

# Reliable Transport Protocol

Shankar

June 16, 2014

Overview

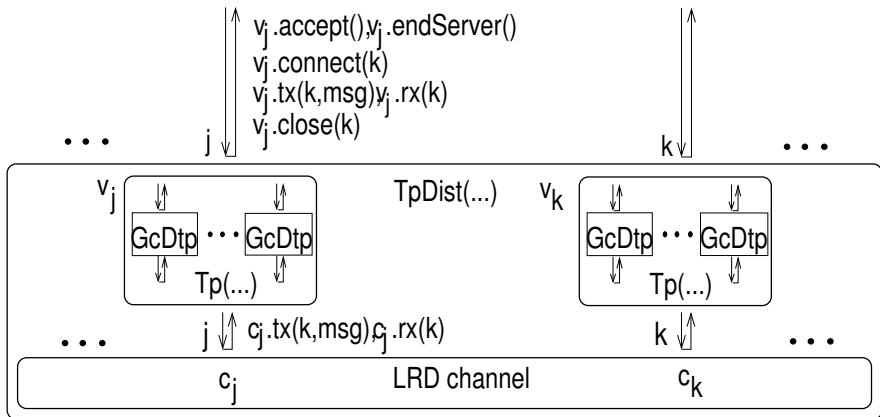
Graceful-closing data transfer protocol

Transport protocol description

Transport protocol program: unbounded endpoint numbers

Transport protocol program: cyclic endpoint numbers

Transport protocol with abort



- Implements reliable transport service using LRD channel
- $Tp$  system at each address

```
■ program TpDist ( ADDR )  
  // implements RelTransportService ( ADDR )  
  {cj} ← start LrdChannel(ADDR)  
  for j in ADDR  
    vj ← start Tp ( ADDR, j, cj )  
  return {vj}
```

- Connection establishment involves 3-way handshake
- Tp j maintains an **endpoint** for k only while interacting with k
  - **opening, open, open-and-closing, closed**
  - opening: **active** (if client) or **passive** (if server)
- Each endpoint gets a unique **endpoint number** when created
  - same role as TCP's initial sequence numbers
  - increasing but **need not be consecutive** // clock, counter
  - when open, maintains both local and remote endpoint numbers
- Endpoint's 4-tuple: [j,k,n,m] // addrs j, k; endpt numbers n, m

- Clean separation of connection establishment and data transfer
- Conn establishment provides dedicated **virtual** channel
  - by tagging msgs with endpoint's 4-tuple
- Can run any dtp (data transfer protocol) over this
  - for concreteness, use **graceful-closing dtp** (GcDtp)
- Tp system starts a dtp system when endpoint becomes open
  - relays msgs: dtp system  $\leftrightarrow$  user, LRD channel
- Connect req msg indicates whether sender is client or server
  - TCP does not do this, and hence has a flaw

## Overview

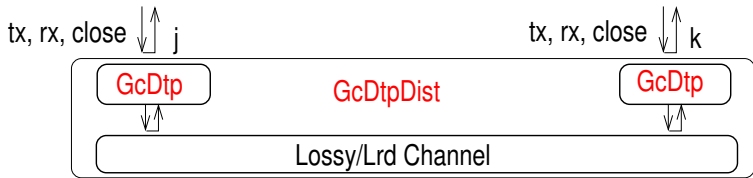
Graceful-closing data transfer protocol

Transport protocol description

Transport protocol program: unbounded endpoint numbers

Transport protocol program: cyclic endpoint numbers

Transport protocol with abort



## ■ Input fns

- tx(k, data)
- rx(): returns [0, data] or [-1] (if in-data done)
- close(): returns if data xfr done, remote closing

## ■ Msgs: DAT, ACK, FIN, FINACK

// FINs sent by close

## ■ Local thread fns:

- doTxDat(): sends DAT
- doRxDataAck(): rcvs; sends ACK, FINACK



- $A_1$  : if  $j.rx$  returns  $[0, db]$  then  $j.drXH \circ [db]$  prefix-of  $k.dtxh$
- $A_2$  : if  $j.rx$  returns  $[-1]$  then  $j.drXH = k.dtxh$ ,  $k$  closing/closed
- $A_3$  : if  $j$  is closed then no ongoing  $j.tx$  or  $j.rx$ ,  $j.drXH = k.dtxh$ , and  $k$  closing/closed
- $A_4$  : if  $j.tx$  is ongoing then  $j.tx$  returns
- $A_5$  : if  $j.rx$  ongoing,  $j.drXH \neq k.dtxh$ , then  $j.rx$  returns
- $A_6$  : if  $j.rx$  is ongoing,  $j.drXH = k.dtxh$ ,  $k$  closing/closed then  $j.rx$  returns
- $A_7$  : if  $j$  closing,  $k$  closing/closed then  $j$  becomes closed or  $j.rx$  not ongoing

- Terminate GcDtp system when it becomes closed:
  - fn close: insert endSystem() before the return
  - original: doRxDataAck() responds to FIN even when closed
  - now: tp system handles takes care of this
  
- Rename output calls: c.tx, c.rx → x.dTx, x.dRx
  - original: c is Lrd channel sid
  - now: c is tp sid // already has input fns tx, rx

Overview

Graceful-closing data transfer protocol

Transport protocol description

Transport protocol program: unbounded endpoint numbers

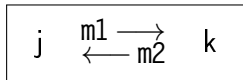
Transport protocol program: cyclic endpoint numbers

Transport protocol with abort

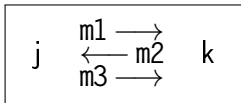
- Endpoint ( $ep$ ) [ $j, k, n, m$ ]
  - maintained at  $tp\ j$  while interacting with  $tp\ k$
  - $n$ :  $ep$ 's **local** number;  $\geq 0$ 
    - set from **increasing** clock/counter when  $ep$  created
  - $m$ :  $ep$ 's **remote** number;  $-1$  if unknown
    - rcvd from remote endpoint [ $k, j, m, n$ ]
- $Ep$  [ $j, k, n, m$ ]
  - **active opening** (**aop**):  $j.connect(k)$  ongoing
  - **passive opening** (**pop**):  $j.accept$  responding to  $k.connect(j)$
  - **open**: connected to [ $k, j, m, n$ ]; send/rcv data
  - **closing**:  $j.close(k)$  ongoing; rcv data // still open
  - **closed**: endpoint no longer exists

- Interaction between  $j$  and  $k$  is a succession of **handshakes**
  - each connect/close req initiates a handshake

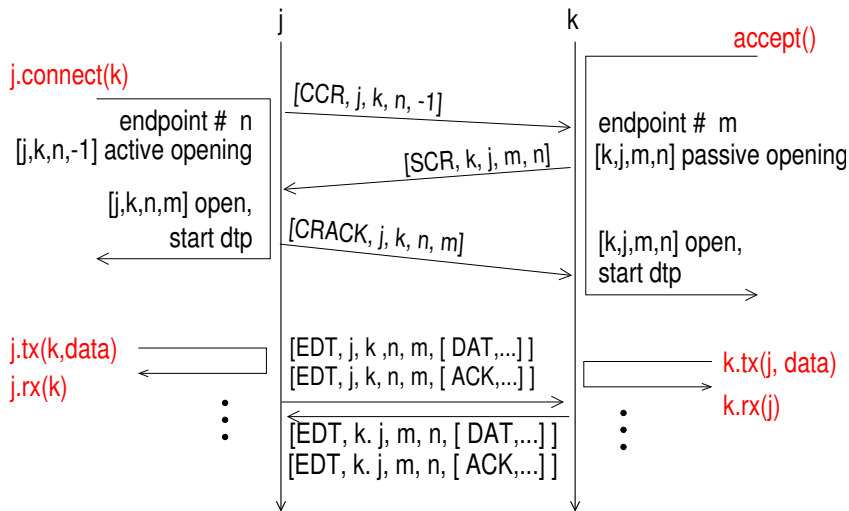
2-way handshake

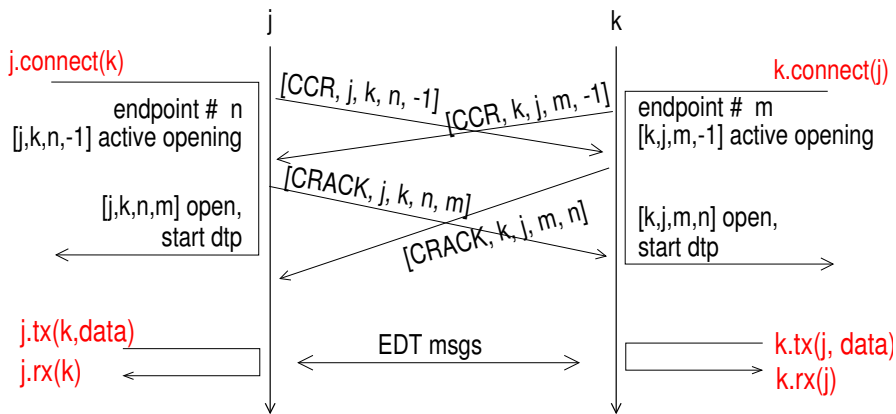


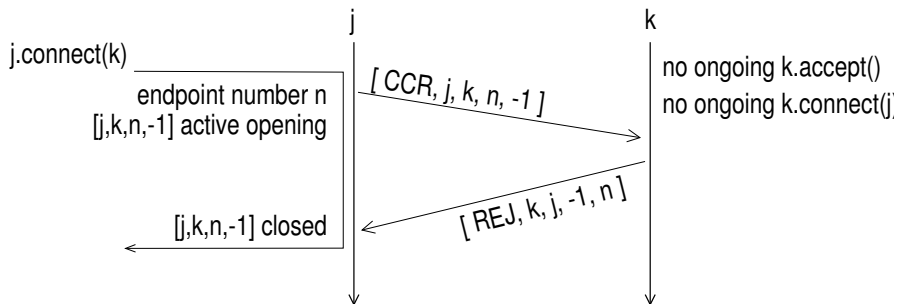
3-way handshake



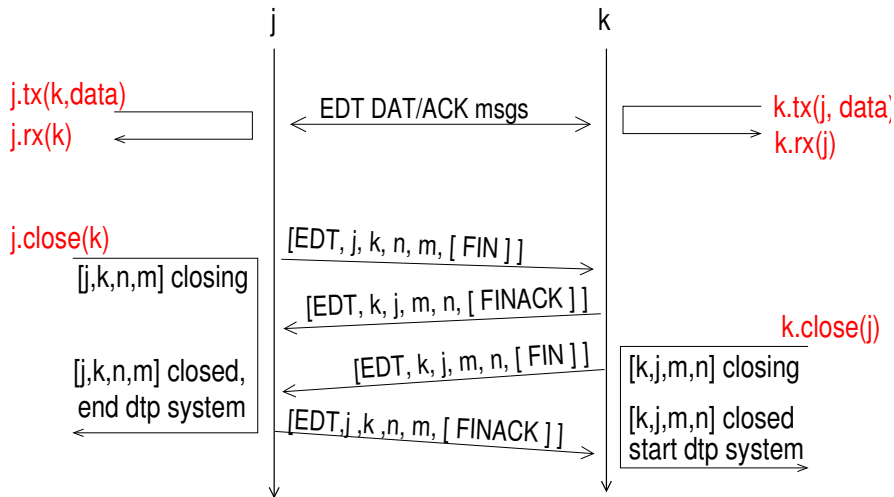
- Client-server connection: one 3-way handshake
- Client-client connection: two simultaneous 2-way handshakes
- Connect rejection: two simultaneous 2-way handshakes
- Data transfer: 2-way handshakes
- Close: two 2-way handshakes // one for each direction











- To overcome msg loss
  - non-final msgs resent until response rcvd // primary msg
  - final msg sent only in response // secondary msg
  - final msg sent even if handshake over from sender's perspective
- Old msgs can start handshakes, which need to be ended
  - handle via reject msgs / increasing ep numbers
  - example
    - accepting k rcvs old  $[CCR, j, k, m, -1]$
    - starts pop  $[k, j, n, m]$ , (re)sends  $[SCR, k, j, n, m]$
    - if j not aop to k: j rejects SCR
    - if j aop  $[j, k, p, -1]$ : k becomes pop  $[k, j, q, p]$

- [CCR, j, k, n, m]: **client conn-req**; primary; sent by aop [j, k, n, m]
- [SCR, j, k, n, m]: **server conn-req**; primary; sent by pop [j, k, n, m]
- [CRACK, j, k, n, m]: **conn-req ack**; secondary
  - response to [SCR|CCR, k, j, m, n] by aop/open [j, k, n, m]
- [REJ, j, k, n, m]: **reject**; secondary
  - response to [CCR, k, j, m, n] when not accepting or aop to k
  - response to [SCR, k, j, m, n] when not aop to k
- [EDT, j, k, n, m, dtmsg]: **encapsulates** dtmsg
  - primary if dtmsg is DAT or FIN
  - secondary if dtmsg is ACK or FINACK

- Handshakes can overlap
  - j starts a handshake while k is still in previous handshake
  - msg of later handshake can end the earlier handshake
  - example
    - opening [j,k,n,m] rcvs [EDT,k,j,m,n,.]
    - becomes open

- `ccrBuff`: map  $\langle \text{addr}, \text{ep \#} \rangle$  to store CCRs in server mode
- `ccrBuffk` exists and equals  $n$  iff
  - $j$  in server mode
  - no  $[j,k]$  socket exists
  - at least one  $[CCR,k,j,\dots]$  rcvd
  - $n$  is highest sender ep  $\#$  in these CCRs
- An ongoing `j.accept()` gets its next CCR from `ccrBuff`
  - fifo or priority queuing

Overview

Graceful-closing data transfer protocol

Transport protocol description

Transport protocol program: unbounded endpoint numbers

Transport protocol program: cyclic endpoint numbers

Transport protocol with abort

- Parameters: ADDR,  $j$ ,  $c_j$
- Input fns // called by service users
  - `accept()`, `endServer()`, `connect(k)`, `close(k)` // cm
  - `tx(k, data)`, `rx(k)` // dt
- Input fns // called by  $dtp_k$ 
  - `dtx(k, dtmsg)`, `drx(k)` // send/rcv dtmsg
- Local fns
  - `doRx()` // rcv msg, update state, tx secondary msg
- Output calls
  - `c_j.tx(k, msg)`, `c_j.rx()` // send/rcv msg
- atomicity assumption: `await`
- progress assumption: weak fairness for all threads

- Main

- **ngen**: endpoint number generator // clock/counter, initially 0
- For every **k** st ep  $[j,k]$  exists // initially none
  - **st<sub>k</sub>**: status: aop, pop, rejected, open // no closed
  - **ln<sub>k</sub>**: local number,  $\geq 0$
  - **rn<sub>k</sub>**: remote number,  $-1$  if null
- lff in server mode, no ep  $[j,k,..]$ ,  $\geq 1$  CCR rcvd from k
  - **ccrBuff<sub>k</sub>**: highest sender ep # in these CCRs
- For every open ep  $[j,k,..]$ 
  - **dtp<sub>k</sub>**: sid of dtp system for k
  - **dtpRxQ<sub>k</sub>**: rcvd dtp msgs for dtp<sub>k</sub>
- `startThread ( doRx() )`



- input mysid.accept():
  - while (true)
    - await (not server-mode or ccrBuff not empty)
    - if (not server-mode)
      - return []
    - $k \leftarrow$  earliest key in ccrBuff
    - $(st_k, ln_k, rn_k) \leftarrow$  (pop, ngen++, ccrBuff<sub>k</sub> delete)
    - while (true)
      - if ( $st_k$  is pop)
        - $c_j.tx(k, [SCR, j, k, n, rn_k])$
      - else if ( $st_k = open$ )
        - return [k]
      - else //  $st_k = RJCT$ 
        - delete  $st_k, ln_k, rn_k$ , break

- input mysid.endServer():
  - await (true)
  - exit server-mode
  - for (k in ccrBuff.keys)
    - $c_j.tx(k, [REJ, j, k, ccrBuff_k, -1])$
  - empty ccrBuff
  - await (no ongoing accpet)
  - return

- input mysid.connect(k):
  - await (true)
  - $(st_k, ln_k, rn_k) \leftarrow (aop, ngen++, -1)$
  - while (true)
    - await (true)
    - if ( $st_k$  is aop)
      - $c_j.tx(k, [CCR, j, k, n, rn_k])$
    - else if ( $st_k = open$ )
      - return [k]
    - else //  $st_k = RJCT$ 
      - delete  $st_k, ln_k, rn_k$
      - return []

- input `mysid.tx(k, data)`:  
    `dtpk.tx(k, data)`  
    return
- input `mysid.rx(k)`:  
    `rval ← dtpk.rx()`  
    return `rval`
- input `mysid.dtx(k, dtmsg)`:  
    await (true)  
    `cj.tx([EDT, j, k, 1nk, rnk, dtmsg])`  
    return `rval`
- input `mysid.drx(k)`:  
    await (`dtpRxQk` not empty)  
    remove head of `dtpRxQk`  
    return it

- input `mysid.close(k)`:  
    `dtpk.close()`  
    delete `dtpk, dtpRxQk`  
    return
- `doRx()`: // executed by local thread  
    while (true)  
        (`type, k, m, n, dtmsg`) ← `cj.rx()`  
        `handle<type>(k, m, n, dtmsg)` // dtmsg only for EDT
- helper fn `startDtp(k)`:  
    `dtpk ← start GcDtp(j, k, mysid, ...)`  
    `dtpRxQk ← []`

- helper fn handleCCR(k, m, n):
  - if (no  $st_k$ )
    - if (not server-mode)  $c_j.tx([REJ, j, k, n, m])$
    - else if (no  $ccrBuff_k$  or  $m > ccrBuff_k$ )  $ccrBuff_k \leftarrow m$
  - else if ( $st_k$ ) is aop and  $n = 1n_k$ )
    - $(st_k, rn_k) \leftarrow (open, m)$
    - startDtp(k)
  - else if ( $st_k$ ) is pop and  $m > rn_k$ )
    - $rn_k \leftarrow m$
  - else if ( $st_k$ ) is open)
    - if ( $[m, n] = [rn_k, 1n_k]$ )
      - $c_j.tx([CRACK, j, k, n, m])$
    - else if ( $m > rn_k$ )
      - $c_j.tx([REJ, j, k, n, m])$  //  $dtpRxQ_k.append([FINACK])$  ??

- helper fn handleSCR(k, m, n):
  - if (no  $st_k$ )
    - $c_j.tx([REJ, j, k, n, m])$
  - else if ( $st_k$  is aop and  $n = 1n_k$ )
    - $(st_k, rn_k) \leftarrow (open, m)$
    - startDtp(k)
    - $c_j.tx([CRACK, j, k, n, m])$
  - else if ( $st_k$  is pop and  $m > rn_k$ )
    - $c_j.tx([REJ, j, k, n, m])$
  - else if ( $st_k$  is open)
    - if ( $[m, n] = [rn_k, 1n_k]$ )
      - $c_j.tx([CRACK, j, k, n, m])$
    - else if ( $m > rn_k$ )
      - $c_j.tx([REJ, j, k, n, m])$  //  $dtpRxQ_k.append([FINACK])$  ??

- helper fn handleCRACK(k, m, n):
  - if ( $st_k$  exists and  $st_k$  is aop or pop and  $[m, n] = [rn_k, ln_k]$ )
    - $st_k \leftarrow \text{open}$
    - startDtp(k)
- helper fn handleREJ(k, m, n):
  - if ( $st_k$  exists and  $st_k$  is aop or pop and  $[m, n] = [rn_k, ln_k]$ )
    - $st_k \leftarrow \text{rjctd}$



- helper fn handleEDT(k, m, n, dtmsg):
  - if (no  $st_k$ )
    - if (dtmsg = [FIN])
      - $c_j.tx([EDT, j, k, n, m, , [FINACK]])$
    - else if ( $st_k$  is aop or pop and
      - $[m, n] = [rn_k, ln_k]$ )
        - $st_k \leftarrow open$
        - startDtp(k)
        - $dtpRxQ_k.append(dtmsg)$
    - else if ( $st_k$  is open and
      - $[m, n] = [rn_k, ln_k]$ )
        - $dtpRxQ_k.append(dtmsg)$

Overview

Graceful-closing data transfer protocol

Transport protocol description

Transport protocol program: unbounded endpoint numbers

Transport protocol program: cyclic endpoint numbers

Transport protocol with abort

- Above protocol can use modulo- $N$  endpoint numbers if

$$N \geq \frac{6L + 4W + 2C}{\delta}$$

- $L$ : max message lifetime of the LRD channel
  - $\delta$ : min time between ngen increases (new endpoints) at an addr
  - $W$ : max opening duration of an endpoint
  - $C$ : max open duration of an endpoint
- 
- Note
    - $L, \delta$  arise as in sliding window protocol
    - $W, C$  arise because  $j$  tracks  $k$  only while opening/open to it
    - set  $C$  to 0 for correctness with  $\text{Pr} \approx 1 - (1/N^2)$

- Let  $j$  rcv msg  $[., k, j, m, n]$  when it has ep  $[j, k, \tilde{n}, \tilde{m}]$

Message	Receiving endpoint	Possible tests
$[CCR SCR, k, j, m, n]$	$[pop aop, \tilde{n}, \tilde{m}]$	$n = \tilde{n}; m > \tilde{m};$ $m > ccrBuff[k]$
$[CCR SCR, k, j, m, n]$	$[open, \tilde{n}, \tilde{m}]$	$m = \tilde{m}; n = \tilde{n};$ $m > \tilde{m}$
$[CRACK REJ, k, j, m, n]$	$[aop pop, \tilde{n}, \tilde{m}]$	$m = \tilde{m}; n = \tilde{n}$
$[EDT, k, j, m, n, .]$	$[aop pop open, \tilde{n}, \tilde{m}]$	$m = \tilde{m}; n = \tilde{n}$

- Need  $K$  st  $m - \tilde{m}$  and  $n - \tilde{n}$  wrt above tests

- Following are invariant

$$F_1 : ([j, k, \tilde{n}, \tilde{m}] \text{ exists}) \Rightarrow \tilde{m} \leq k.nGen$$

$$F_2 : ([j, k, \tilde{n}, \tilde{m}] \text{ exists}) \text{ and } ([., k, j, m, n] \text{ rcvbl}) \Rightarrow m \leq k.nGen$$

$$F_3 : ([j, k, \tilde{n}, \tilde{m}] \text{ exists}) \Rightarrow \tilde{n} \leq j.nGen$$

$$F_4 : ([j, k, \tilde{n}, \tilde{m}] \text{ exists}) \text{ and } ([., k, j, m, n] \text{ rcvbl}) \Rightarrow n \leq \tilde{n} b$$

$$F_5 : ([j, k, \tilde{n}, \tilde{m}] \text{ opening}) \Rightarrow k.nGen \leq \tilde{m} + (L + 2W)/\delta$$

$$F_6 : ([j, k, \tilde{n}, \tilde{m}] \text{ open}) \Rightarrow k.nGen \leq \tilde{m} + (C + L + 2W)/\delta$$

- Following are invariant

$$G_1 : [j,k,\tilde{n},\tilde{m}] \text{ opening and } [.,k,j,m,n] \text{ rcvbl} \Rightarrow \\ (m = -1 \text{ or } \tilde{m} = -1 \text{ or } m \leq \tilde{m} + (L + 2W)/\delta)$$

$$G_2 : ([j,k,\tilde{n},\tilde{m}] \text{ open}) \text{ and } ([.,k,j,m,n] \text{ rcvbl}) \Rightarrow \\ m \leq \tilde{m} + (C + L + 2W)/\delta$$

$$G_3 : ([j,k,\tilde{n},\tilde{m}] \text{ opening or open}) \text{ and } ([CCR|SCR,k,j,m,n] \text{ rcvbl}) \Rightarrow \\ m \geq \tilde{m} - (L + W)/\delta \text{ and } n \geq \tilde{n} - (2L + 2W)/\delta$$

$$G_5 : ([j,k,\tilde{n},\tilde{m}] \text{ opening}) \text{ and } ([CRACK,k,j,m,n] \text{ rcvbl}) \Rightarrow \\ m \geq \tilde{m} - (2L + 2W)/\delta \text{ and } n \geq \tilde{n} - (2L + W)/\delta$$

$$G_6 : [j,k,\tilde{n},\tilde{m}] \text{ opening and } [REJ,k,j,m,n] \text{ rcvbl} \Rightarrow \\ m \geq \tilde{m} - (3L + 2W)/\delta \text{ and } n \geq \tilde{n} - (2L + W)/\delta$$

$$G_7 : [j,k,\tilde{n},\tilde{m}] \text{ opening or open and } [EDT,k,j,m,n,..] \text{ rcvbl} \Rightarrow \\ m \geq \tilde{m} - (L + C + W)/\delta \text{ and } n \geq \tilde{n} - (L + C + 2W)/\delta$$

- From the above, the following hold

$$H_1 : m \leq \tilde{m} + (C + L + 2W)/\delta \quad // G_1, G_2$$

$$H_2 : m \geq \tilde{m} - \max(3L + 2W, L + C + W)/\delta \quad // G_3-G_7$$

$$H_3 : n \leq \tilde{n} \quad // F_4$$

$$H_4 : n \geq \tilde{n} - \max(2L + 2W, L + C + 2W)/\delta \quad // G_3-G_7$$

$$H_5 : -\max(3L + 2W, L + C + W)/\delta \leq m - \tilde{m} \leq (C + L + 2W)/\delta \quad // H_1, H_2$$

$$H_6 : -\max(2L + 2W, L + C + 2W)/\delta \leq n - \tilde{n} \leq \tilde{n} \quad // H_3, H_4$$

$$H_7 : K \leq m - \tilde{m}, n - \tilde{n} \leq K, \text{ where } K = (3L + 2W + C)/\delta \quad // H_5, H_6$$

- Ep #s in msgs and vars now range over  $-1..N-1$
- Optional: ngen is now modulo-N
- Tests involving these values are now as follows



Old test	New test
$m = rn[k]$	no change
$n = \lceil n[k]$	no change
$m > ccrBuff[k]$	$1 \leq \text{mod}(m - ccrBuff[a1]) \leq N/2$
$m > rn[k]$	$(rn[k] = -1 \text{ and } m \neq -1)$ or $1 \leq \text{mod}(m - rn[k]) \leq N/2$



Overview

Graceful-closing data transfer protocol

Transport protocol description

Transport protocol program: unbounded endpoint numbers

Transport protocol program: cyclic endpoint numbers

Transport protocol with abort

- Tp now aborts endpoint if response to a primary message not rcvd after K resends
  - returns of functions distinguish between closing (or rejection) and abort
  - use the abortable dtp program and service