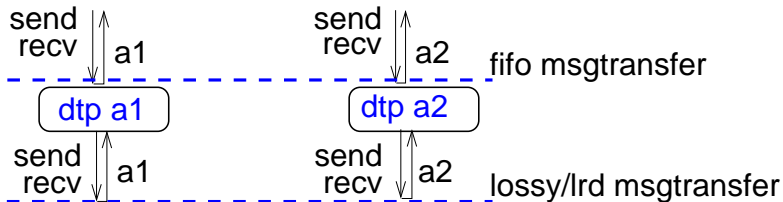


# Implementing fifo msgtransfer2 using lossy or LRD msgtransfer2

Shankar

March 12, 2019

# Data transfer protocol: DtpDist



- Implements  $\text{fifo msgtransfer}(a1, a2)$  using  $\text{lossy|lrd msgtransfer}$
- Sliding window protocol (Swp)
  - $a1 \rightarrow a2$  fifo transfer
  - define programs at “algorithm-level”: atomic rules
  - prove correctness with modulo-N sequence numbers
  - correctness-preserving refinement to await program
- Obtain dtp via correctness-preserving merge of two Swps

Sliding Window Protocol

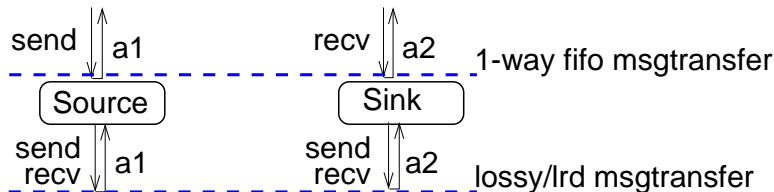
Analysis of Sliding Window Protocol

Await-structured Source and Sink Programs

Data transfer Protocol and Proof

Graceful-closing Data transfer Protocol

Abortable Data transfer Protocol



### ■ Swp(a1, a2)

```

xa1, xa2 ← startSystem( Lossy | LrdChannel(a1, a2) )
ya1 ← startSystem( Source(a1, a2, xa1) )
ya2 ← startSystem( Sink(a2, a1, xa2) )
return ya1, ya2

```

- $d_0, d_1, \dots$ : data blocks sent into source by its user
- Source
  - send  $[d_k, k]$  repeatedly until acked //  $k$ : seq #
- Sink
  - respond with awaited seq # // cumulative ack
    - has  $0 \dots 4 \sqcup 6, \dots 8$ : recv 7 / send 5; recv 5 / send 9
  - passes data blocks to its user in order
- For good throughput
  - $> 1$  outstanding at source, buffer out-of-sequence at sink
- Above requires unbounded seq #s (un)
- not good for hw implementation
- Instead use cyclic seq #s (csn)

- Use mod- $N$   $csn$  instead of  $usn$
- Send  $csn \bmod(k, N)$  instead of  $usn \bmod k$
- Receiver of  $csn$  has to infer the corresponding  $usn$ 
  - maintains **window** of “possible”  $usn$ 
    - say  $L \cdots U$
  - maps rcvd  $csn$  to  $usn$  with same cyclic value
    - $usn \leftarrow L + \text{mod}(csn - L, N)$ ;  
if  $usn > U$  ignore rcvd  $csn$
- Seq #s in transit must remain close to window

## ■ Source

- $ng$ : # blks from user
- $ns$ : # blks sent at least once
- $na$ : # blks acked
- map  $sbuff$ : for blks  $na \dots ng-1$
- $send\ window$ :  $na \dots ns+SW-1$  //  $SW < N$
- $outstanding\ window$ :  $na \dots ns$  // ok if low end is  $na+1$

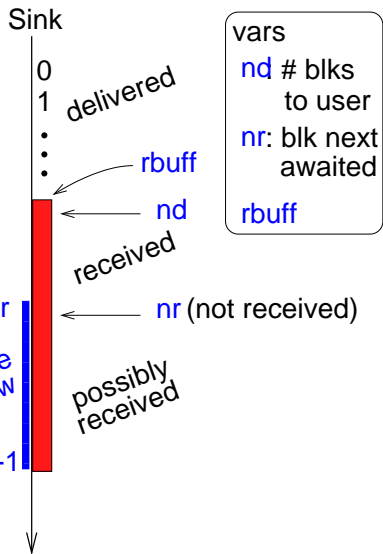
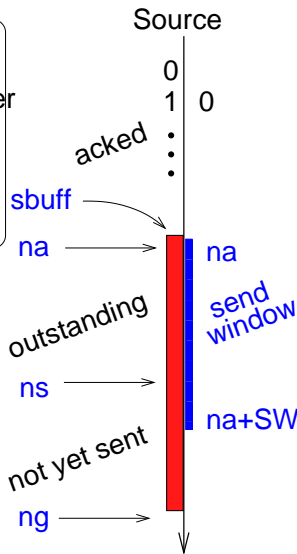
## ■ Sink

- $nd$ : # blks to user
- $nr$ : # contiguous blks rcvd
- map  $rbuff$ : for blks  $nd \dots nd+RW-1$  //  $RW < N$
- $receive\ window$ :  $nr \dots nd+RW-1$

- Over time, windows slide to increasing seq #s

vars

- $ng$  # blks from user
- $ns$  # sent
- $na$  # acked
- $sbuf$



vars

- $nd$  # blks to user
- $nr$ : blk next awaited
- $rbuf$



## Source

```
na, ns, ng  $\leftarrow$  0; sbuff  $\leftarrow$  []  
■ db from user:  
    sbuffng  $\leftarrow$  db; ng++  
■ ns < min(na+SW, ng):  
    send [sbuffns,  $\overline{ns}$ ]; ns++  
■ k in na..ns-1:  
    send [sbuffk,  $\overline{k}$ ]  
■ recv [cn]:  
    j  $\leftarrow$  na +  $\overline{cn-na}$   
    if na < j  $\leq$  ns:  
        sbuff.remove(na..j-1)  
        na  $\leftarrow$  j
```

## Sink

```
nd, nr  $\leftarrow$  0; rbuff  $\leftarrow$  []  
■ nd < nr:  
    rbuffnd to user; nd++  
■ recv [db, cn]:  
    j  $\leftarrow$  nr +  $\overline{cn-nr}$   
    if nr  $\leq$  j < nd+RW:  
        rbuffj  $\leftarrow$  db  
        while nr in rbuff.keys:  
            nr  $\leftarrow$  nr+1  
    send [ $\overline{nr}$ ]  


$\overline{k}$ : mod(k,N)


```

## Sliding Window Protocol

Analysis of Sliding Window Protocol

Await-structured Source and Sink Programs

Data transfer Protocol and Proof

Graceful-closing Data transfer Protocol

Abortable Data transfer Protocol

- **dbh**: auxiliary var indicating seq of blks sent by user

- Desired properties

$X_1 : Inv \ (nd < nr \Rightarrow rbuff_{nd} = dbh_{nd})$

$X_2 : (nd = k < ng \ \textit{leads-to} \ nd > k)$

assuming wfair send sbuff<sub>ns</sub>, resend sbuff<sub>na</sub>,  
 deliver rbuff<sub>nd</sub> to user,  
 source/sink rcv msg,  
 channel progress

- Intermediate desired  $Inv \ A_1 - A_3$

$A_0 : nd \dots nr-1 \text{ in } rbuff.keys$

$A_1 : (k \text{ in } rbuff.keys) \Rightarrow rbuff_k = dbh_k \quad // \ A_{0,1} \Rightarrow X_1.pred$

$A_2 : na \leq nr \leq ns \leq na+SW \quad \text{and} \quad ns \leq ng$

- Add *auxiliary* *usr* field to msgs // not read
  - source sends  $[sbuf_j, \bar{j}, j]$
  - sink sends  $[\overline{nr}, nr]$
- Sink maps rcvd  $[\bar{j}, j]$  wrt rcv window  $nr \dots nr+RW-1$ 
  - $j$  in window: mapped correctly
  - $j < window$ :  $nr-1$  ✓,  $nr-2$  ✓,  $\dots$  ✓,  $nr-N+RW-1$  ✗
  - $j > window$ :  $nr+RW$  ✓,  $nr+RW+1$  ✓,  $\dots$  ✓,  $nr+N$  ✗
  - desired  $Inv A_3$ 

$A_3 : \text{data } j \text{ rcvable} \Rightarrow j \text{ in } nr-N+RW \dots nr+N-1$
- Source maps rcvd  $[\bar{j}, j]$  wrt outstanding window  $na \dots ns$ 
  - desired  $Inv A_3$ 

$A_4 : \text{ack } j \text{ rcvable} \Rightarrow j \text{ in } ns-N+1 \dots na+N$

## ■ Goal

 $A_0 : \text{nd} \dots \text{nr}-1 \text{ in } \text{rbuff.keys}$  $A_1 : (k \text{ in } \text{rbuff.keys}) \Rightarrow \text{rbuff}_k = \text{dbh}_k$  $A_2 : \text{na} \leq \text{nr} \leq \text{ns} \leq \text{na} + \text{SW}$ 

## ■ Correct interpretation

 $A_3 : \text{data } j \text{ rcvble} \Rightarrow j \text{ in } \text{nr}-N+\text{RW} \dots \text{nr}+N-1$  $A_4 : \text{ack } j \text{ rcvble} \Rightarrow j \text{ in } \text{ns}-N+1 \dots \text{na}+N$ ■ For every step  $e$ :  $\{A_0 - A_4\} \text{ e } \{A_0, A_1, A_2\}$  holds

## ■ Suffices to establish

■  $\{A_{2,3,4}\} \text{ e } \{A_3\}$  for  $e$ : send data; rcv data affecting nr■  $\{A_{2,3,4}\} \text{ e } \{A_4\}$  for  $e$ : send ack; rcv ack affecting ns

## Sliding Window Protocol

### Analysis of Sliding Window Protocol

Inv  $A_3-A_4$  for lossy channel

Inv  $A_3-A_4$  for LRD channel

Progress for lossy/LRD channel

### Await-structured Source and Sink Programs

### Data transfer Protocol and Proof

### Graceful-closing Data transfer Protocol

### Abortable Data transfer Protocol

- For SwpDist with lossy channel
  - now show that  $Inv A_3 - A_4$  holds if  $N \geq SW + RW$

$$A_2 : na \leq nr \leq ns \leq na + SW$$

$$A_3 : \text{data } j \text{ rcvble} \Rightarrow j \text{ in } nr - N + RW \dots nr + N - 1$$

- Send data  $k$  preserves  $A_3$ .rhs lower bound if  $N \geq SW + RW$

- $k \geq na$  // guards,  $A_2$   
 $\geq ns - SW$  //  $A_2$   
 $\geq nr - SW$  //  $A_2$   
 $\geq nr - N + RW$  //  $N \geq SW + RW$

- Send data  $k$  preserves  $A_3$ .rhs upper bound

- $k \leq na + SW - 1$  // guards,  $A_2$   
 $\leq nr + SW - 1$  //  $A_2$   
 $\leq nr + N - 1$  //  $N \geq SW$



$$A_2 : na \leq nr \leq ns \leq na + SW$$

$$A_3 : \text{data } j \text{ rcvable} \Rightarrow j \text{ in } nr - N + RW \dots nr + N - 1$$

■  $nr$  increase preserves  $A_3$  if  $N \geq SW + RW$

■ let  $nr$  become  $k$  at time  $t_0$ .

■ so  $k-1$  rcvd at  $t_0$  or prior, so sent at some  $t_1$  //  $t_1 < t_0$

■  $ns(t_1) \geq k+1$  // guards

■  $na(t_1) \geq k+1 - SW$  ★★ //  $A_2$

■ let  $j$  be rcvable after  $t_0$ , so sent at some  $t_2$  after  $t_1$

■  $j \geq na(t_2) \geq na(t_1)$  //  $na$  non-decreasing

$\geq k+1 - SW$  // ★★

$\geq k+1 - N + RW$  //  $N \geq SW + RW$

■ Above is an operational proof; see text for an assertional proof

$$A_2 : na \leq nr \leq ns \leq na+SW$$

$$A_4 : \text{ack } j \text{ rcvable} \Rightarrow j \text{ in } ns-N+1 \dots na+N$$

■ Proof of  $Inv A_4$

1. acks sent have non-decreasing usn // nr non-decreasing
2. acks in transit have non-decreasing usn // no reordering
3. ack usn lower bounded by na // na is an ack usn
4.  $Inv (\text{ack } j \text{ rcvable} \Rightarrow na \leq j \leq nr)$  // 1, 2, 3
5.  $Inv A_4$  // 4,  $A_2$

■ Above is an operational proof; see text for an assertional proof

## Sliding Window Protocol

Analysis of Sliding Window Protocol

Inv  $A_3-A_4$  for lossy channel

Inv  $A_3-A_4$  for LRD channel

Progress for lossy/LRD channel

Await-structured Source and Sink Programs

Data transfer Protocol and Proof

Graceful-closing Data transfer Protocol

Abortable Data transfer Protocol

- For SwpDist with LRD channel
  - obviously  $Inv A_3 - A_4$  does not hold for any  $N$
- But  $Inv A_3 - A_4$  holds if
  - LRD channel has max msg lifetime  $L$
  - min time  $\delta$  between  $ns$  increments
  - $N \geq SW + RW + \frac{L}{\delta}$

$$A_2 : na \leq nr \leq ns \leq na+SW$$

$$A_3 : \text{data } j \text{ rcvble} \Rightarrow j \text{ in } nr-N+RW \dots nr+N-1$$

■  $A_3$ .rhs upper bound holds exactly as in lossy channel case

■  $A_3$ .rhs lower bound holds as follows

1 let data  $j$  be rcvble at  $t_0$

2 data  $j$  was sent at some  $t_1 > t_0 - L$  // 1

3  $j \geq na(t_1) \geq ns(t_1) - SW$  // guards,  $A_2$

4 during  $[t_1, t_0]$ :  $ns$  increases by at most  $L/\delta$  // guards

5  $j \geq ns(t_0) - L/\delta - SW$  // 3, 4

6  $j \geq ns(t_0) - N + RW$  // 5,  $N \geq SW + RW + L/\delta$

7  $j \geq nr(t_0) - N + RW$  // 6,  $A_2$

$$A_2 : na \leq nr \leq ns \leq na+SW$$

$$A_4 : \text{ack } j \text{ rcvable} \Rightarrow j \text{ in } ns-N+1 \dots na+N$$

■  $A_4$ .rhs upper bound holds exactly as in lossy channel case

■  $A_4$ .rhs lower bound holds as follows

1 let ack  $j$  be rcvable at  $t_0$

2 ack  $j$  was sent at some  $t_1 > t_0 - L$  // 1

3  $j = nr(t_1) \geq na(t_1) \geq ns(t_1) - SW$  // 2, guard,  $A_2$

4 during  $[t_1, t_0]$ :  $ns$  increases by at most  $L/\delta$

5  $j \geq ns(t_0) - L/\delta - SW$  // 3, 4

6  $j \geq ns(t_0) - N + 1$  // 5,  $N \geq SW + RW + L/\delta$ ,  $RW > 0$

## Sliding Window Protocol

### Analysis of Sliding Window Protocol

- Inv  $A_3-A_4$  for lossy channel

- Inv  $A_3-A_4$  for LRD channel

- Progress for lossy/LRD channel

### Await-structured Source and Sink Programs

### Data transfer Protocol and Proof

### Graceful-closing Data transfer Protocol

### Abortable Data transfer Protocol

$X_2 : (nd = k < ng \text{ leads-to } nd > k)$

assuming

- wfair send (of  $sbuff_{ns}$ )
- wfair resend (of  $sbuff_{na}$ )
- wfair sink user recv (of  $rbuff_{nd}$ )
- wfair source recv (of msg from lossy/lrd channel)
- wfair sink recv (of msg from lossy/lrd channel)
- lossy/lrd channel progress



$P_1 : na = k < nr \leq ns \text{ leads-to } na > k$

// src resend, snk recv, src recv, channel progress

$P_2 : (na = k = nr < ns \text{ and } nr < nd+RW) \text{ leads-to } nr > k$

// src resend, snk recv, channel progress

$P_3 : (na = k = nr < ns \text{ and } nr = nd+RW)$

$\text{leads-to } nr > k$  // sink user recv,  $P_2$

$P_4 : na = k = nr < ns \text{ leads-to } nr > k$  //  $P_2, P_3$

$P_5 : na = k < ns \text{ leads-to } na > k$  //  $P_4, P_1$

$P_6 : nr = k < ns \text{ leads-to } nr > k$  //  $P_5, \text{Inv } nr \geq na$

$P_7$  :  $(ns = k < ng \text{ and } ns < na+SW)$  *leads-to*  $ns > k$   
// src send

$P_8$  :  $(ns = k < ng \text{ and } ns = na+SW)$  *leads-to*  
 $(ns = k < ng \text{ and } ns < na+SW)$  //  $P_5$

$P_9$  :  $ns = k < ng$  *leads-to*  $ns > k$  //  $P_7, P_8$

$P_{10}$  :  $nr = k < ng$  *leads-to*  $nr > k$  //  $P_9, P_6$

$P_{11}$  :  $nd = k < ng$  *leads-to*  $nd > k$  //  $P_{10}$ , sink user recv

Sliding Window Protocol

Analysis of Sliding Window Protocol

Await-structured Source and Sink Programs

Data transfer Protocol and Proof

Graceful-closing Data transfer Protocol

Abortable Data transfer Protocol

- Given: algorithm-level program  $A$ 
  - init, vars, atomic rules, fairness for rules
- Goal: program  $B$  that implements  $A$ 
  - init, vars, threads, await statements, fairness for threads
- Construct  $B$  as follows
  - include  $A$ -vars
  - each  $A$ -rule  $\rightarrow$  await statement // preserves
  - additional  $B$ -steps do not affect  $A$ -vars // safety
  - allocate local/guest threads to awaits // preserves
  - fairness of  $A$ -rules  $\rightarrow$  fairness of threads // progress

- program Source(aL, aR, xL, N, SW, RW)
  - parameters
    - aL: local addr, aR: remote addr, xL: channel access sid
  - constants: DAT, ACK, RTO // data msg, ack msg, timeout
  - functions
    - input mysid.send: get db from user
    - doSendDat: send ns, resend k // calls xL.send
    - doRecvAck: rcv ack msg // calls xL.recv
  - main
    - ng, ns, na  $\leftarrow$  0, sbuff  $\leftarrow$  []
    - timer  $\leftarrow$  OFF // OFF, 0, 1, ...
    - tSrcSend  $\leftarrow$  startThread( doSendDat() )
    - tSrcRecv  $\leftarrow$  startThread( doRecvAck() )
    - return mysid

```
■ input mysid.send(aR, db):  
    await True:  
        sbuffng ← db  
        ng++  
    return  
  
■ function doSendDat():  
    while True:  
        • await (timer = OFF and ns < min(ng, na+SW)) or  
            (timer > RTO and na < ns):  
            ns ← min(ng, na+SW)  
            for j in na..ns-1:  
                xL.send(aR, [DAT, sbuffj, mod(j,N)])  
            timer ← 0 // (re)start timer
```

- function void doRecvAck():
  - while True:
    - Seq msg  $\leftarrow$  xL.recv();
    - ia {msg in Tuple<ACK, 0..N-1>}
    - await True:
      - int j  $\leftarrow$  na + mod(msg[1] - na, N)
      - if na < j  $\leq$  ns:
        - for k in na..j-1:
          - sbuff.remove(k)
        - na  $\leftarrow$  j
      - if na = ns:
        - timer  $\leftarrow$  OFF; // stop timer
- atomicity assumption: awaits
- progress assumption: weak fairness of threads

- Send new data blocks asap, instead of upon timeout
- Reduce extent of blocking at awaits
  - concurrent access to sbuff
  - duplicate counters
  - ...
- Adapt RTO to measured roundtrip time // congestion control



- program Sink(*aL*, *aR*, *xL*, *N*, *SW*, *RW*)
  - parameters
    - *aL*: local addr, *aR*: remote addr, *xL*: channel access sid
  - constants: DAT, ACK // data msg, ack msg
  - functions
    - input *mysid.recv*: deliver db to user
    - *doRecvDatSendAck*: rcv data, send ack // calls *xL.recv*, *xL.send*
  - main
    - *nd*, *nr*  $\leftarrow$  0, *rbuff*  $\leftarrow$  []
    - *tSnkRecv*  $\leftarrow$  startThread( *doRecvDatSendAck*() )
    - return *mysid*

```
■ input mysid.recv():  
  ● await nd < nr:  
    Seq db ← rbuffnd  
    rbuff.remove(nd)  
    nd++  
    return db
```

```
■ function doRecvDatSendAck():  
  while True:  
    • Seq msg ← xL.recv();  
      ia {msg in Tuple<DAT, 0..N-1, Seq>}  
    await True:  
      int j ← nr + mod(msg[2] - nr, N)  
      if (nr ≤ j < nd+RW)  
        and (not j in rbuff.keys):  
          rbuff[j] ← msg[1]  
          while nr in rbuff.keys:  
            nr++  
          xL.send(aR, [ACK, mod(nr, N)])
```

■ atomicity assumption: awaits

■ progress assumption: weak fairness of threads

Sliding Window Protocol

Analysis of Sliding Window Protocol

Await-structured Source and Sink Programs

Data transfer Protocol and Proof

Graceful-closing Data transfer Protocol

Abortable Data transfer Protocol

- **DtpDist**: merge two SwpDist, but one lossy/lrd channel
- **Dtp**: Source + Sink, but share channel access
  - params, ia, constants: // same as src, snk
  - Source
    - ng, ns, na, sbuff, timer
    - mysid.send()
    - doSendDat(): thread tSrcSend // await{.. xL.send..}
    - doRecvAck(): thread tSrcRecv // xL.recv
  - Sink
    - nd, nr, rbuff
    - mysid.recv
    - ~~doRecvDatSendAck~~: thread tSnkRecv // xL.recv,  
await{.. xL.send..}
  - **doRecvDatAck()**: thread tRecv // xL.recv,  
await{.. xL.send..}
    - rcv data/ack msg, do doRecvAck / doRecvDatSendAck

- params: aL, aR, xL, N, SW, RW
- main
  - ng, ns, na, sbuff, timer
  - nd, nr, rbuff
  - tSrcSend  $\leftarrow$  startThread ( doSendDat() )
  - tSnkRecv  $\leftarrow$  startThread ( doRecvDatAck() )
- input mysid.send ( aR, db)  
    <as in Source(>
- input mysid.recv ()  
    <as in Sink(>
- function doSendDat()  
    <as in Source(>

- function doRecvDatAck():
  - while True:
    - Seq msg  $\leftarrow$  xL.recv();
      - ia {msg in union (Tuple<ACK, 0..N-1>, Tuple<DAT, 0..N-1, Seq>}
    - if msg[0] = ACK:
      - <Source.doRecvAck() update>
    - elif msg[0] = ACK:
      - <Sink.doRecvDatSendAck() update>
- atomicity assumption: awaits
- progress assumption: weak fairness of threads

- Obtain fifo channel inverse
  - only two addresses, so no internal nondeterminism
- Define program Z of DtpDist and fifo channel inverse
- Define assertions that Z must satisfy
- Proof that Z satisfies them follows from Swp properties
  - fifo-channel inverse  $\text{txh}_{a1,a2} = \text{aux var Dtp}(a1).dbh$
  - see text for details



Sliding Window Protocol

Analysis of Sliding Window Protocol

Await-structured Source and Sink Programs

Data transfer Protocol and Proof

Graceful-closing Data transfer Protocol

Abortable Data transfer Protocol

- Allow user to close the data transfer
  - open → closing → closed
  - after call: user cannot send, but can recv
  - returns when
    - remote user also has requested closing
    - all data in both directions have been delivered to users
- New messages
  - [FIN]: indicates closing
  - [FINACK]: ack to [FIN]
- New flags, all initially false:
  - finRcvd
  - finAckRcvd
  - closed

- input mysid.close()
  - wait for na = ng; // all outgoing data acked
  - repeatedly send [FIN] until finAckRcvd is true;
  - wait for finRcvd to be true;
  - wait for nd = nr; // no more incoming data
  - set closed to true;
  - return;

- Augment doRecvDatAck() to handle FIN and FINACK

Change part after “msg ← xL.recv()” to:

- if (not closed)
  - if msg is DAT or ACK: handle as before
  - if msg is FIN: set finRcvd, send [FINACK]
  - if msg is FINACK: set finAckRcvd
- else if (closed and msg is [FIN])
  - send [FINACK]

- Modify mysid.recv

- return [0,data] if nr > nd // incoming data
- return [-1] if finRcvd and nr = nd // no more incoming data

## ■ Safety properties

- if  $j.\text{recv}$  returns  $[\emptyset, \text{data}]$ :  $j.\text{drxh} \circ [\text{data}]$  prefix-of  $k.\text{dtxh}$
- if  $j.\text{recv}$  returns  $[-1]$  or  $j$  is closed:  
 $j.\text{drxh} = k.\text{dtxh}$  and  $k$  is closing or closed.

## ■ Progress properties

- $j.\text{send}$  ongoing *leads-to*  $j.\text{send}$  returns
- $j.\text{recv}$  ongoing,  $j.\text{drxh} \neq k.\text{dtxh}$  *leads-to*  $j.\text{recv}$  returns
- $j.\text{recv}$  ongoing,  $j.\text{drxh} = k.\text{dtxh}$ ,  $k$  is closing or closed  
*leads-to*  $j.\text{recv}$  returns
- $j$  closing,  $k$  is closing or closed  
*leads-to*  $j$  becomes closed or  $j.\text{recv}$  not ongoing

Sliding Window Protocol

Analysis of Sliding Window Protocol

Await-structured Source and Sink Programs

Data transfer Protocol and Proof

Graceful-closing Data transfer Protocol

Abortable Data transfer Protocol

- Graceful-closing terminates connection
  - but not dtp systems or lossy/lrd channel
  - closed dtp system still needs to respond to FIN
  
- To close dtp systems and channel, need to modify Dtp
  - `abort` if FINACK not rcvd within some  $K$  sends of FIN
  - can do the same for data transfer
  - abort: return guest threads, retrieve local threads, end system
  
- Abortable DtpDist provides a weaker service
  - data delivered is prefix of data sent // still holds
  - ~~all data sent is delivered when close returns~~
  - ~~all data sent is eventually delivered~~