─────────── MODULE *WSAtomicTransaction* ───────────

EXTENDS *TLC*

This is a preliminary version of a formal specification of the Web Services Atomic Transaction protocol, described in the document by Cabrera et al. We are specifying the safety property of the protocol (what is allowed to happen), not its liveness property (what must eventually happen).

The protocol involves an initiator, a transaction coordinator ($TC$), and a set of participants. The $TC$ exchanges messages with the participants. For convenience, we assume that the initiator and $TC$ are actually executed on the same processor, so they can be considered to be a single process.

The protocol allows messages to be lost, duplicated, or received out of order. A process is therefore free to resend a message at any time. An implementation will resend a message if it times out without receiving a reply. Such resending is not described explicitly in the specification. Instead, we represent the communication infrastructure by a set *msgs* of all messages that have ever been sent. Since resending a sent message does not change the specification's state, it is allowed by our specification.

An action that, in an implementation, would be enabled by the receipt of certain messages is here enabled by the existence of those messages in *msgs*. Loss of a message is represented by simply not executing that action, even though it is enabled. This works because we are specifying only safety properties, so there is no requirement that an enabled action is ever executed.

Since messages are never removed from *msgs*, receipt of the same message twice is allowed. This can happen in an implementation either because the communication infrastructure delivers a duplicate copy, or because the sender mistakenly believed the original copy had been lost and resent the message. In most cases, the specification says that such duplicate copies are ignored because the response has already been sent, so *msgs* already contains the response. However, in an implementation, receipt of a duplicate message may indicate that the sender resent the message because it never received the response. Hence, an implementation would resend the response. The specification sometimes asserts that such responses are in *msgs*, indicating that they would probably be resent in an implementation.

The CONSTANT and VARIABLES sections define constant parameters and variables.

CONSTANT *Participant*   The set of all participants.

VARIABLES  *iState*,     The state of the initiator.
           *tcData*,     The data maintained by the coordinator.
           *pData*,      $pData[p]$ is the data maintained by participate $p$.
           *msgs*        The set of all sent messages.

*Message* $\triangleq$

The set of all possible messages. A message sent from the $TC$ to a participant has a *dest* field indicating its destination. A message from a participant to the $TC$ has a *src* field indicating its sender. A participant's "Register" message also has a *reg* field indicating if it's registering as volatile or durable.

$\quad$ [*type* : { "RegisterResponse" , "Prepare" , "Commit" , "Rollback" },
$\quad$ *dest* : *Participant*]
$\cup$ [*type* : { "Prepared" , "ReadOnly" , "Committed" , "Aborted" },
$\quad$ *src* : *Participant*]
$\cup$ [*type* : { "Register" },
$\quad$ *reg* : { "volatile" , "durable" },
$\quad$ *src* : *Participant*]

1

$TypeOK \triangleq$

$\land iState \in \{$ "active", "committed", "aborted", "completing" $\}$

$\land tcData \in$
$\quad [st \;: \{$ "active", "preparingVolatile", "preparingDurable",
$\qquad\qquad$ "aborting", "committing" $\}$,
$\quad reg : [Participant \rightarrow \{$ "unregistered", "volatile", "durable",
$\qquad\qquad\qquad\qquad$ "prepared", "readOnly", "committed" $\}]]$
$\quad \cup \;\; [st \;: \{$ "ended" $\}$,
$\qquad\quad res : \{$ "committed", "aborted" $\}]$

$\land pData \in [Participant \rightarrow \quad [st : \{$ "unregistered", "prepared" $\}]$
$\qquad\qquad\qquad\qquad\quad \cup$
$\qquad\qquad\qquad\qquad [st \;: \{$ "registering", "active", "preparing" $\}$,
$\qquad\qquad\qquad\qquad\; reg : \{$ "volatile", "durable" $\}]$
$\qquad\qquad\qquad \cup \;\; [st \;: \{$ "ended" $\}$,
$\qquad\qquad\qquad\qquad\; res : \{$ "?", "committed", "aborted" $\}]]$

$\land\ msgs \subseteq Message$

$Consistency\ \triangleq$

$\land\ (iState = \text{``committed''})$
$\qquad \Rightarrow\ \lor\ \land\ tcData.st = \text{``ended''}$
$\qquad\qquad\qquad \land\ tcData.res = \text{``committed''}$
$\qquad\qquad\qquad \land\ \forall\, p \in Participant :$
$\qquad\qquad\qquad\qquad \lor\ pData[p].st = \text{``unregistered''}$
$\qquad\qquad\qquad\qquad \lor\ \land\ pData[p].st\ =\ \text{``ended''}$
$\qquad\qquad\qquad\qquad\qquad \land\ pData[p].res \in\ \{\,\text{``?''},\ \text{``committed''}\,\}$
$\qquad\qquad\ \lor\ \land\ tcData.st = \text{``committing''}$
$\qquad\qquad\qquad \land\ \forall\, p \in Participant :$
$\qquad\qquad\qquad\qquad \lor\ pData[p].st \in \{\,\text{``unregistered''},\ \text{``prepared''}\,\}$
$\qquad\qquad\qquad\qquad \lor\ \land\ pData[p].st\ =\ \text{``ended''}$
$\qquad\qquad\qquad\qquad\qquad \land\ pData[p].res \in\ \{\,\text{``?''},\ \text{``committed''}\,\}$

$\land\ \forall\, p \in Participant :$
$\qquad \land\ pData[p].st\ =\ \text{``ended''}$
$\qquad \land\ pData[p].res = \text{``committed''}$
$\qquad \Rightarrow\ \land\ iState = \text{``committed''}$
$\qquad\qquad \land\ \lor\ \land\ tcData.st\ =\ \text{``ended''}$
$\qquad\qquad\qquad\quad \land\ tcData.res = \text{``committed''}$
$\qquad\qquad\qquad\quad \land\ iState\ =\ \text{``committed''}$
$\qquad\qquad\quad \lor\ tcData.st = \text{``committing''}$
$\qquad\qquad \land\ \forall\, pp \in Participant :$
$\qquad\qquad\qquad \lor\ pData[pp].st \in \{\,\text{``unregistered''},\ \text{``prepared''}\,\}$
$\qquad\qquad\qquad \lor\ \land\ pData[pp].st\ =\ \text{``ended''}$
$\qquad\qquad\qquad\qquad \land\ pData[pp].res \in\ \{\,\text{``?''},\ \text{``committed''}\,\}$

$Init \triangleq$

$\land iState\ = $ "active"
$\land tcData = [st\ \ \mapsto$ "active",
$\qquad\qquad\quad reg \mapsto [p \in Participant \mapsto$ "unregistered"$]]$
$\land pData\ = [p \in Participant \mapsto [st \mapsto$ "unregistered"$]]$
$\land msgs\ \ = \{\}$

---

THE ACTIONS

The next-state action is the disjunction of the four actions $TCInternal$, $TCRcvMsg$, $PInternal$, and $PRcvMsg$. The major part of the specification consists of the definitions of these four actions, which follow.

$TCInternal \triangleq$

This action describes the actions of the initiator and the internal actions of the $TC$–that is, the actions of the initiator prompted by timeouts or by a spontaneous action of the initiator. It also describes actions enabled by the $TC$ having received enough messages. (Those actions could be described as occurring when the the $TC$ receives the last message needed to enable it, but it's more convenient to let it be done as a separate internal action.)

$\land\ \lor$  The initiator decides to commit the transaction. It sets its state to "completing". At the same time, the $TC$ sets its state to "preparingVolatile" and sends "Prepare" messages to every participant that has registered as volatile.

$\qquad\ \land iState =$ "active"
$\qquad\ \land \forall\, p \in Participant :$
$\qquad\qquad (pData[p].st =$ "registering"$) \Rightarrow$
$\qquad\qquad\quad \land\ \ pData[p].reg =$ "durable"
$\qquad\qquad\quad \land\ \ \exists\, pp \in Participant :$
$\qquad\qquad\qquad\qquad \land pData[pp].st \in \{$ "active", "preparing" $\}$
$\qquad\qquad\qquad\qquad \land pData[pp].reg =$ "volatile"

The only participants that may be registering are durable ones that are being installed by a volatile participant that is not yet prepared. It is up to the application to ensure that this condition is met.

$\qquad\ \land iState'\ =$ "completing"
$\qquad\ \land tcData' = [tcData \text{ EXCEPT } !.st =$ "preparingVolatile"$]$
$\qquad\ \land msgs' = msgs \cup [type : \{$ "Prepare" $\},$
$\qquad\qquad\qquad\qquad dest : \{p \in Participant :$
$\qquad\qquad\qquad\qquad\qquad\qquad tcData.reg[p] =$ "volatile" $\}]$

$\quad\ \lor$  Either the initiator decides to abort the transaction while in its "active" state, or else the $TC$ decides to abort it while in a "preparing" state–presumably because of some timeout. The initiator's state is set to "aborted", the $TC$ state is set to "aborting", and "Rollback" messages are sent to every registered participant from which the $TC$ did not already receive a "ReadOnly" message.

$\qquad\ \land\ \lor iState =$ "active"
$\qquad\qquad \lor tcData.st \in \{$ "preparingVolatile", "preparingDurable" $\}$
$\qquad\ \land iState'\ =$ "aborted"
$\qquad\ \land tcData' = [tcData \text{ EXCEPT } !.st =$ "aborting"$]$

4

$$\land msgs' = msgs \cup [type : \{\text{"Rollback"}\},$$
$$dest : \{p \in Participant :$$
$$tcData.reg[p] \notin \{\text{"unregistered"},$$
$$\text{"readOnly"}\}\}]$$

$\lor$ The $TC$ ends the volatile prepare and begins the durable prepare. It does this when it has received "Prepared" or "ReadOnly" messages from every participant that registered as volatile, and it sends "Prepare" message to every participant that registered as durable.

$\quad \land tcData.st = \text{"preparingVolatile"}$
$\quad \land \forall p \in Participant : tcData.reg[p] \neq \text{"volatile"}$
$\quad \land tcData' = [tcData \text{ EXCEPT } !.st = \text{"preparingDurable"}]$
$\quad \land msgs' = msgs \cup [type : \{\text{"Prepare"}\},$
$\qquad\qquad\qquad dest : \{p \in Participant :$
$\qquad\qquad\qquad\qquad tcData.reg[p] = \text{"durable"}\}]$
$\quad \land \text{UNCHANGED } iState$

$\lor$ The $TC$ finishes the durable prepare and decides to commit the transaction. It can do this when it has received a "Prepared" or "ReadOnly" message from every durable particpant. It sets its state to "committing", notifies the initiator that the transaction has committed, and sends "Commit" messages to every participant that replied with a "Prepared" message (instead of a "ReadOnly" message).

$\quad \land tcData.st = \text{"preparingDurable"}$
$\quad \land \forall p \in Participant : tcData.reg[p] \neq \text{"durable"}$
$\quad \land tcData' = [tcData \text{ EXCEPT } !.st = \text{"committing"}]$
$\quad \land msgs' = msgs \cup [type : \{\text{"Commit"}\},$
$\qquad\qquad\qquad dest : \{p \in Participant :$
$\qquad\qquad\qquad\qquad tcData.reg[p] = \text{"prepared"}\}]$
$\quad \land iState' = \text{"committed"}$

$\lor$ The action by which the $TC$ forgets the transaction, entering the "ended" state. It can do this if it is in the "aborting" state, or if it is in the "committing" state and has received a "ReadOnly" or "Committed" message from every registered participant.

$\quad \land \lor tcData.st = \text{"aborting"}$
$\qquad \lor \land tcData.st = \text{"committing"}$
$\qquad\quad \land \forall p \in Participant :$
$\qquad\qquad tcData.reg[p] \in \{\text{"unregistered"}, \text{"readOnly"}, \text{"committed"}\}$
$\quad \land tcData' = [st \mapsto \text{"ended"},$
$\qquad\qquad\qquad res \mapsto \text{IF } tcData.st = \text{"aborting"} \text{ THEN } \text{"aborted"}$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{ELSE } \text{"committed"}]$
$\quad \land \text{UNCHANGED } \langle iState, msgs \rangle$

$\land \text{UNCHANGED } pData$

The participants' states are left unchanged.

The following action definition uses the TLA+ construct

$\text{CASE } B1 \to P1 \;\square\; \ldots \;\square\; Bn \to Pn$

5

This construct is used here (and in most places) when the state predicates $Bi$, which are called the *guards* , are mutually disjoint–that is, no two of them can be true in the same state. When the guards are mutually disjoint, the CASE expression equals $Pi$ if guard $Bi$ is true. If none of the guards are true, then the value of the expression is undefined. If $TLC$ ever evaluates such an undefined expression, it reports an error. Thus, this CASE statement is equivalent to the formula

$$(B1 \wedge P1) \vee \ldots \vee (Bn \wedge Pn)$$

except that when none of the guards are true, the formula equals FALSE while the value of the CASE statement is undefined.

In this following action, the guards of each CASE formula describe all the possible cases in which the $TC$ can receive a particular message.

$TCRcvMsg \triangleq$

  The action in which the $TC$ receives a message from a participant.

  $\exists\, m \in msgs :$

    $m$ is the message being received.

    LET $Reply(tp) \triangleq msgs' = msgs \cup \{[type \mapsto tp,\ dest \mapsto m.src]\}$

        Locally defines $Reply(tp)$ to be the action of sending a message of type $tp$ to the sender of message $m$.

       $HaveSent(tp) \triangleq \exists\, mm \in msgs : (mm.type = tp) \wedge (mm.dest = m.src)$

        Locally defines $HaveSent(tp)$ to be true iff the $TC$ has sent a message of type $tp$ to the sender of $m$.

    IN   $\vee$  $m$ is a "Register" message.

          $\wedge\ m.type =$ "Register"

          $\wedge$ CASE  The normal case, in which the $TC$ state is either "active", or else this is a durable participant registering while the $TC$ is performing the volatile prepare.

               $\vee\ tcData.st =$ "active"

               $\vee\ \wedge\ tcData.st =$ "preparingVolatile"

                  $\wedge\ m.reg =$ "durable"

            $\rightarrow$  The $TC$ sends a "RegisterResponse" reply to the sender, and sets the appropriate $tcData.reg$ component (the one corresponding to the sender) to the contents of the $reg$ field of $m$.

               $\wedge\ Reply(\text{"RegisterResponse"})$

               $\wedge\ tcData' = [tcData \text{ EXCEPT } !.reg[m.src] = m.reg]$

           $\square$  If the $TC$ is in a "preparing" or "committing" state or has forgotten the transaction, then $m$ is a duplicate message to which the $TC$ has already responded and so is ignored.

             $\wedge\ \vee\ \wedge\ tcData.st =$ "preparingVolatile"

                  $\wedge\ m.reg =$ "volatile"

               $\vee\ tcData.st \in \{\text{"preparingDurable"},\ \text{"committing"}\}$

             $\wedge\ HaveSent(\text{"RegisterResponse"})$

          $\rightarrow$

           UNCHANGED $\langle tcData,\ msgs \rangle$

6

□ If the $TC$ is in the "aborting" state, then if the sender is not already registered, then the decision to abort was made before the sender could register, in which case a "Rollback" message is sent. (Is this correct?) Otherwise, this is a duplicate message to which the $TC$ has already responded, and it has already sent a "Rollback" message to the sender (unless the participant responded "ReadOnly" to a "Prepare" message).

$\wedge\ tcData.st =$ "aborting"
$\wedge\ \vee\ tcData.reg[m.src] \in \{$ "unregistered", "readOnly" $\}$
$\quad\ \vee\ \wedge\ tcData.reg[m.src] \in \{$ "volatile", "durable", "prepared" $\}$
$\qquad\ \wedge\ HaveSent($ "Rollback" $)$
$\rightarrow$
$\quad \wedge\ \text{IF}\ tcData.reg[m.src] =$ "unregistered"
$\qquad\quad \text{THEN}\ Reply($ "Rollback" $)$
$\qquad\quad \text{ELSE}\ \ \text{UNCHANGED}\ msgs$
$\quad \wedge\ \text{UNCHANGED}\ tcData$

□ If the $TC$ is in the "ended" state, then either the transaction has been aborted before the sender had a chance to register, or else the "Register" message is old and, if it was committed, then the sender has already forgotten the transaction and hence will ignore any message it receives. Therefore, it's safe to send a "Rollback" message.

$\wedge\ tcData.st =$ "ended"
$\wedge\ \vee\ tcData.res =$ "aborted"
$\quad\ \vee\ pData[m.src].st =$ "ended"
$\rightarrow$
$\quad \wedge\ Reply($ "Rollback" $)$
$\quad \wedge\ \text{UNCHANGED}\ tcData$

$\wedge\ \text{UNCHANGED}\ \langle iState,\ pData\rangle$

The initiator's state and the participants' data are not changed.

$\vee$ $m$ is a "Prepared" message.

$\wedge\ m.type =$ "Prepared"
$\wedge\ \text{CASE}$ The normal case, in which the $TC$ has sent a "Prepare" message and is waiting for the sender's reply.

$\quad\ \vee\ \wedge\ tcData.st =$ "preparingVolatile"
$\qquad\ \wedge\ tcData.reg[m.src] =$ "volatile"
$\quad\ \vee\ \wedge\ tcData.st =$ "preparingDurable"
$\qquad\ \wedge\ tcData.reg[m.src] =$ "durable"
$\rightarrow$
$\quad \wedge\ tcData' = [tcData\ \text{EXCEPT}\ !.reg[m.src] =$ "prepared" $]$
$\quad \wedge\ \text{UNCHANGED}\ msgs$

□ The $TC$ has forgotten the transaction.

$tcData.st =$ "ended"
$\rightarrow$

The $TC$ sends a "Rollback" reply to the sender. The transaction could either have aborted or committed. However, if the transaction has committed, then the sender will know that it has and will ignore the "Rollback" message.

$\land\ Reply(\text{"Rollback"})$
$\land\ \text{UNCHANGED}\ tcData$

□    If the $TC$ is either in a "preparing" state, or in the "aborting" or "committing" state, then it has already received and acted on a copy of $m$, so it does nothing.

$\lor\ \land\ tcData.st \in \{\,\text{"preparingVolatile"},\ \text{"preparingDurable"}\,\}$
     $\land\ tcData.reg[m.src] = \text{"prepared"}$
$\lor\ \land\ tcData.st = \text{"aborting"}$
     $\land\ HaveSent(\text{"Rollback"})$
$\lor\ \land\ tcData.st = \text{"committing"}$
     $\land\ HaveSent(\text{"Commit"})$
$\rightarrow$
     $\text{UNCHANGED}\ \langle tcData,\ msgs\rangle$

$\land\ \text{UNCHANGED}\ \langle iState,\ pData\rangle$

The initiator's state is not changed.

$\lor$    $m$ is a "ReadOnly" message.

$\land\ m.type = \text{"ReadOnly"}$
$\land\ \text{CASE}$    The normal case, in which the $TC$ has sent a "Prepare" message and is waiting for the sender's reply.

     $\lor\ \land\ tcData.st = \text{"preparingVolatile"}$
       $\land\ tcData.reg[m.src] = \text{"volatile"}$
     $\lor\ \land\ tcData.st = \text{"preparingDurable"}$
       $\land\ tcData.reg[m.src] = \text{"durable"}$
$\rightarrow$

The $TC$ sets its $tcData.reg$ component corresponding to the sender to "readOnly".

$\land\ tcData' = [tcData\ \text{EXCEPT}\ !.reg[m.src] = \text{"readOnly"}]$
$\land\ \text{UNCHANGED}\ msgs$

□    If the $TC$ has forgotten the transaction, then either $m$ is a duplicate message, or else it was decided to abort the transaction before the $TC$ received the response to the "Prepare" message it sent to the sender of $m$. In either case, the message is ignored.

$tcData.st = \text{"ended"}$
$\rightarrow$
$\text{UNCHANGED}\ \langle tcData,\ msgs\rangle$

□    In the following cases, $m$ is a duplicate of a message that the $TC$ has already received and it is ignored.

8

$$\lor \land tcData.st \in \{\,\text{``preparingVolatile''},\ \text{``preparingDurable''}\,\}$$
$$\quad\ \land tcData.reg[m.src] = \text{``readOnly''}$$
$$\lor \land tcData.st = \text{``aborting''}$$
$$\quad\ \land \lor tcData.reg[m.src] = \text{``readOnly''}$$
$$\qquad\ \lor \land tcData.reg[m.src] \in \{\,\text{``volatile''},\ \text{``durable''}\,\}$$
$$\qquad\qquad \land HaveSent(\text{``Rollback''})$$
$$\lor tcData.st = \text{``committing''}$$
$$\to$$
$$\quad \text{UNCHANGED } \langle tcData,\ msgs \rangle$$

$$\land \text{UNCHANGED } \langle iState,\ pData \rangle$$

> The initiator's state and the participants' data are not changed.

$$\lor \quad$$ $m$ is an "Aborted" message.

$$\land m.type = \text{``Aborted''}$$

$\land$ CASE   The normal case, in which the $TC$ receives the message when it is "active" or in a "preparing" state and the sender has not replied to a "Prepare" message. In this case, the transaction is aborted, the $TC$ sends "Rollback" messages to all registered participants from which it has not already received a "ReadOnly" message, and the initiator is notified that the transaction has been aborted.

$$\land tcData.st \in \{\,\text{``active''},\ \text{``preparingVolatile''},\ \text{``preparingDurable''}\,\}$$
$$\land tcData.reg[m.src] \in \{\,\text{``unregistered''},\ \text{``volatile''},\ \text{``durable''}\,\}$$
$$\to$$
$$\quad \land iState' = \text{``aborted''}$$
$$\quad \land tcData' = [tcData \text{ EXCEPT } !.st = \text{``aborting''}]$$
$$\quad \land msgs' = msgs \cup$$
$$\qquad\qquad [type : \{\,\text{``Rollback''}\,\},$$
$$\qquad\qquad\ \ dest : \{p \in Participant :$$
$$\qquad\qquad\qquad\quad tcData.reg[p] \notin$$
$$\qquad\qquad\qquad\qquad \{\,\text{``unregistered''},\ \text{``readOnly''}\,\}\}]$$

> 11 *Dec* 2003, changed "ReadOnly" $\to$ "readOnly"
>
> typo found by *Colin Campbell*

$$\land \text{UNCHANGED } pData$$

☐   If the $TC$ is already in the "aborting" state or it has forgotten an aborted transaction, then the $TC$ has already sent a "Rollback" message to the sender (perhaps because $m$ is a duplicate of a message the $TC$ already received). If the $TC$ has forgotten a committed transaction, then this "Aborted" message was sent because the sender received an obsolete "Prepare" message after it had forgotten the transaction. In either case, $m$ is ignored.

$$\land \lor tcData.st = \text{``aborting''}$$
$$\quad \lor \land tcData.st = \text{``ended''}$$
$$\qquad\ \land \lor \land tcData.res = \text{``aborted''}$$
$$\qquad\qquad\quad \land HaveSent(\text{``Rollback''})$$
$$\qquad\qquad \lor \land tcData.res = \text{``committed''}$$

$$\wedge\, pData[m.src].st\ =\ \text{``ended''}$$
$$\wedge\, pData[m.src].res =\ \text{``committed''}$$
$$\rightarrow$$

UNCHANGED $\langle tcData,\, pData,\, iState,\, msgs \rangle$

$\square$

$\wedge \vee tcData.st \in \{\,\text{``committing''},\ \text{``ended''}\,\}$
$\quad \vee \wedge tcData.st \in\ \{\,\text{``active''},\ \text{``preparingVolatile''},$
$\qquad\qquad\qquad\qquad\qquad \text{``preparingDurable''}\,\}$
$\qquad \wedge tcData.reg[m.src] \in \{\,\text{``prepared''},\ \text{``readOnly''},$
$\qquad\qquad\qquad\qquad\qquad\qquad \text{``committed''}\,\}$

$\rightarrow$

UNCHANGED $\langle tcData,\, pData,\, iState,\, msgs \rangle$

$\vee$

$\wedge\, m.type =\ \text{``Committed''}$

$\wedge$ CASE

$\qquad tcData.st =\ \text{``committing''}$

$\rightarrow$

$\qquad tcData' = [tcData\ \text{EXCEPT}\ !.reg[m.src] =\ \text{``committed''}]$

$\square$

$\qquad \wedge\, tcData.st\ =\ \text{``ended''}$
$\qquad \wedge\, tcData.res =\ \text{``committed''}$

$\rightarrow$

$\qquad$ UNCHANGED $tcData$

$\wedge$ UNCHANGED $\langle iState,\, pData,\, msgs \rangle$

$PInternal\ \triangleq$

$\exists\, p \in Participant :$

LET $SendMsg(tp)\ \triangleq\ msgs' =\ msgs \cup \{[type \mapsto tp,\ src \mapsto p]\}$
$\quad SendRegisterMsg(rg)\ \triangleq\ msgs' =\ msgs \cup \{[type \mapsto \text{``Register''},$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad src\ \mapsto p,$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad reg\ \mapsto rg]\}$

Locally defined action expressions. $SendMsg(tp)$ sends a message of type $tp$ from participant $p$ to the $TC$. $SendRegisterMsg(rg)$ sends a $Register\ rg$ message, where $rg$ is either "durable" or "volatile".

IN $\quad\lor\quad$ $p$ registers as a volatile participant. It can do this only if it is *unregistered* and the initiator is in the "active" state.

$\quad\land pData[p].st =$ "unregistered"
$\quad\land iState =$ "active"
$\quad\land pData' = [pData \text{ EXCEPT } ![p] = [st \mapsto$ "registering",
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad reg \mapsto$ "volatile"$]]$
$\quad\land SendRegisterMsg($ "volatile" $)$
$\quad\land$ UNCHANGED $\langle iState,\ tcData \rangle$

$\lor\quad$ $p$ registers as a durable participant. It can do this only if it is *unregistered* and, if either the initiator is in the "active" state, or there is some volatile participant that is willing to wait for $p$ to register before preparing. Since we don't model "willingness to wait", we allow the participant to register as long as there is some volatile participant that can wait for it.

$\quad\land pData[p].st =$ "unregistered"
$\quad\land \lor iState =$ "active"
$\qquad\lor \exists pp \in Participant :$
$\qquad\qquad\land pData[pp].st \in \{$ "active", "preparing" $\}$
$\qquad\qquad\land pData[pp].reg =$ "volatile"
$\quad\land pData' = [pData \text{ EXCEPT } ![p] = [st \mapsto$ "registering",
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad reg \mapsto$ "durable"$]]$
$\quad\land SendRegisterMsg($ "durable" $)$
$\quad\land$ UNCHANGED $\langle iState,\ tcData \rangle$

$\lor\quad$ $p$ spontaneously aborts and forgets about the transaction. We do not allow it to abort in the "registering" state. If we allowed this, then we could wind up with a situation in which a participant aborted before it registered, and the transaction committed anyway. In practice, there will have to be some way for a participant to give up when it hasn't received a *RegisterResponse* message. However, to do this, it must learn from the initiator or a volatile participant that the transaction aborted so it can forget about it. Since we are not modeling this kind of inter-participant communication, we do not model this procedure.

$\quad\land pData[p].st \in \{$ "active", "preparing" $\}$

$\quad\land pData' = [pData \text{ EXCEPT } ![p] = [st \mapsto$ "ended",
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad res \mapsto$ "aborted"$]]$
$\quad\land SendMsg($ "Aborted" $)$
$\quad\land$ UNCHANGED $\langle iState,\ tcData \rangle$

$\lor\quad$ $p$ either prepares or becomes read-only. If $p$ is volatile, then it cannot do this if there is a durable participant that is in the "registering" state, and there is no other volatile participant to wait for it to register.

$\quad\land pData[p].st =$ "preparing"
$\quad\land \lor pData[p].reg =$ "durable"
$\qquad\lor \neg\exists dp \in Participant : \land pData[dp].st =$ "registering"
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \land pData[dp].reg =$ "durable"

$\lor \, \exists \, vp \in Participant \setminus \{p\} : \land pData[vp].st \ = \text{``active''}$
$\qquad\qquad\qquad\qquad\qquad\quad \land pData[vp].reg = \text{``volatile''}$

$\land \, \lor \, \land pData' = [pData \text{ EXCEPT } ![p] = [st \ \mapsto \text{``prepared''}]]$
$\qquad \land SendMsg(\text{``Prepared''})$
$\quad \lor \, \land pData' = [pData \text{ EXCEPT } ![p] = [st \ \mapsto \text{``ended''},$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad res \mapsto \text{``?''}]]$
$\qquad \land SendMsg(\text{``ReadOnly''})$
$\land \text{ UNCHANGED } \langle iState, tcData \rangle$

$PRcvMsg \triangleq$

The action in which a participant receives a message from the $TC$.

$\exists \, m \in msgs :$

$m$ is the message being received.

LET $Reply(tp) \ \triangleq \ msgs' = \ msgs \cup \{[type \mapsto tp, \ src \mapsto m.dest]\}$

Locally defines $Reply(tp)$ to be the action of sending a message of type $tp$ from the sender of message $m$.

$HaveSent(tp) \ \triangleq \ \exists \, mm \in msgs : (mm.type = tp) \land (mm.src = m.dest)$

Locally defines $HaveSent(tp)$ to be true iff the participant receiving $m$ has sent a message of type $tp$ to the $TC$.

IN $\quad \land \lor \quad m$ is a "RegisterResponse" message
$\qquad\qquad\quad \land m.type = \text{``RegisterResponse''}$
$\qquad\qquad\quad \land$ CASE $\quad$ The normal case.

$\qquad\qquad\qquad\qquad pData[m.dest].st = \text{``registering''}$
$\qquad\qquad\qquad \rightarrow$
$\qquad\qquad\qquad\qquad \land pData' = [pData \text{ EXCEPT } ![m.dest].st = \text{``active''}]$
$\qquad\qquad\qquad\qquad \land \text{ UNCHANGED } msgs$

$\qquad\qquad\quad \square \quad$ This is a duplicate of an already-received message, or else it is ignorable because another message to the participant arrived ahead of it.

$\qquad\qquad\qquad\qquad pData[m.dest].st \in$
$\qquad\qquad\qquad\qquad\qquad \{\text{``active''}, \text{``preparing''}, \text{``prepared''}, \text{``ended''}\}$
$\qquad\qquad\qquad \rightarrow$
$\qquad\qquad\qquad\qquad \text{UNCHANGED } \langle pData, msgs \rangle$

$\qquad\quad \lor \quad m$ is a "Prepare" message.
$\qquad\qquad \land m.type = \text{``Prepare''}$
$\qquad\qquad \land \quad$ This is either the normal case (the participant is in the "active" state), or else this "Prepare" message has arrived before the "RegisterResponse" message.

$\qquad\qquad\quad$ CASE $pData[m.dest].st \in \{\text{``registering''}, \text{``active''}\}$
$\qquad\qquad\qquad \rightarrow \ \land pData' = [pData \text{ EXCEPT } ![m.dest].st = \text{``preparing''}]$
$\qquad\qquad\qquad\qquad \land \text{ UNCHANGED } msgs$

$\qquad\qquad\quad \square \quad$ This is a duplicate of a message already received.

12

$$pData[m.dest].st \in \{\text{``preparing''}, \text{``prepared''}\}$$
$\rightarrow$

   UNCHANGED $\langle pData, msgs \rangle$

□   The transaction has been forgotten. The participant responds with an "Aborted" message–even though the transaction might have committed or be in the process of committing. In the latter case, the "Aborted" message will be ignored by the $TC$.

   $\land\ pData[m.dest].st = \text{``ended''}$
   $\land\ \lor\ \land\ pData[m.dest].res = \text{``committed''}$
   $\qquad\ \land\ HaveSent(\text{``Committed''})$
   $\quad\ \lor\ \land\ pData[m.dest].res = \text{``aborted''}$
   $\qquad\ \land\ \lor\ HaveSent(\text{``Aborted''})$
   $\qquad\qquad\ \lor\ \land\ tcData.st = \text{``ended''}$
   $\qquad\qquad\qquad\ \land\ tcData.res = \text{``committed''}$
   $\quad\ \lor\ \land\ pData[m.dest].res = \text{``?''}$
   $\qquad\ \land\ HaveSent(\text{``ReadOnly''})$
   $\rightarrow$

   $\land\ Reply(\text{``Aborted''})$
   $\land\ $ UNCHANGED $pData$

$\lor$   $m$ is a "Commit" message.

$\land\ m.type = \text{``Commit''}$

$\land$   The normal case.

   CASE $pData[m.dest].st = \text{``prepared''}$
   $\rightarrow$

   $\land\ pData' =$
   $\qquad [pData\ \text{EXCEPT}\ ![m.dest] = [st\ \mapsto\ \text{``ended''},$
   $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\ res \mapsto \text{``committed''}]]$
   $\land\ Reply(\text{``Committed''})$

   □   The transaction has ended, so this must be a message that was already received.

   $\land\ pData[m.dest].st = \text{``ended''}$
   $\land\ pData[m.dest].res \in \{\text{``?''}, \text{``committed''}\}$
   $\rightarrow$

   UNCHANGED $\langle pData, msgs \rangle$

$\lor$   $m$ is a "Rollback" message.

$\land\ m.type = \text{``Rollback''}$

$\land$   The participant can be in any registered state. If it hasn't ended, then this causes it to abort. It also sends an "Aborted" message to the $TC$. This message isn't needed, but it is apparently used as an acknowledgement.

   CASE $pData[m.dest].st \in$
   $\qquad \{\text{``registering''}, \text{``active''}, \text{``preparing''}, \text{``prepared''}\}$
   $\rightarrow$

$$\land\ pData' =$$
$$[pData \text{ EXCEPT } ![m.dest] = [st \ \mapsto \text{ "ended"},$$
$$res \mapsto \text{ "aborted"}]]$$
$$\land\ Reply(\text{"Aborted"})$$

☐ If the participant has already finished, then this is ignored. It is possible for this message to arrive even though the participant has ended by committing the transaction. In this case, the "Rollback" message was generated by a duplicate "Register" message arriving at the $TC$ after it had forgotten the transaction.

$$pData[m.dest].st\ =\ \text{"ended"}$$
$$\rightarrow$$
$$\text{UNCHANGED}\ \langle pData,\ msgs\rangle$$

$$\land\ \text{UNCHANGED}\ \langle iState,\ tcData\rangle$$

---

$Next \ \triangleq\ TCInternal \lor TCRcvMsg \lor PInternal \lor PRcvMsg$

The specification's next-state action.

---

$vars \ \triangleq\ \langle iState,\ tcData,\ pData,\ msgs\rangle$

The tuple of all variables.

---

$Spec \ \triangleq\ Init \land \Box[Next]_{vars}$

The complete spec of the two-phase $Commit$ protocol.

THEOREM $Spec \Rightarrow \Box(TypeOK \land Consistency)$

This theorem asserts that the predicates $TypeOK$ and $Consistency$ are invariants of the specification. $TLC$ checked this with 4 participants. It generated 10269919 states, 504306 of them were distinct. The longest non-repeating behavior had 45 states. It took $TLC$ about 4-1/4 minutes on a 2-processor, $2.4GHz$ machine.

Last modified on *Thu Dec* 11 13:27:41 *PT* 2003 by lamport