

---

MODULE *ZeusOwnership*

---

EXTENDS *ZeusOwnershipMeta*

This Module specifies the full slow-path of the *Zeus* ownership protocol as appears in the according paper of *Eurosys'21* without faults. It model checks its properties in the face of concurrent conflicting requests of changing ownerships and emulated reliable commits.

Faults are added on top with the *ZeusOwnershipFaults.tla* spec

---

WARNING: We need to make sure that requester *REQs* are executed at most once; this requires: an *APP* node to be sticky to its *LB* driver for a Key (unless failure) and send *REQ* msgs via a *FIFO REQ* channel so that driver does not re-issues *REQs* that have been already completed in the past!

We emulate executing only once via *committedRTS* (*committedREQs* is used to check INVARIANTS)

$$\begin{aligned} \text{commit\_REQ}(o\_ts, r\_ts) &\triangleq \\ &\wedge \text{committedRTS}' = \text{committedRTS} \cup \{r\_ts\} \\ &\wedge \text{committedREQs}' = \text{committedREQs} \cup \{o\_ts\} \end{aligned}$$

$$\begin{aligned} \text{upd\_t\_meta}(n, \text{version}, \text{state}, t\_acks) &\triangleq \\ &\wedge tState' = [tState \text{ EXCEPT } ![n] = \text{state}] \\ &\wedge tVersion' = [tVersion \text{ EXCEPT } ![n] = \text{version}] \\ &\wedge tRcvACKs' = [tRcvACKs \text{ EXCEPT } ![n] = t\_acks] \end{aligned}$$

$$\begin{aligned} \text{upd\_r\_meta}(n, \text{ver}, tb, id, \text{type}) &\triangleq \\ &\wedge rID' = [rID \text{ EXCEPT } ![n] = id] \\ &\wedge rEID' = [rEID \text{ EXCEPT } ![n] = mEID] \quad \text{always update to latest } mEID \\ &\wedge rType' = [rType \text{ EXCEPT } ![n] = \text{type}] \\ &\wedge rTS' = [rTS \text{ EXCEPT } ![n].\text{ver} = \text{ver}, ![n].tb = tb] \end{aligned}$$

to update the epoch *id* of last message issue

$$\text{upd\_rEID}(n) \triangleq \text{upd\_r\_meta}(n, rTS[n].\text{ver}, rTS[n].tb, rID[n], rType[n])$$

$$\begin{aligned} \text{upd\_o\_meta}(n, \text{ver}, tb, \text{state}, \text{driver}, \text{vec}, \text{ACKs}) &\triangleq \\ &\wedge oVector' = [oVector \text{ EXCEPT } ![n] = \text{vec}] \\ &\wedge oRcvACKs' = [oRcvACKs \text{ EXCEPT } ![n] = \text{ACKs}] \\ &\wedge oState' = [oState \text{ EXCEPT } ![n] = \text{state}] \\ &\wedge oDriver' = [oDriver \text{ EXCEPT } ![n] = \text{driver}] \\ &\wedge oTS' = [oTS \text{ EXCEPT } ![n].\text{ver} = \text{ver}, ![n].tb = tb] \end{aligned}$$

$$\text{upd\_o\_meta\_driver}(n, \text{ver}, tb) \triangleq \text{upd\_o\_meta}(n, \text{ver}, tb, \text{"drive"}, n, oVector[n], \{\})$$

$$\begin{aligned} \text{upd\_o\_meta\_add\_ack}(n, \text{sender}) &\triangleq \\ &\text{upd\_o\_meta}(n, oTS[n].\text{ver}, oTS[n].tb, oState[n], oDriver[n], oVector[n], oRcvACKs[n] \cup \{\text{sender}\}) \end{aligned}$$

$$\begin{aligned} \text{upd\_o\_meta\_apply\_val}(n, m) &\triangleq \\ &\wedge \text{IF } rTS[n].tb \notin m\text{AliveNodes} \\ &\quad \text{THEN } \text{upd\_o\_meta}(n, oTS[n].\text{ver}, oTS[n].tb, \text{"valid"}, 0, \text{post\_oVec}(n, 0, oVector[n]), \{\}) \\ &\quad \text{ELSE } \text{upd\_o\_meta}(n, oTS[n].\text{ver}, oTS[n].tb, \text{"valid"}, 0, \text{post\_oVec}(n, rTS[n].tb, oVector[n]), \{\}) \end{aligned}$$

$$\begin{aligned} \text{upd\_o\_meta\_apply\_val\_n\_reset\_o\_state}(n) &\triangleq \\ \text{upd\_o\_meta}(n, 0, 0, \text{"valid"}, 0, [\text{readers} \mapsto \{\}, \text{owner} \mapsto 0], \{\}) \end{aligned}$$


---

REQUESTER Helper operators

$$\begin{aligned} \text{choose\_req}(n) &\triangleq \\ \text{LET } \text{choice} &\triangleq \text{CHOOSE } x \in \{0, 1\} : \text{TRUEIN} \\ \text{IF } \text{is\_reader}(n) & \\ \text{THEN } \wedge \text{IF } \text{choice} = 0 & \\ \text{THEN "change-owner"} & \\ \text{ELSE "remove-reader"} & \\ \text{ELSE } \wedge \text{IF } \text{choice} = 0 & \\ \text{THEN "add-owner"} & \\ \text{ELSE "add-reader"} & \\ \text{max\_committed\_ver}(S, n) &\triangleq \text{IF } \forall i \in S : i.tb \neq n \text{ THEN } [ver \mapsto 0, tb \mapsto 0] \\ &\quad \text{ELSE CHOOSE } i \in S : \wedge i.tb = n \\ &\quad \quad \wedge \forall j \in S : \vee j.tb \neq n \\ &\quad \quad \quad \vee j.ver \leq i.ver \\ \text{next\_rTS\_ver}(n) &\triangleq \text{max\_committed\_ver}(\text{committedRTS}, n).ver + 1 \\ \text{upd\_rs\_meta\_n\_send\_req}(n, r\_type) &\triangleq \\ \wedge \text{upd\_r\_meta}(n, \text{next\_rTS\_ver}(n), n, 0, r\_type) & \\ \wedge \text{upd\_o\_meta}(n, 0, 0, \text{"request"}, 0, [\text{readers} \mapsto \{\}, \text{owner} \mapsto 0], \{\}) & \\ \wedge \text{o\_send\_req}([ver \mapsto \text{next\_rTS\_ver}(n), tb \mapsto n], 0, r\_type) & \end{aligned}$$


---

REQUESTER ACTIONS

$$\begin{aligned} O\text{RequesterREQ}(n) &\triangleq \text{Requester issues a REQ} \\ &\wedge \text{is\_valid\_requester}(n) \\ &\wedge \text{is\_reader}(n) \\ &\wedge \text{next\_rTS\_ver}(n) \leq O\_MAX\_VERSION \quad \text{bound execution -> Bound this in reachable states} \\ &\wedge \text{upd\_rs\_meta\_n\_send\_req}(n, \text{"change-owner"}) \quad \text{to limit the state space only choose change ownership} \\ &\wedge \text{upd\_rs\_meta\_n\_send\_req}(n, \text{choose\_req}(n)) \\ &\wedge \text{unchanged\_mtc} \\ \text{Requester receives NACK and replays REQ w/ higher rID} \\ O\text{RequesterNACK}(n) &\triangleq \\ &\wedge \text{is\_in\_progress\_requester}(n) \\ &\wedge rID[n] < O\_MAX\_VERSION \quad \text{TODO: may Bound rID to number of APP\_NODES instead} \\ &\wedge \exists m \in oMsgs : o\_rcv\_nack(m, n) \\ &\wedge \text{upd\_r\_meta}(n, rTS[n].ver, n, rID[n] + 1, rType[n]) \\ &\wedge \text{o\_send\_req}([ver \mapsto rTS[n].ver, tb \mapsto n], rID[n] + 1, rType[n]) \\ &\wedge \text{unchanged\_mtco} \end{aligned}$$

$ORequesterRESP(n) \triangleq$  Requester receives a *RESP* and sends a *VAL* to arbiters  
 $\exists m \in oMsgs :$   
 $\wedge o\_rcv\_resp(m, n)$   
 $\wedge is\_in\_progress\_requester(n)$   
 $\wedge commit\_REQ(m.oTS, rTS[n])$   
 $\wedge upd\_t\_meta(n, m.tVersion, "valid", tRcvACKs[n])$  todo this is optional  
 $\wedge upd\_o\_meta(n, m.oTS.ver, m.oTS.tb, "valid", 0, post\_oVec(n, n, m.oVector), \{\})$   
 $\wedge o\_send\_val(m.oTS)$   
 $\wedge unchanged\_mtr$

$ORequesterActions \triangleq$   
 $\exists n \in APP\_LIVE\_NODES :$   
 $\vee ORequesterREQ(n)$   
 $\vee ORequesterNACK(n)$   
 $\vee ORequesterRESP(n)$

---

**DRIVER ACTIONS**

$ODriverINV(n, m) \triangleq$   
 $\wedge o\_rcv\_req(m)$   
 $\wedge oState[n] = "valid"$   
 $\wedge oTS[n].ver < O\_MAX\_VERSION$  bound execution - > Bound this in reachable states  
 $\wedge upd\_t\_meta(n, 0, tState[n], tRcvACKs[n])$   
 $\wedge upd\_r\_meta(n, m.rTS.ver, m.rTS.tb, m.rID, m.rType)$   
 $\wedge upd\_o\_meta\_driver(n, oTS[n].ver + 1, n)$   
 $\wedge o\_send\_inv(n, n, [ver \mapsto oTS[n].ver + 1, tb \mapsto n], oVector[n], m.rTS, m.rID, m.rType)$   
 $\wedge unchanged\_mc$

$ODriverNACK(n, m) \triangleq$   
 $\wedge o\_rcv\_req(m)$   
 $\wedge rTS[n] \neq m.rTS$   
 $\wedge oState[n] \neq "valid"$   
 $\wedge msg\_not\_exists(o\_rcv\_nack, m.rTS.tb)$  NACK does not exist (bound state space)  
 $\wedge o\_send\_nack(m.rTS, m.rID)$   
 $\wedge unchanged\_mtrco$

$ODriverACK(n, m) \triangleq$   
 $\wedge o\_rcv\_ack(m, n)$   
 $\wedge upd\_o\_meta\_add\_ack(n, m.sender)$   
 $\wedge IF\ m.tVersion \neq 0$   
 $\quad THEN\ upd\_t\_meta(n, m.tVersion, tState[n], tRcvACKs[n])$   
 $\quad ELSE\ unchanged\_t$   
 $\wedge unchanged\_Mmrc$

$ODriverRESP(n) \triangleq$   
 $\wedge oState[n] = "drive"$

$\wedge has\_rcvd\_all\_ACKs(n)$   
 $\wedge requester\_is\_alive(n)$   
 $\wedge msg\_not\_exists(o\_rcv\_resp, rTS[n].tb)$  *RESP does not exist (bound state space)*  
 $\wedge o\_send\_resp(rTS[n], oTS[n], post\_oVec(n, rTS[n].tb, oVector[n]), tVersion[n])$   
 $\wedge unchanged\_mtrco$

$ODriverActions \triangleq$   
 $\exists n \in LB\_LIVE\_NODES :$   
 $\vee ODriverRESP(n)$   
 $\vee \exists m \in oMsgs :$   
 $\vee ODriverINV(n, m)$   
 $\vee ODriverNACK(n, m)$   
 $\vee ODriverACK(n, m)$

---

**LB ARBITER ACTIONS**

$inv\_to\_be\_applied(n, m) \triangleq$   
 $\vee o\_rcv\_inv\_greater\_ts(m, n)$   
 $\vee (o\_rcv\_inv\_equal\_ts(m, n) \wedge oState[n] = \text{"invalid"} \wedge m.epochID > rEID[n])$

$check\_n\_apply\_inv(n, m) \triangleq$   
 $\wedge inv\_to\_be\_applied(n, m)$   
 $\wedge upd\_r\_meta(n, m.rTS.ver, m.rTS.tb, m.rID, m.rType)$   
 $\wedge upd\_o\_meta(n, m.oTS.ver, m.oTS.tb, \text{"invalid"}, m.driver, m.oVector, \{\})$

*We do not model lost messages thus arbiter need not respond w/ INV when ts is smaller*

$OLBArbiterINV(n, m) \triangleq$   
 $\wedge check\_n\_apply\_inv(n, m)$   
 $\wedge \vee oState[n] \neq \text{"drive"}$   
 $\vee o\_send\_nack(rTS[n], rID[n])$   
 $\wedge o\_send\_ack(n, m.oTS, 0)$   
 $\wedge unchanged\_mtc$

$OLBArbiterVAL(n, m) \triangleq$   
 $\wedge o\_rcv\_val(m, n)$   
 $\wedge upd\_o\_meta\_apply\_val(n, m)$   
 $\wedge unchanged\_Mmtrc$

$OLBArbiterActions \triangleq$   
 $\exists n \in LB\_LIVE\_NODES : \exists m \in oMsgs :$   
 $\vee OLBArbiterINV(n, m)$   
 $\vee OLBArbiterVAL(n, m)$

---

**(O)wner or (R)eader ARBITER ACTIONS**

*reader doesn't apply an INV but always responds with an ACK*

(and data if non-sharing  $rType$  and in  $tValid$  state)

$$\begin{aligned}
ORArbiterINV(n, m) &\triangleq \\
&\wedge is\_reader(n) \\
&\wedge tState[n] = \text{"valid"} \\
&\wedge o\_rcv\_inv(m, n) \\
&\wedge o\_send\_ack(n, m.oTS, tVersion[n]) \\
&\wedge unchanged\_mtrco \\
\\
OOArbiterINV(n, m) &\triangleq \\
&\wedge is\_owner(n) \\
&\wedge m.type = \text{"S\_INV"} \\
&\wedge m.oVector.owner = n \quad \text{otherwise owner lost a VAL} \rightarrow SFMOArbiterINVLostVAL \\
&\wedge tState[n] = \text{"valid"} \\
&\wedge check\_n\_apply\_inv(n, m) \\
&\wedge o\_send\_ack(n, m.oTS, tVersion[n]) \\
&\wedge unchanged\_mtc \\
\\
OOArbiterVAL(n, m) &\triangleq \\
&\wedge o\_rcv\_val(m, n) \\
&\wedge \text{IF } oVector[n].owner = n \\
&\quad \text{THEN } \wedge upd\_o\_meta\_apply\_val(n, m) \\
&\quad \text{ELSE } \wedge upd\_o\_meta\_apply\_val\_n\_reset\_o\_state(n) \\
&\wedge unchanged\_Mmtrc \\
\\
OAPPArbiterActions &\triangleq \\
&\exists n \in APP\_LIVE\_NODES : \exists m \in oMsgs : \\
&\quad \vee ORArbiterINV(n, m) \\
&\quad \vee OOArbiterINV(n, m) \\
&\quad \vee OOArbiterVAL(n, m)
\end{aligned}$$


---

Owner actions emulating tx updates

$$\begin{aligned}
TOwnerINV(n) &\triangleq \\
&\wedge upd\_t\_meta(n, tVersion[n] + 1, \text{"write"}, \{\}) \\
&\wedge t\_send(n, \text{"T\_INV"}, tVersion[n] + 1) \\
&\wedge unchanged\_mrco \\
\\
TOwnerACK(n) &\triangleq \\
&\exists m \in oMsgs : \\
&\quad \wedge t\_rcv\_ack(m, n) \\
&\quad \wedge upd\_t\_meta(n, tVersion[n], tState[n], tRcvACKs[n] \cup \{m.sender\}) \\
&\quad \wedge unchanged\_Mmrco \\
\\
TOwnerVAL(n) &\triangleq \\
&\wedge oVector[n].readers \subseteq tRcvACKs[n] \quad \text{has received all acks from readers} \\
&\wedge upd\_t\_meta(n, tVersion[n], \text{"valid"}, \{\})
\end{aligned}$$

$$\wedge t\_send(n, \text{"T\_VAL"}, tVersion[n])$$

$$\wedge unchanged\_mrco$$

Reader actions emulating tx updates

$$TReaderINV(n) \triangleq$$

$$\exists m \in oMsgs :$$

$$\wedge t\_rcv\_inv(m, n)$$

$$\wedge m.tVersion > tVersion[n]$$

$$\wedge upd\_t\_meta(n, m.tVersion, \text{"invalid"}, \{\})$$

$$\wedge t\_send(n, \text{"T\_ACK"}, m.tVersion)$$

$$\wedge unchanged\_mrco$$

$$TReaderVAL(n) \triangleq$$

$$\exists m \in oMsgs :$$

$$\wedge t\_rcv\_val(m, n)$$

$$\wedge m.tVersion = tVersion[n]$$

$$\wedge upd\_t\_meta(n, tVersion[n], \text{"valid"}, \{\})$$

$$\wedge unchanged\_Mmrco$$

$$TOwnerReaderActions \triangleq$$

$$\exists n \in APP\_LIVE\_NODES :$$

$$\vee \wedge is\_valid\_owner(n)$$

$$\wedge \vee TOwnerINV(n)$$

$$\vee TOwnerACK(n)$$

$$\vee TOwnerVAL(n)$$

$$\vee \wedge is\_reader(n)$$

$$\wedge \vee TReaderINV(n)$$

$$\vee TReaderVAL(n)$$


---

Modeling *Sharding* protocol (Requester and Arbiter actions)

$$ONext \triangleq$$

$$\vee OInit\_min\_owner\_rest\_readers$$

$$\vee ORequesterActions$$

$$\vee ODriverActions$$

$$\vee OLBArbiterActions$$

$$\vee OAPPArbiterActions$$

$$\vee TOwnerReaderActions$$

The complete definition of the algorithm

$$Spec \triangleq OInit \wedge \Box [ONext]_{vars}$$

$$Invariants \triangleq \wedge (\Box OTypeOK)$$

$$\wedge (\Box CONSISTENT\_DATA) \wedge (\Box ONLY\_ONE\_CONC\_REQ\_COMMITTS)$$

$$\wedge (\Box AT\_MOST\_ONE\_OWNER) \wedge (\Box OWNER\_LATEST\_DATA)$$

$$\wedge (\Box CONSISTENT\_SHARERS) \wedge (\Box CONSISTENT\_OVECTORS)$$

THEOREM  $Spec \Rightarrow Invariants$

$$LSpec \triangleq Spec \wedge WF\_vars(ONext)$$

$$LIVENESS \triangleq \exists i \in LB\_NODES: \Box\Diamond(oState[i] = \text{"valid"} \wedge oTS[i].ver > 3)$$

THEOREM  $LSpec \Rightarrow LIVENESS$