetadataFile' = [MetadataFile \text{ EXCEPT } !.lastTailLSN = E_HighLSN',$
 $\quad \quad !.lastTailVersion = MetadataFile.lastTailVersion + 1,$
 $\quad \quad !.fileIds = GetFileIds(validExtents)]$
 $\wedge \text{UNCHANGED } \langle E_LowLSN, WE, REs, MaxNum, TornWrite, E_THIP, E_NWEIP, New_WE \rangle$

Now, actual delete file/zero *WE* file's tail in Phase 2.

TruncateTailP2 \triangleq

$\wedge PrevState = \text{"truncate_tail_p1"}$
 $\wedge PrevState' = \text{"truncate_tail"}$
 $\wedge \text{IF } WE.start < WE.end \text{ } WE \text{ has data}$


```

    THEN  $\wedge WE' = [WE \text{ EXCEPT } !.end = WE.end - 1]$ 
       $\wedge REs' = REs$ 
    ELSE  $\wedge WE' = \text{LET } lastRE \triangleq REs[Len(REs)]$ 
      IN  $[lastRE \text{ EXCEPT } !.end = lastRE.end - 1]$ 
       $\wedge REs' = SubSeq(REs, 1, Len(REs) - 1)$ 
 $\wedge E\_TTIP' = \text{FALSE}$ 
 $\wedge \text{UNCHANGED } \langle E\_LowLSN, E\_HighLSN, MaxNum, MetadataFile, TornWrite, E\_THIP, E\_NWEIP, N$ 

Next  $\triangleq$ 
 $\vee AppendToFile$ 
 $\vee WriteExtentFullNewWE$ 
 $\vee NewWriteExtentAddToMetadataFile$ 
 $\vee CrashWhileAppend$ 
 $\vee CrashNoDataLoss$ 
 $\vee Close$ 
 $\vee Recovery$ 
 $\vee TruncateHeadP1$ 
 $\vee TruncateHeadP2$ 
 $\vee TruncateTailP1$ 
 $\vee TruncateTailP2$ 
  Not modelling Data Loss.
  I am not sure, if we should just fail to open if we find we lost data
  so that we build from new replica.
 $\vee CrashDataLost$ 

Invariants:

Invariant 1: NoDataLoss
All read only extents have non missing LSN
 $REs[1].start < REs[1].end \triangleq REs[2].start < REs[2].end \triangleq REs[3].start < \dots$ 
write extent has latest data
 $E\_HighLSN \triangleq WE.end \geq WE.start \triangleq REs[last].end$ 

[dangling_extent] [lowLSN, highLSN - valid range] [dangling_extent]
OrderedExtent(ex1, ex2, highLSN)  $\triangleq$ 
 $\wedge ex1.start < ex1.end$ 
 $\wedge ex1.end = ex2.start$ 
 $\wedge ex1.end \leq highLSN$ 

ValidReadOnlyExtents  $\triangleq$ 
 $\wedge \forall i \in 1 \dots Len(REs) - 1 : \wedge OrderedExtent(REs[i], REs[i + 1], E\_HighLSN)$ 
 $\wedge REs[i].end < E\_HighLSN$ 
 $\wedge \text{IF } Len(REs) > 0$ 
  THEN OrderedExtent(REs[Len(REs)], WE, E\_HighLSN)
  ELSE 1 = 1

ValidWriteExtent  $\triangleq$ 

```

$$\begin{aligned}
& \wedge WE.start \leq WE.end \\
& \wedge WE.end = E_HighLSN \\
MetadataExtentsCoverDataRange & \triangleq \\
\text{LET } allFiles & \triangleq GetMetadataFiles \\
firstFile & \triangleq allFiles[1] \\
lastFile & \triangleq allFiles[Len(allFiles)] \\
\text{IN } & \wedge firstFile.start \leq E_LowLSN \\
& \wedge lastFile.end \geq E_HighLSN \\
NoDataLoss & \triangleq \\
& \text{Not valid state during crash or } truncate_tail_phase1 \\
& \vee PrevState \in \{\text{"crash"}, \text{"truncate_tail_p1"}\} \\
& \vee E_TTIP = \text{TRUE } TruncateTail \text{ in progress} \\
& \vee \wedge ValidReadOnlyExtents \\
& \wedge ValidWriteExtent \\
& \wedge MetadataExtentsCoverDataRange \\
& \wedge E_LowLSN \leq E_HighLSN \\
\text{Invariant 2: No dangling files on disk} & \\
\text{No file/extent present on disk which are not required.} & \\
NotDanglingExtent(ex, lowLSN, highLSN) & \triangleq \\
& \neg(\vee ex.start \geq highLSN \\
& \vee ex.end \leq lowLSN) \\
NoDanglingExtents & \triangleq \\
& \vee PrevState = \text{"crash"} \\
& TH/TT \text{ is in progress - so some files are dangling.} \\
& \vee E_THIP = \text{TRUE} \\
& \vee E_TTIP = \text{TRUE} \\
& \vee \wedge \forall i \in 1 \dots Len(REs) : NotDanglingExtent(REs[i], E_LowLSN, E_HighLSN) \\
& \wedge \vee WE.start = WE.end \text{ } WE \text{ is empty} \\
& \vee E_LowLSN = E_HighLSN \text{ } \text{There is no data in log} \\
& \text{If there is some data, } WE \text{ should be valid} \\
& \vee NotDanglingExtent(WE, E_LowLSN, E_HighLSN) \\
\text{Invariant 3 : Correctness of } MetadataFile: & \\
1. FileIds \text{ should be in increasing order} & \\
2. HeadLSN \text{ should be same as Expected } E_LowLSN & \\
IsFileIdPresent(fileIds, id) & \triangleq \\
\text{LET } SameId(fid) & \triangleq fid = id \\
\text{IN } Len(SelectSeq(fileIds, SameId)) & = 1 \\
AllMetadataFilesPresentOnDisk & \triangleq \\
\text{LET } allFiles & \triangleq Append(REs, WE) \\
allFileIds & \triangleq [i \in 1 \dots Len(allFiles) \mapsto allFiles[i].id] \\
\text{IN } & \wedge \forall i \in 1 \dots Len(MetadataFile.fileIds) :
\end{aligned}$$

```

    IsFileIdPresent(allFileIds, MetadataFile.fileIds[i])
  ∧ IF New_WE.exist
    if New_WE is present, it should not be in RE, WE and mentioned in MetadataFile
    i.e New_WE is transient file
  THEN  ∧ ¬IsFileIdPresent(allFileIds, New_WE.id)
        ∧ ¬IsFileIdPresent(MetadataFile.fileIds, New_WE.id)
    ELSE 1 = 1

```

```

MetadataFileCorrect ≜
  ∧ No missing file - files in increasing order
  ∀ i ∈ 1 .. Len(MetadataFile.fileIds) - 1 :
    MetadataFile.fileIds[i] < MetadataFile.fileIds[i + 1]
  ∧ MetadataFile.headLSN = E_LowLSN
  ∧ AllMetadataFilesPresentOnDisk even during crash
  ∧ IF MetadataFile.cleanShutdown
    THEN MetadataFile.lastTailLSN = E_HighLSN
    Todo: What should still be correct in clean shutdown case ?
  ELSE 1 = 1

```

```

Invariant 4: Valid version number
CorrectVersionNumber ≜
  IF MetadataFile.lastTailLSN < E_HighLSN some write finished after last TT
  THEN WE.version = MetadataFile.lastTailVersion
  Multiple TT can happen one after another increasing version no. Thanks TLA+
  ELSE WE.version ≤ MetadataFile.lastTailVersion

```

```

If we have clean shutdown state - except after close, we are in bad state.
CleanShutdownOnlyAfterClose ≜
  ¬( ∧ PrevState ≠ "close"
    ∧ MetadataFile.cleanShutdown = TRUE
  )

```

```

Change below value to see different steps taken for particular test run.
LSNSteps ≜
  E_HighLSN < MaxNum

```

Spec Ends

MAK: Do you check that the algorithm makes progress? Remember, a system that does nothing doesn't violate any safety prop.

Todo:

Spec after this is WIP and not used.

Invariants that should fail - Signifies that we have handled these cases.

TruncateTail is not called on empty WE for truncating data upto REs

Todo: This is not failing - This case is not handled.

Using E_TTIP is not correct as that means that TT has already begin.

TruncateTailCalledOnEmptyWE \triangleq
 $\neg(\wedge E_TTIP = \text{TRUE}$
 $\wedge WE.start = WE.end$
 $)$

TruncateHeadCalledOnEmptyWE \triangleq
 $\neg(\wedge E_THIP = \text{TRUE}$
 $\wedge WE.start = WE.end$
 $)$

Crash: Not modelling case of data lost.

we lost all data after some *LSN*

Todo: need to model - data lost in between *E_LowLSN* and *E_HighLSN*

In that case, we will Fail replica in real world
and rebuild from another source.

MaxOf2(a, b) \triangleq
IF $a < b$
THEN b
ELSE a

CrashDataLost \triangleq
 $\wedge PrevState' = \text{"crash"}$
 $\wedge E_HighLSN' = \text{IF } E_HighLSN > (MaxNum \div 2)$
 THEN $MaxOf2(E_LowLSN, MaxNum \div 2)$
 ELSE IF $E_HighLSN > 3$
 THEN $MaxOf2(E_LowLSN, 3)$
 ELSE $MaxOf2(E_LowLSN, 1)$
 $\wedge MetadataFile' = [MetadataFile \text{ EXCEPT } !.cleanShutdown = \text{FALSE}]$
 $\wedge \text{UNCHANGED } \langle E_LowLSN, MaxNum, REs, WE, TornWrite \rangle$

\ * Modification History
\ * Last modified Tue Nov 17 09:44:38 PST 2020 by markus
\ * Last modified Tue Nov 17 09:22:58 PST 2020 by markus
\ * Last modified Mon Nov 16 16:32:52 PST 2020 by asnegi
\ * Created Wed Oct 28 17:55:29 PDT 2020 by asnegi