

# **Distributed Algorithms**

**Communication Channels  
in Practice**

**Distributed Programming Laboratory**



ÉCOLE POLYTECHNIQUE  
FÉDÉRALE DE LAUSANNE

# **Processes/Channels**

**Processes communicate by message passing through communication channels**

**Messages are uniquely identified and the message identifier includes the sender's identifier**

# Fair-loss links

- ***FL1. Fair-loss:***
- ***FL2. Finite duplication:***
- ***FL3. No creation:***

# Fair-loss links

- ***FL1. Fair-loss:*** If a message is sent infinitely often by  $p_i$  to  $p_j$ , and neither  $p_i$  or  $p_j$  crashes, then  $m$  is delivered infinitely often by  $p_j$
- ***FL2. Finite duplication:*** If a message  $m$  is sent a finite number of times by  $p_i$  to  $p_j$ ,  $m$  is delivered a finite number of times by  $p_j$
- ***FL3. No creation:*** No message is delivered unless it was sent

# Stubborn links

- ***SL1. Stubborn delivery:*** if a process  $p_i$  sends a message  $m$  to a correct process  $p_j$ , and  $p_i$  does not crash, then  $p_j$  delivers  $m$  an infinite number of times
- ***SL2. No creation:*** No message is delivered unless it was sent

# Algorithm (sl)

- **Implements: StubbornLinks (sp2p).**
- **Uses: FairLossLinks (flp2p).**
- **upon event  $\langle \text{sp2pSend}, \text{dest}, m \rangle$  do**
  - **while (true) do**
    - **trigger  $\langle \text{flp2pSend}, \text{dest}, m \rangle$ ;**
- **upon event  $\langle \text{flp2pDeliver}, \text{src}, m \rangle$  do**
  - **trigger  $\langle \text{sp2pDeliver}, \text{src}, m \rangle$ ;**

# Reliable (Perfect) links

- ***Properties***

- ***PL1. Validity:***

- ***PL2. No duplication:*** No message is delivered (to a process) more than once

- ***PL3. No creation:*** No message is delivered unless it was sent

# Reliable (Perfect) links

## • *Properties*

- ***PL1. Validity:*** If  $p_i$  and  $p_j$  are correct, then every message sent by  $p_i$  to  $p_j$  is eventually delivered by  $p_j$
- ***PL2. No duplication:*** No message is delivered (to a process) more than once
- ***PL3. No creation:*** No message is delivered unless it was sent



# Algorithm (pl)

- **Implements: PerfectLinks (pp2p).**
- **Uses: StubbornLinks (sp2p).**
- **upon event < Init> do delivered :=  $\emptyset$ ;**
- **upon event < pp2pSend, dest, m> do**
  - **trigger < sp2pSend, dest, m>;**
- **upon event < sp2pDeliver, src, m> do**
  - **if  $m \notin$  delivered then**
    - **trigger < pp2pDeliver, src, m>;**
    - **add m to delivered;**

# Reliable links

- ✓ **We shall assume reliable links (also called perfect) throughout this course (unless specified otherwise)**
- ✓ **Roughly speaking, reliable links ensure that messages exchanged between correct processes are not lost**

# Reliable FIFO links

- ✓ Ensures properties PL1 to PL3 of perfect links
- ✓ *FIFO*. The messages are delivered in the same order they were sent.

# Algorithm (f11)

- ✓ **Implements: Reliable FIFO links (fp2p).**
- ✓ **Uses: Reliable links (pp2p).**
- ✓ **Relies on acknowledgements messages.**
- ✓ **Acknowledgements are control messages.**

# Algorithm (f11)

- ✓ upon event <init> do
  - ✓ nb\_acks[\*] := 0
  - ✓ nb\_sent[\*] := 0
  
- ✓ upon event <fp2pSend, dest, m> do
  - ✓ wait nb\_acks[dest] = nb\_sent[dest]
  - ✓ nb\_sent[dest] := nb\_sent[dest]+1
  - ✓ trigger <p2pSend, dest, m>

# Algorithm (f11)

- ✓ upon event <pp2pDeliver, src, m> do
  - ✓ trigger <pp2pSend, src, ack>
  - ✓ trigger <fp2pDeliver, src, m>
  
- ✓ upon event <pp2pDeliver, src, ack> do
  - ✓ nb\_ack[src] := nb\_ack[src]+1

# Algorithm (f12)

- ✓ **Implements: Reliable FIFO links (fp2p).**
- ✓ **Uses: Reliable links (pp2p).**
- ✓ **Relies on sequence numbers attached to each message.**
  
- ✓ **upon event <init> do**
  - ✓ **seq\_nb[\*] := 0**
  - ✓ **next[\*] := 0**

# Algorithm (f12)

- ✓ upon event **<fp2pSend, dest, m>** do
  - ✓ **fifo\_m := ( seq\_nb[dest], m )**
  - ✓ **trigger <pp2pSend, dest, fifo\_m>**
  - ✓ **seq\_nb[dest] := seq\_nb[dest]+1**
  
- ✓ upon event **<pp2pDeliver, src, (sn,m)>** do
  - ✓ **wait next[src] = sn**
  - ✓ **trigger <fp2pDeliver, src, m>**
  - ✓ **next[src] := next[src]+1**



# **(f1) vs. (f2)**

- ✓ **(f1) uses 2 messages per applicative message.**
- ✓ **(f1) artificially limits bandwidth if latency is high.**
  
- ✓ **(f2) increases the size of messages.**
- ✓ **Sequence numbers in (f2) have an unbounded size.**

# Algorithm (fl3)

- ✓ **Implements: Reliable FIFO links (fp2p).**
- ✓ **Uses: Reliable links (pp2p).**
- ✓ **Combines acknowledgements and sequence numbers mechanisms.**
- ✓ **An acknowledgement is sent every `ack_int` messages received.**
- ✓ **The sequence numbers are reset when they reach `ack_int x win_size`.**
- ✓ **The sender has to block at the right moment.**

# Algorithm (f13)

- ✓ upon event <init> do
  - ✓ seq\_nb[\*] := 0
  - ✓ next[\*] := 0
  - ✓ ack\_nb[\*] := 0

# Algorithm (f13)

- ✓ upon event  $\langle \text{fp2pSend}, \text{dest}, m \rangle$  do
  - ✓ **wait  $\text{ack\_nb}[\text{dest}] > \text{seq\_nb}[\text{dest}] - \text{win\_size}$**
  - ✓  **$\text{fifo\_m} := (\text{seq\_nb}[\text{dest}], m)$**
  - ✓ **trigger  $\langle \text{pp2pSend}, \text{dest}, \text{fifo\_m} \rangle$**
  - ✓  **$\text{seq\_nb}[\text{dest}] := \text{seq\_nb}[\text{dest}] + 1$**

# Algorithm (f13)

- ✓ upon event  $\langle \text{pp2pDeliver}, \text{src}, (\text{sn}, \text{m}) \rangle$  do
  - ✓ wait  $\text{next}[\text{src}] = \text{sn}$
  - ✓ **trigger  $\langle \text{pp2pSend}, \text{src}, \text{ack} \rangle$**
  - ✓  $\text{next}[\text{src}] := \text{next}[\text{src}] + 1$
  - ✓ **trigger  $\langle \text{fp2pDeliver}, \text{src}, \text{m} \rangle$**
  
- ✓ upon event  $\langle \text{pp2pDeliver}, \text{src}, \text{ack} \rangle$  do
  - ✓  $\text{ack\_nb}[\text{src}] := \text{ack\_nb}[\text{src}] + 1$

# Algorithm (f14)

- ✓ upon event <init> do
  - ✓ seq\_nb[\*] := 0
  - ✓ next[\*] := 0
  - ✓ ack\_nb[\*] := 0

# Algorithm (f14)

- ✓ upon event  $\langle \text{fp2pSend}, \text{dest}, m \rangle$  do
  - ✓ wait  $\text{ack\_nb}[\text{dest}] \times \text{ack\_int} >$   
 $\text{seq\_nb}[\text{dest}] - \text{win\_size} \times \text{ack\_int}$
  - ✓  $\text{fifo\_m} := ( \text{seq\_nb}[\text{dest}] \bmod (\text{win\_size} \times \text{ack\_int}), m )$
  - ✓ trigger  $\langle \text{pp2pSend}, \text{dest}, \text{fifo\_m} \rangle$
  - ✓  $\text{seq\_nb}[\text{dest}] := \text{seq\_nb}[\text{dest}] + 1$

# Algorithm (f14)

- ✓ upon event  $\langle \text{pp2pDeliver}, \text{src}, (\text{sn}, \text{m}) \rangle$  do
  - ✓ wait  $\text{next}[\text{src}] = \text{sn}$
  - ✓ **if  $(\text{sn}+1) \bmod \text{ack\_int} = 0$** 
    - ✓ **trigger  $\langle \text{pp2pSend}, \text{src}, \text{ack} \rangle$**
  - ✓  **$\text{next}[\text{src}] := (\text{next}[\text{src}] + 1) \bmod (\text{win\_size} \times \text{ack\_int})$**
  - ✓ **trigger  $\langle \text{fp2pDeliver}, \text{src}, \text{m} \rangle$**
  
- ✓ upon event  $\langle \text{pp2pDeliver}, \text{src}, \text{ack} \rangle$  do
  - ✓  **$\text{ack\_nb}[\text{src}] := \text{ack\_nb}[\text{src}] + 1$**



# Fair-loss links

- ***FL1. Fair-loss:*** If a message is sent infinitely often by  $p_i$  to  $p_j$ , and neither  $p_i$  or  $p_j$  crashes, then  $m$  is delivered infinitely often by  $p_j$
- ***FL2. Finite duplication:*** If a message  $m$  is sent a finite number of times by  $p_i$  to  $p_j$ ,  $m$  is delivered a finite number of times by  $p_j$
- ***FL3. No creation:*** No message is delivered unless it was sent

# Stoppable Stubborn links

- ***SL1. Stubborn delivery:*** if a process  $p_i$  sends a message  $m$  to a correct process  $p_j$ , and  $p_i$  does not crash, then  $p_j$  delivers  $m$  an infinite number of times **unless  $p_i$  receives a stop event for  $m$**
- ***SL2. No creation:*** No message is delivered unless it was sent

# Algorithm (ssl)

- **Implements:**  
**StoppableStubbornLinks (ssp2p).**
- **Uses: FairLossLinks (flp2p).**
- **upon event <init> do**
  - **sending =  $\emptyset$**

# Algorithm (ssl)

- upon event  $\langle \text{ssp2pSend}, \text{dest}, m \rangle$  do
  - **add  $m$  to sending**
  - **while ( $m$  in sending) do**
    - **trigger  $\langle \text{flp2pSend}, \text{dest}, m \rangle$ ;**
- upon event  $\langle \text{flp2pDeliver}, \text{src}, m \rangle$  do
  - **trigger  $\langle \text{ssp2pDeliver}, \text{src}, m \rangle$ ;**

# Algorithm (ssl)

- upon event **< flp2pDeliver, src, m > do**
  - **trigger < ssp2pDeliver, src, m >;**
- upon event **<ssp2pStop, m >**
  - **remove m from sending**

# Perfect Stoppable Links

## • *Properties*

- ***PL1. Validity:*** If  $p_i$  and  $p_j$  are correct, then every message sent by  $p_i$  to  $p_j$  is eventually delivered by  $p_j$  **unless  $p_i$  receives a stop event for  $m$**
- ***PL2. No duplication:*** No message is delivered (to a process) more than once
- ***PL3. No creation:*** No message is delivered unless it was sent

# Algorithm (psl)

- **Implements: PerfectStoppableLinks (psp2p).**
- **Uses: StubbornStoppableLinks (ssp2p).**
- **upon event < Init> do delivered :=  $\emptyset$ ;**
- **upon event < psp2pSend, dest, m> do**
  - **trigger < ssp2pSend, dest, m>;**
- **upon event < ssp2pDeliver, src, m> do**
  - **if  $m \notin$  delivered then**
    - **trigger < psp2pDeliver, src, m>;**
    - **add m to delivered;**

# Algorithm (psl)

- upon event  $\langle \text{psp2pStop}, m \rangle$  do
  - trigger  $\langle \text{ssp2pStop}, m \rangle$



# Algorithm (f15)

- ✓ **Implements: Reliable FIFO links (fp2p).**
- ✓ **Uses: Perfect Stoppable Links (psp2p).**
- ✓ **Relies on acknowledgements messages.**
- ✓ **Acknowledgements are control messages.**

# Algorithm (f15)

- ✓ upon event  $\langle \text{psp2pDeliver}, \text{src}, (\text{sn}, \text{m}) \rangle$  do
  - ✓ wait  $\text{next}[\text{src}] = \text{sn}$
  - ✓ if  $(\text{sn}+1) \bmod \text{ack\_int} = 0$ 
    - ✓ trigger  $\langle \text{psp2pSend}, \text{src}, \text{ack} \rangle$
  - ✓  $\text{next}[\text{src}] := (\text{next}[\text{src}] + 1) \bmod (\text{win\_size} \times \text{ack\_int})$
  - ✓ trigger  $\langle \text{fp2pDeliver}, \text{src}, \text{m} \rangle$
- ✓ upon event  $\langle \text{psp2pDeliver}, \text{src}, \text{ack} \rangle$  do
  - ✓  $\text{ack\_nb}[\text{src}] := \text{ack\_nb}[\text{src}] + 1$
  - ✓ **trigger psp2pStop for all messages associated with ack**

# Reliable Broadcast in Practice

- ✓ **What is the problem with (rb) on top of (beb) in practice ?**
  - **> scalability**

# Reliable Broadcast in Practice

- ✓ **What is the problem with (rb) on top of (beb) in practice ?**
  - **> scalability**
  
- ✓ **upon event <bebBroadcast, m> do**
  - ✓ **forall pi in S do**
    - **trigger <pp2pSend, pi, m>**

# Problem with rb/beb

- ✓ **1 process does all the work !**
- ✓ **We need to parallelize**

# Algorithm (gossip)

- ✓ **Implements: ReliableBroadcast (rb).**
- ✓ **Uses: Perfect Links (pp2p).**
- ✓ **Relies on spreading messages in a randomized way**
- ✓ **Every process forwards messages to random peers**
- ✓ **Probabilistic guarantees**
  - > **liveness with probability 1**

# Algorithm (gossip)

- ✓ **upon event <init> do**
  - ✓ **delivered =  $\emptyset$**
  - ✓ **while (true)**
    - **for each m in delivered do**
      - **p = random process**
      - **trigger <pp2pSend, p, m>**

# Algorithm (gossip)

- ✓ upon event **<rbBroadcast, m>**
  - ✓ add m to delivered
  - ✓ trigger **<rbDeliver, self, m>**
  
- ✓ upon event **<pp2pDeliver, src, m>** do
  - ✓ if **m ∉ delivered** then
    - add m to delivered
    - trigger **<rbDeliver, src, m>**



# **Gossip**

## **Experiment**