

An algorithm for replicating multi-threaded applications as done in replicated Popcorn. The application is made deterministic through the use of logical time. Any inter-thread synchronization operation must be protected by calls to *EnterSync* and *ExitSync*. Reads of the socket *API* are modeled by *EnterRead*. The scheduler processes (one per kernel) makes sure that the different copies of the application are consistent.

EXTENDS *Naturals, Sequences, Integers, Library*

CONSTANTS *Pid, MaxTime, Kernel, SchedulerPID, Requests*

ASSUME *SchedulerPID* \notin *Pid*

InitLogTime $\triangleq 1$

LogTime \triangleq *InitLogTime* .. *MaxTime* The set of logical time values

Processes are of the form $\langle k, pid \rangle$, where k is the kernel the process is running on.

P \triangleq *Kernel* \times *Pid*

Primary \triangleq CHOOSE $k \in$ *Kernel* : TRUE

Ker(*p*) \triangleq *p*[1]

PID(*p*) \triangleq *p*[2]

OnKernel(*kernel*) \triangleq {*kernel*} \times *Pid*

Logical time comparison, using *PIDs* to break ties.

Less(*p*, *tp*, *q*, *tq*) \triangleq
 $tp < tq \vee (tp = tq \wedge PID(p) \leq PID(q))$

The sequence of *TCP* packets that will be received. No duplicates allowed (therefore the set *TcpData* must be big enough) so that any misordering of the threads will lead to a different data read. For *TCPMultiStream*, each stream has different data.

StreamLength $\triangleq 3$

TcpData $\triangleq 1 \dots StreamLength * Requests$

TcpStream $\triangleq [i \in 1 \dots StreamLength \mapsto i]$

TcpMultiStream $\triangleq [r \in 1 \dots Requests \mapsto$
 $([i \in 1 \dots StreamLength \mapsto i + (StreamLength * (r - 1))])]$

ASSUME *NoDup*(*TcpStream*)

ASSUME *Len*(*TcpStream*) = *StreamLength*

Shifts a sequence by 1: $Shift(\langle 1, 2, 3 \rangle) = \langle 2, 3 \rangle$ and $Shift(\langle \rangle) = \langle \rangle$.

Shift(*s*) \triangleq

IF *Len*(*s*) > 1
 THEN $[i \in 1 \dots (Len(s) - 1) \mapsto s[i + 1]]$
 ELSE $\langle \rangle$

Shiftn(*s*) \triangleq

IF *Len*(*s*) > 1

```

THEN  $[i \in 1 \dots (Len(s) - 1) \mapsto s[i + 1]]$ 
ELSE  $\langle -1 \rangle$ 

```

The algorithm *ReadAppend* models a set of worker threads being scheduled by the deterministic scheduler and executing the following code.

```

Code of worker  $w$ : While(true){
   $x = read(socket)$ ;
  append( $\langle w, x \rangle$ , file);
}

```

Variables:

The variable *bumps* records all logical time bumps executed by the primary in order for the secondaries to do the same, *i.e.* the initial logical time, the new logical time, and the value read from the tcp buffer. $\langle t1, t2, d \rangle \in bumps[pid]$ means that the primary bumped process *pid* from logical time *t1* to *t2* and delivered the data *d*. Note that the scheduler set *bumps*[*pid*] to a value that depends on the logical time of the processes on all replicas, and this value is then immediately available to all replicas. A more detailed model would instead include a distributed implementation of the choice of the logical time to bump the process to.

reads[*p*] stores the last value read by *p* from the socket.

tcpBuff[*k*] represents the state of the tcp buffer on kernel *k*. Each time a process reads from the buffer, the buffer shrinks by 1.

Definitions:

Bumped(kernel) is the set of processes running on the kernel “kernel” which are waiting to execute a “bump” decided by the primary.

If $p \in Bumped(Ker(p))$ then *BumpedTo*(*p*) is the logical time to which *p* should be bumped to.

If $p \in WaitingSync(Ker(p))$ then *IsNextProc*(*kernel*, *p*) is true iff *p* is the process to be scheduled next, that is: (1) *p* has the lowest *ltime* among running and waiting-sync processes and (2) if *q* is on the same kernel and *q* is waiting for a read and the primary has already decided to which logical time *tq* to bump *q*, then *ltime*[*p*] must be less than *tq*.

BumpTo is the logical time to which to bump a process that needs bumping. It is some logical time greater than all the logical times reached by any process on any kernel.

```

--algorithm ReadAppend{
  variables
    status =  $[p \in P \mapsto \text{“running”}]$ ,
    ltime =  $[p \in P \mapsto InitLogTime]$ ,
    file =  $[k \in Kernel \mapsto \langle \rangle]$ ,
    bumps =  $[p \in Pid \mapsto \{\}]$ ,
    reads =  $[p \in P \mapsto -1]$ ,
    tcpBuff =  $[k \in Kernel \mapsto TcpMultiStream]$ ,
    Queue for accepted connections
    socketQueue =  $[k \in Kernel \mapsto \langle \rangle]$ ,
    Queue for unhandled connections
    requestQueue =  $[k \in Kernel \mapsto [r \in 1 \dots Requests \mapsto r]]$ ,
    The socket that is handled by a process
    handledSocket =  $[p \in P \mapsto -1]$ ,

```

```

connections = [k ∈ Kernel ↦ ⟨⟩]

define {
  Run(p)  $\triangleq$  status[p] = "running"
  Running(kernel)  $\triangleq$  {p ∈ OnKernel(kernel) : Run(p)}
  WaitingSync(kernel)  $\triangleq$ 
    {p ∈ OnKernel(kernel) : status[p] = "waiting sync"}
  WaitingRead(kernel)  $\triangleq$ 
    {p ∈ OnKernel(kernel) : status[p] = "waiting read"}
  Bumped(kernel)  $\triangleq$  {p ∈ OnKernel(kernel) :
    ∧ status[p] = "waiting read"
    ∧ ∃ t ∈ LogTime : ∃ d ∈ TcpData : ⟨ltime[p], t, d⟩ ∈ bumps[PID(p)]}
  BumpedTo(p)  $\triangleq$ 
    CHOOSE t ∈ LogTime : ∃ d ∈ TcpData : ⟨ltime[p], t, d⟩ ∈ bumps[PID(p)]
  BumpData(p)  $\triangleq$ 
    CHOOSE d ∈ TcpData : ∃ t ∈ LogTime : ⟨ltime[p], t, d⟩ ∈ bumps[PID(p)]
  IsNextProc(kernel, p)  $\triangleq$ 
    ∧ ∀ q ∈ Running(kernel) ∪ WaitingSync(kernel) :
      q ≠ p ⇒ Less(p, ltime[p], q, ltime[q])
    ∧ ∀ q ∈ Bumped(kernel) : Less(p, ltime[p], q, BumpedTo(q))
  BumpTo  $\triangleq$  CHOOSE i ∈ LogTime : ∀ p ∈ P : ltime[p] < i
}

macro EnterRead(p){
  status[p] := "waiting read" ;
}

macro EnterSync(p){
  status[p] := "waiting sync" ;
}

macro ExitSync(p){
  ltime[p] := ltime[p] + 1 ;
}

Processes consume a connection
process (worker ∈ P){
  ww1: while (TRUE){
    EnterSync(self) ;
  ww2: await Run(self) ;
  ww3: if (Len(requestQueue[Ker(self)]) > 0){
    handledSocket[self] := requestQueue[Ker(self)][1] ;
    requestQueue[Ker(self)] := Shift(requestQueue[Ker(self)]) ;
  } ;
  ww4: ExitSync(self) ;
  ww5: if (handledSocket[self] ≠ -1){
    ww9: while (Len(tcpBuff[Ker(self)][handledSocket[self]]) > 0){
      EnterRead(self) ;
    } ;
  } ;
}

```


$$\begin{aligned}
Running(kernel) &\triangleq \{p \in OnKernel(kernel) : Run(p)\} \\
WaitingSync(kernel) &\triangleq \\
&\quad \{p \in OnKernel(kernel) : status[p] = \text{"waiting sync"}\} \\
WaitingRead(kernel) &\triangleq \\
&\quad \{p \in OnKernel(kernel) : status[p] = \text{"waiting read"}\} \\
Bumped(kernel) &\triangleq \{p \in OnKernel(kernel) : \\
&\quad \wedge status[p] = \text{"waiting read"} \\
&\quad \wedge \exists t \in LogTime : \exists d \in TcpData : \langle ltime[p], t, d \rangle \in bumps[PID(p)]\} \\
BumpedTo(p) &\triangleq \\
&\quad CHOOSE t \in LogTime : \exists d \in TcpData : \langle ltime[p], t, d \rangle \in bumps[PID(p)] \\
BumpData(p) &\triangleq \\
&\quad CHOOSE d \in TcpData : \exists t \in LogTime : \langle ltime[p], t, d \rangle \in bumps[PID(p)] \\
IsNextProc(kernel, p) &\triangleq \\
&\quad \wedge \forall q \in Running(kernel) \cup WaitingSync(kernel) : \\
&\quad \quad q \neq p \Rightarrow Less(p, ltime[p], q, ltime[q]) \\
&\quad \wedge \forall q \in Bumped(kernel) : Less(p, ltime[p], q, BumpedTo(q)) \\
BumpTo &\triangleq CHOOSE i \in LogTime : \forall p \in P : ltime[p] < i
\end{aligned}$$

$$vars \triangleq \langle status, ltime, file, bumps, reads, tcpBuff, socketQueue, requestQueue, handledSocket, connections, pc \rangle$$

$$ProcSet \triangleq (P) \cup (\{ \langle k, SchedulerPID \rangle : k \in Kernel \})$$

$$\begin{aligned}
Init &\triangleq \text{Global variables} \\
&\quad \wedge status = [p \in P \mapsto \text{"running"}] \\
&\quad \wedge ltime = [p \in P \mapsto InitLogTime] \\
&\quad \wedge file = [k \in Kernel \mapsto \langle \rangle] \\
&\quad \wedge bumps = [p \in Pid \mapsto \{\}] \\
&\quad \wedge reads = [p \in P \mapsto -1] \\
&\quad \wedge tcpBuff = [k \in Kernel \mapsto TcpMultiStream] \\
&\quad \wedge socketQueue = [k \in Kernel \mapsto \langle \rangle] \\
&\quad \wedge requestQueue = [k \in Kernel \mapsto [r \in 1 .. Requests \mapsto r]] \\
&\quad \wedge handledSocket = [p \in P \mapsto -1] \\
&\quad \wedge connections = [k \in Kernel \mapsto \langle \rangle] \\
&\quad \wedge pc = [self \in ProcSet \mapsto \text{CASE } self \in P \rightarrow \text{"ww1"} \\
&\quad \quad \square \quad self \in \{ \langle k, SchedulerPID \rangle : k \in Kernel \} \rightarrow \text{"s1"}]
\end{aligned}$$

$$\begin{aligned}
ww1(self) &\triangleq \wedge pc[self] = \text{"ww1"} \\
&\quad \wedge status' = [status \text{ EXCEPT } ![self] = \text{"waiting sync"}] \\
&\quad \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"ww2"}] \\
&\quad \wedge \text{UNCHANGED } \langle ltime, file, bumps, reads, tcpBuff, socketQueue, requestQueue, handledSocket, connections \rangle
\end{aligned}$$

$$\begin{aligned}
ww2(self) &\triangleq \wedge pc[self] = \text{"ww2"} \\
&\quad \wedge Run(self) \\
&\quad \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"ww3"}]
\end{aligned}$$

\wedge UNCHANGED $\langle status, ltime, file, bumps, reads, tcpBuff,$
 $socketQueue, requestQueue, handledSocket,$
 $connections \rangle$

$ww3(self) \triangleq \wedge pc[self] = \text{"ww3"}$
 \wedge IF $Len(requestQueue[Ker(self)]) > 0$
 \quad THEN $\wedge handledSocket' = [handledSocket \text{ EXCEPT } ![self] = requestQueue[Ker(self)][1]]$
 $\quad \wedge requestQueue' = [requestQueue \text{ EXCEPT } ![Ker(self)] = Shift(requestQueue[Ker(self)])]$
 \quad ELSE \wedge TRUE
 $\quad \wedge$ UNCHANGED $\langle requestQueue, handledSocket \rangle$
 $\wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"ww4"}]$
 \wedge UNCHANGED $\langle status, ltime, file, bumps, reads, tcpBuff,$
 $socketQueue, connections \rangle$

$ww4(self) \triangleq \wedge pc[self] = \text{"ww4"}$
 $\wedge ltime' = [ltime \text{ EXCEPT } ![self] = ltime[self] + 1]$
 $\wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"ww5"}]$
 \wedge UNCHANGED $\langle status, file, bumps, reads, tcpBuff, socketQueue,$
 $requestQueue, handledSocket, connections \rangle$

$ww5(self) \triangleq \wedge pc[self] = \text{"ww5"}$
 \wedge IF $handledSocket[self] \neq -1$
 \quad THEN $\wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"ww9"}]$
 \quad ELSE $\wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"ww13"}]$
 \wedge UNCHANGED $\langle status, ltime, file, bumps, reads, tcpBuff,$
 $socketQueue, requestQueue, handledSocket,$
 $connections \rangle$

$ww9(self) \triangleq \wedge pc[self] = \text{"ww9"}$
 \wedge IF $Len(tcpBuff[Ker(self)][handledSocket[self]]) > 0$
 \quad THEN $\wedge status' = [status \text{ EXCEPT } ![self] = \text{"waiting read"}]$
 $\quad \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"ww10"}]$
 \quad ELSE $\wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"ww13"}]$
 $\quad \wedge$ UNCHANGED $status$
 \wedge UNCHANGED $\langle ltime, file, bumps, reads, tcpBuff, socketQueue,$
 $requestQueue, handledSocket, connections \rangle$

$ww10(self) \triangleq \wedge pc[self] = \text{"ww10"}$
 $\wedge Run(self)$
 $\wedge status' = [status \text{ EXCEPT } ![self] = \text{"waiting sync"}]$
 $\wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"ww11"}]$
 \wedge UNCHANGED $\langle ltime, file, bumps, reads, tcpBuff, socketQueue,$
 $requestQueue, handledSocket, connections \rangle$

$ww11(self) \triangleq \wedge pc[self] = \text{"ww11"}$
 $\wedge Run(self)$
 $\wedge file' = [file \text{ EXCEPT } ![Ker(self)] = Append(file[Ker(self)], \langle PID(self), reads[self] \rangle)]$

$$\begin{aligned}
& \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"ww12"}] \\
& \wedge \text{UNCHANGED } \langle status, ltime, bumps, reads, tcpBuff, \\
& \quad socketQueue, requestQueue, handledSocket, \\
& \quad connections \rangle \\
ww12(self) & \triangleq \wedge pc[self] = \text{"ww12"} \\
& \wedge ltime' = [ltime \text{ EXCEPT } ![self] = ltime[self] + 1] \\
& \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"ww9"}] \\
& \wedge \text{UNCHANGED } \langle status, file, bumps, reads, tcpBuff, socketQueue, \\
& \quad requestQueue, handledSocket, connections \rangle \\
ww13(self) & \triangleq \wedge pc[self] = \text{"ww13"} \\
& \wedge handledSocket' = [handledSocket \text{ EXCEPT } ![self] = -1] \\
& \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"ww1"}] \\
& \wedge \text{UNCHANGED } \langle status, ltime, file, bumps, reads, tcpBuff, \\
& \quad socketQueue, requestQueue, connections \rangle \\
worker(self) & \triangleq ww1(self) \vee ww2(self) \vee ww3(self) \vee ww4(self) \\
& \quad \vee ww5(self) \vee ww9(self) \vee ww10(self) \vee ww11(self) \\
& \quad \vee ww12(self) \vee ww13(self) \\
s1(self) & \triangleq \wedge pc[self] = \text{"s1"} \\
& \wedge \vee \wedge \exists p \in \{p \in \text{WaitingSync}(Ker(self)) : \\
& \quad \text{IsNextProc}(Ker(self), p)\} : \\
& \quad status' = [status \text{ EXCEPT } ![p] = \text{"running"}] \\
& \quad \wedge \text{UNCHANGED } \langle ltime, bumps, reads, tcpBuff \rangle \\
& \vee \wedge \exists p \in \text{WaitingRead}(Ker(self)) : \\
& \quad \wedge \text{IF } Ker(self) = \text{Primary} \\
& \quad \quad \text{THEN } \wedge bumps' = [bumps \text{ EXCEPT } ![PID(p)] = bumps[PID(p)] \cup \{\langle ltime[p], BumpTo \rangle\}] \\
& \quad \quad \wedge ltime' = [ltime \text{ EXCEPT } ![p] = BumpTo] \\
& \quad \quad \text{ELSE } \wedge p \in \text{Bumped}(Ker(self)) \wedge BumpData(p) = tcpBuff[Ker(self)][handledSocket[p]] \\
& \quad \quad \wedge ltime' = [ltime \text{ EXCEPT } ![p] = BumpedTo(p)] \\
& \quad \quad \wedge bumps' = bumps \\
& \quad \wedge reads' = [reads \text{ EXCEPT } ![p] = tcpBuff[Ker(self)][handledSocket[p]][1]] \\
& \quad \wedge tcpBuff' = [tcpBuff \text{ EXCEPT } ![Ker(self)][handledSocket[p]] = Shift(tcpBuff[Ker(self)] \\
& \quad \quad \wedge status' = [status \text{ EXCEPT } ![p] = \text{"running"}] \\
& \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"s1"}] \\
& \wedge \text{UNCHANGED } \langle file, socketQueue, requestQueue, handledSocket, \\
& \quad connections \rangle \\
scheduler(self) & \triangleq s1(self) \\
Next & \triangleq (\exists self \in P : worker(self)) \\
& \quad \vee (\exists self \in \{\langle k, SchedulerPID \rangle : k \in \text{Kernel}\} : scheduler(self)) \\
Spec & \triangleq Init \wedge \Box[Next]_{vars}
\end{aligned}$$

END TRANSLATION

$$\begin{aligned} RequestSynchronized &\triangleq \forall k1, k2 \in Kernel : \\ &\quad \vee Prefix(socketQueue[k1], socketQueue[k2]) \\ &\quad \vee Prefix(socketQueue[k2], socketQueue[k1]) \end{aligned}$$
$$\begin{aligned} FilesSynchronized &\triangleq \forall k1, k2 \in Kernel : \\ &\quad \vee Prefix(file[k1], file[k2]) \\ &\quad \vee Prefix(file[k2], file[k1]) \end{aligned}$$
$$\begin{aligned} \text{ConnectionSynchronized} &\triangleq \forall k1, k2 \in \text{Kernel} : \\ &\quad \vee \text{Prefix}(\text{connections}[k1], \text{connections}[k2]) \\ &\quad \vee \text{Prefix}(\text{connections}[k2], \text{connections}[k1]) \end{aligned}$$