



UCL

IGP Convergence in a SP network Sub-second and Beyond Part II : The network

Olivier Bonaventure and Pierre François

Dept. Computing Science and Engineering
Université catholique de Louvain (UCL)

<http://www.info.ucl.ac.be/people/OBO>

<http://www.info.ucl.ac.be/pfr>



INGI

Département
d'ingénierie
informatique

Agenda

- Behaviour of IS-IS in ISP networks
 - ● Internet2
 - Tier-1 ISP
- Simulation study
- Towards sub 50 msec failure recovery

The network



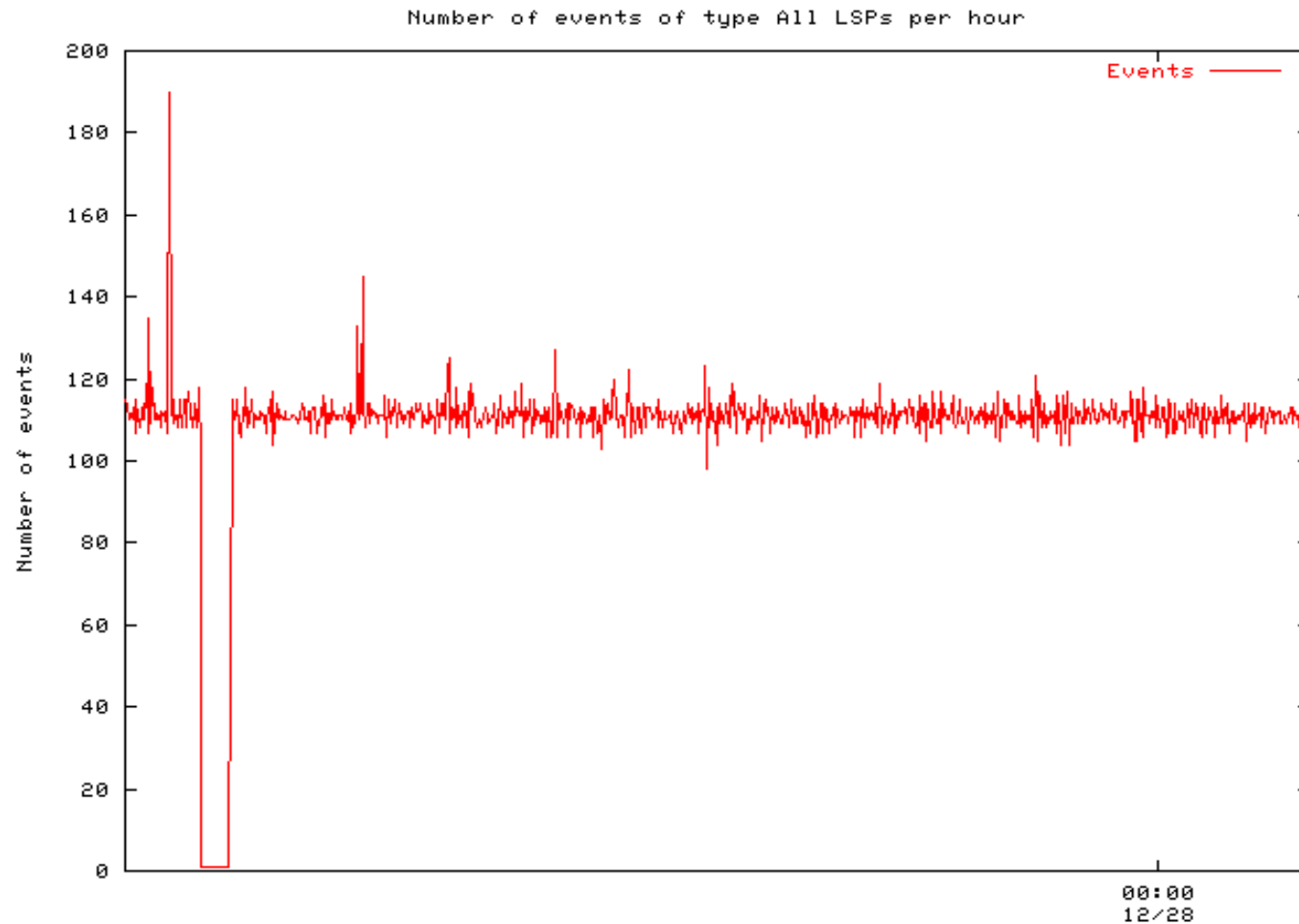
- ISIS data was collected by Abilene on KSCY
- Raw data (December 2004) available from <http://abilene.internet2.edu/observatory/>

Source : <http://abilene.internet2.edu/maps-lists/>

Taxonomy of ISIS events

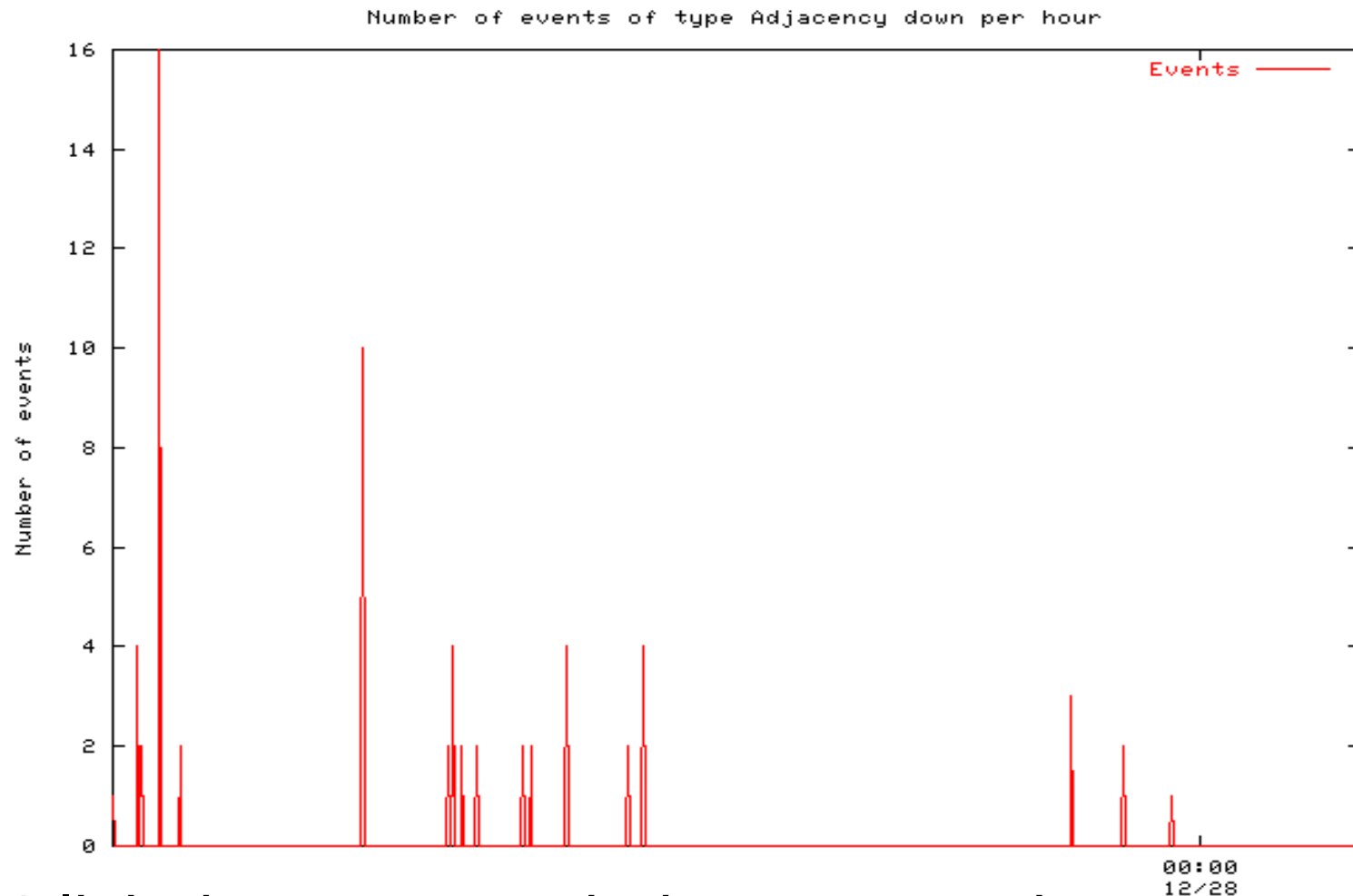
- Refresh LSP
 - ◆ frequency is function of LSP lifetime
- Adjacency down and Adjacency up
 - ◆ link up and link down
- Neighbour metric up or down per hour
 - ◆ Change in link weight for traffic engineering purposes
- IP prefix down or IP prefix up
 - ◆ IP prefix advertised by a router becomes invalid or valid
 - ◆ A change in IP prefix status usually follows link change
- IP prefix metric down or up
 - ◆ Change in the metric associated to a prefix
- Change in Overload bit
 - ◆ Usually set on router reboot during BGP startup
- TE reservation change
 - ◆ Change in reserved bandwidth when MPLS-TE is used
- LSP lifetime set to zero

The ISIS LSPs per hour



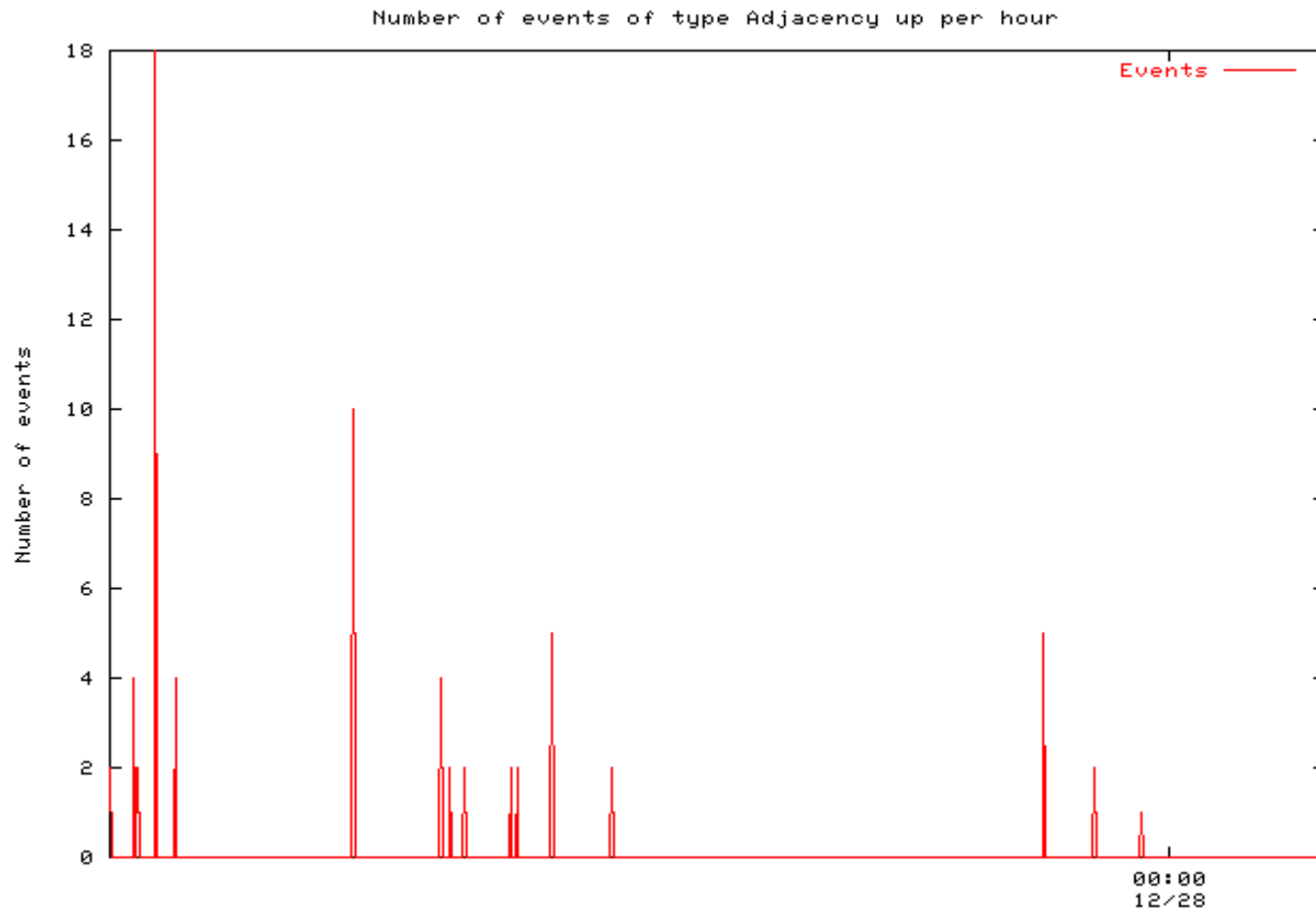
- Most LSPs are simple refresh

The link down events



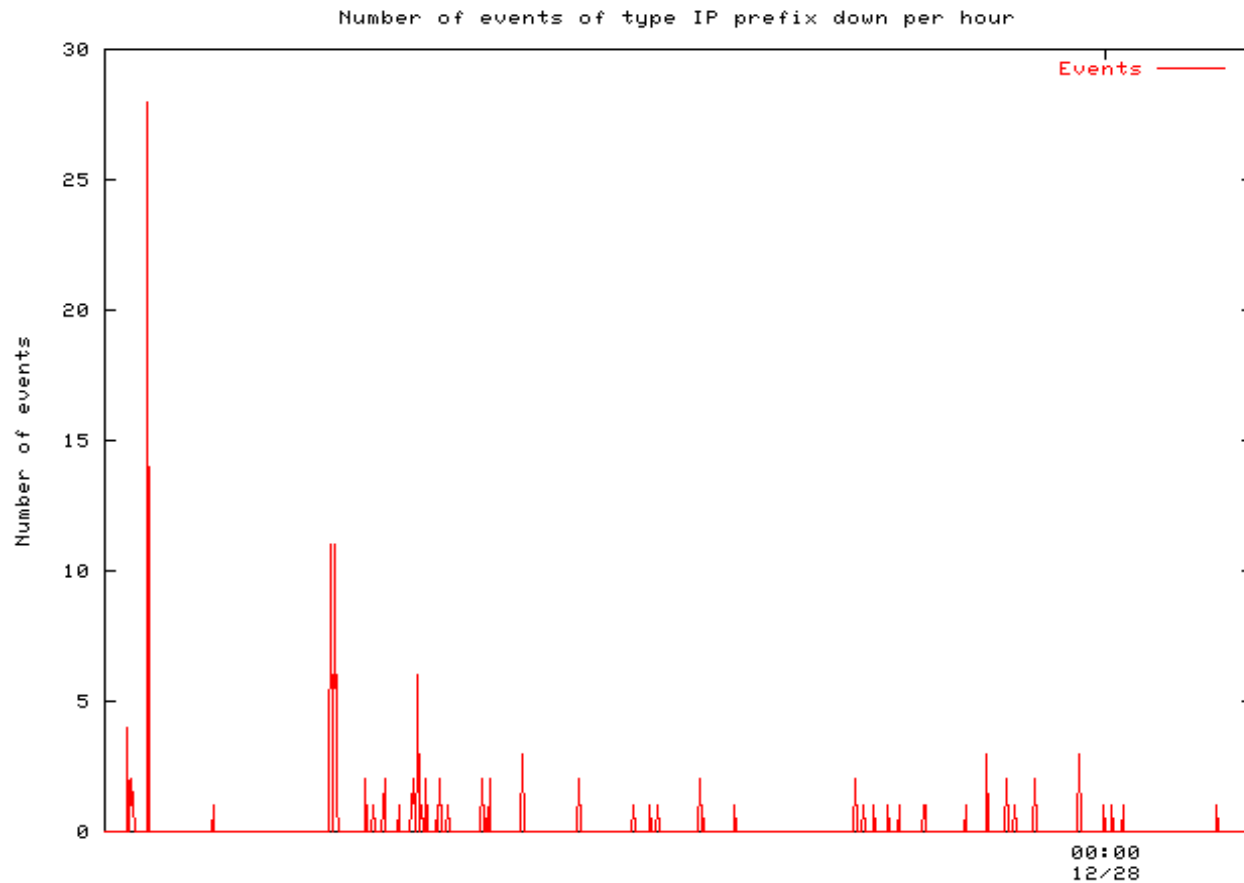
- 71 link down events during one month

The link up events



- 67 link up events during one month

The prefix down events



- 128 prefix down events during one month

Agenda

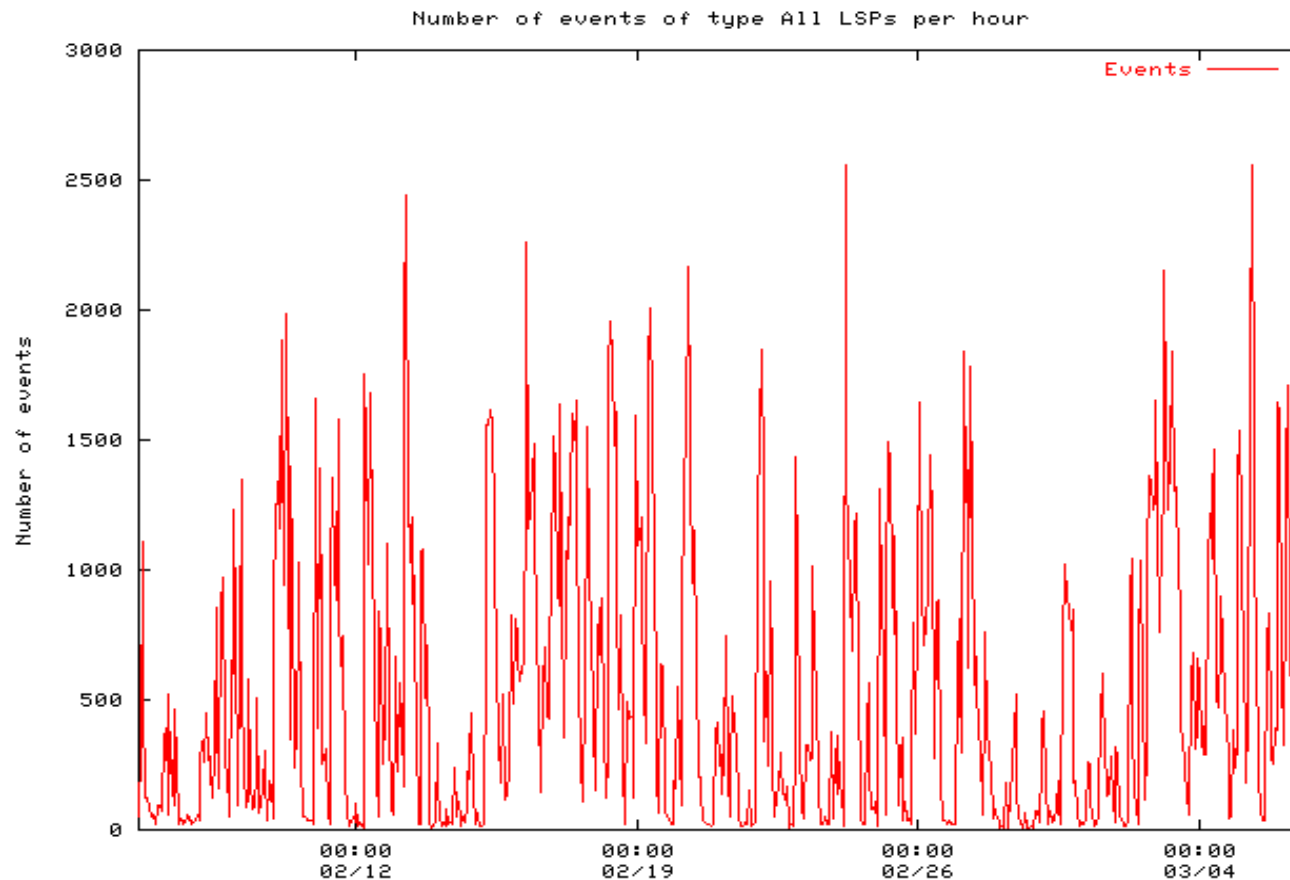
- Behaviour of IS-IS in ISP networks
 - Internet2
 - ● Tier-1 ISP
- Simulation study
- Towards sub 50 msec failure recovery

IS-IS in a tier-1 ISP

- The Network
 - Large tier-1 transit ISP
 - 400 routers in studied ISIS area
 - IS-IS wide metrics and TE extensions are used in the network
 - MPLS traffic engineering is enabled
- The trace
 - IS-IS adjacency between a PC running a modified tcpdump and a router
 - all IS-IS packets logged in libpcap format during one month
 - ◆ analysed with scripts and `lisis`
 - ◆ <http://totem.info.ucl.ac.be/tools.html>

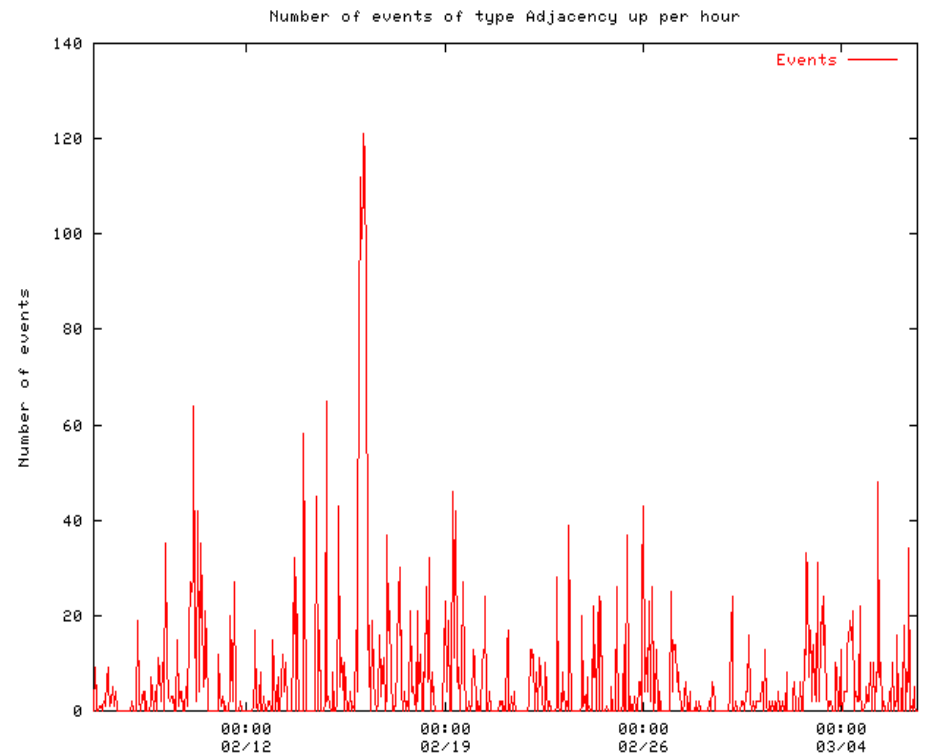
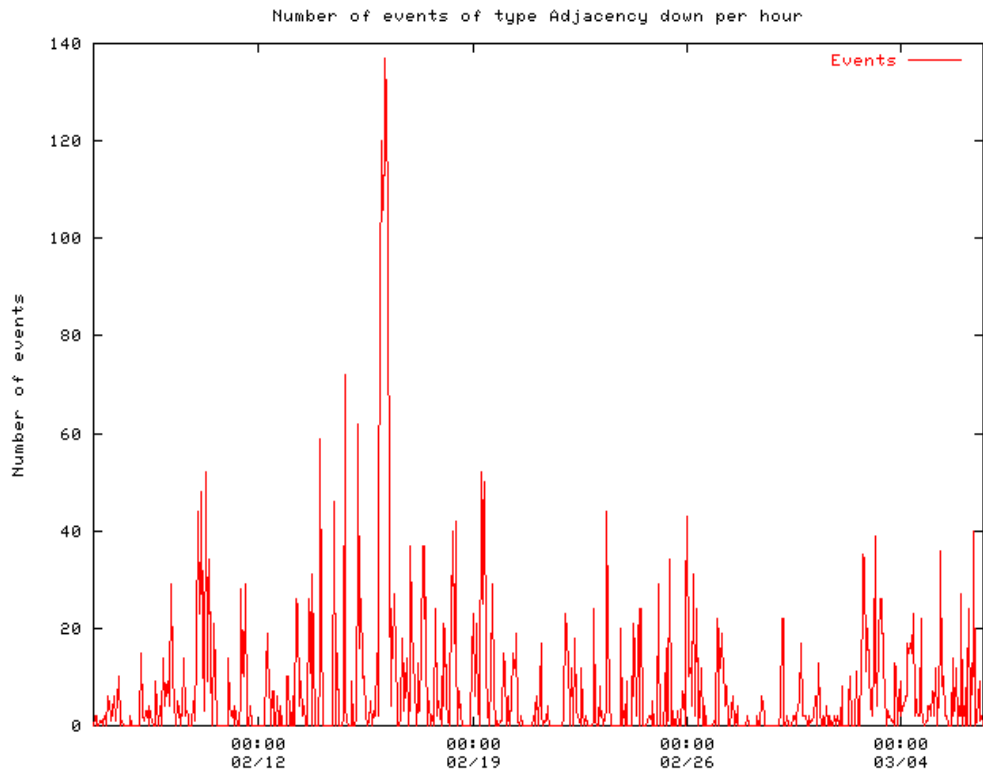
The hourly IS-IS load

- 367383 collected LSPs during one month
 - up to 2500 LSPs per hour...
 - 6% of those LSPs are refresh LSPs
 - ◆ LSP lifetime set to max=65500 seconds



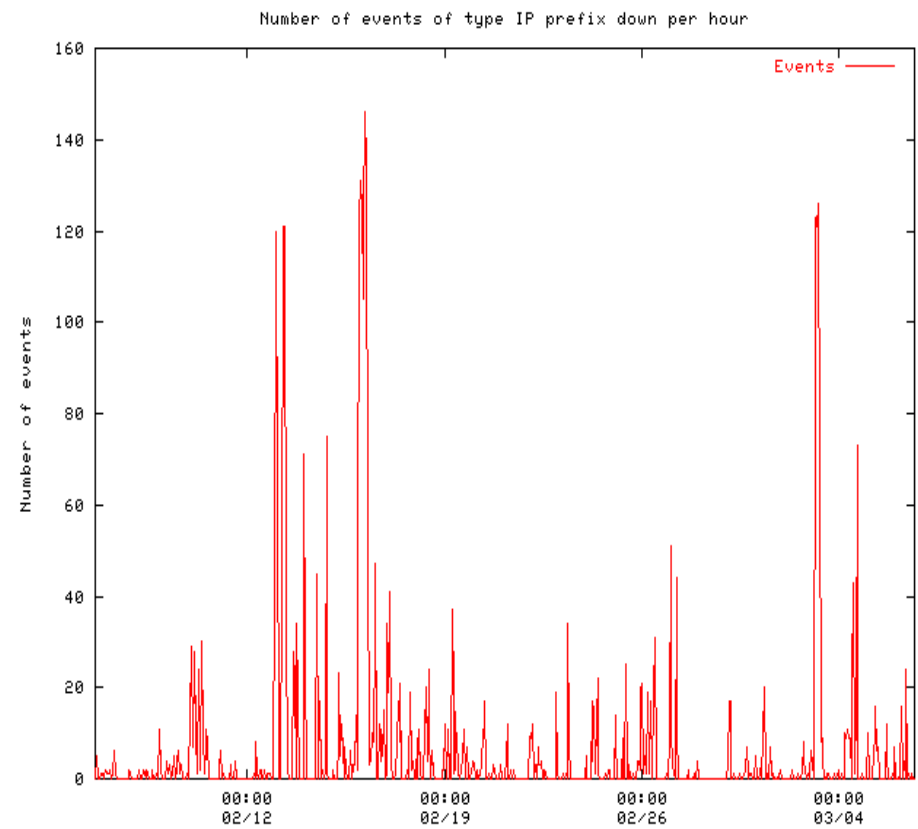
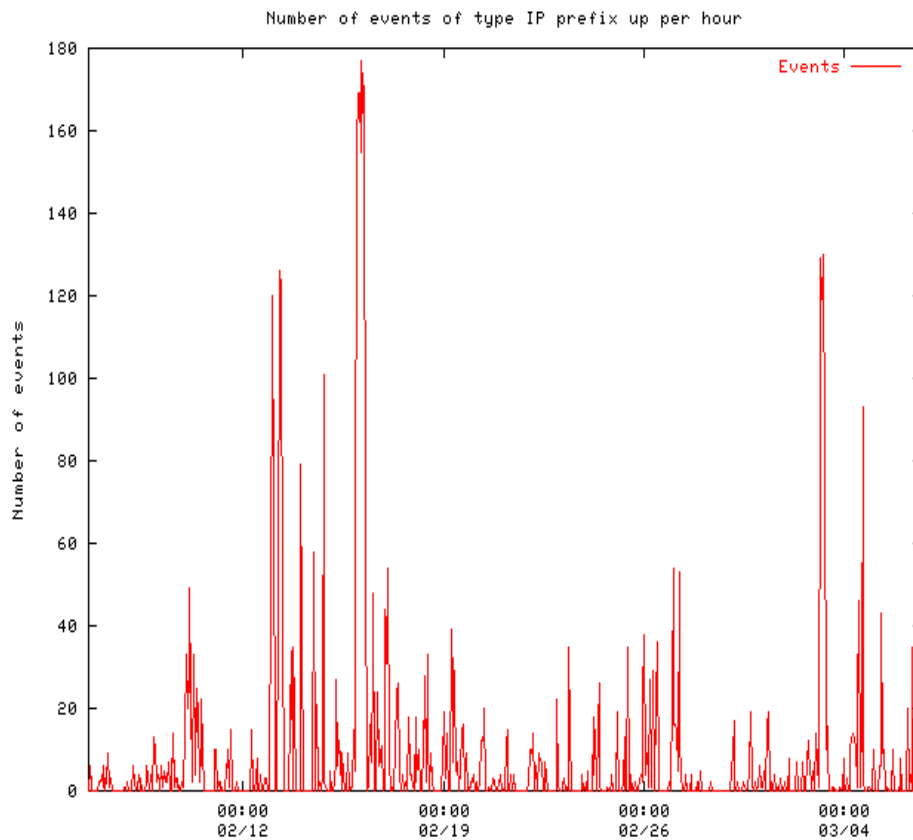
The adjacency changes per hour

- 5276 adjacency down LSPs (left)
 - metric increase events are negligible (40)
- 4487 adjacency up LSPs (right)
 - metric decrease events are negligible



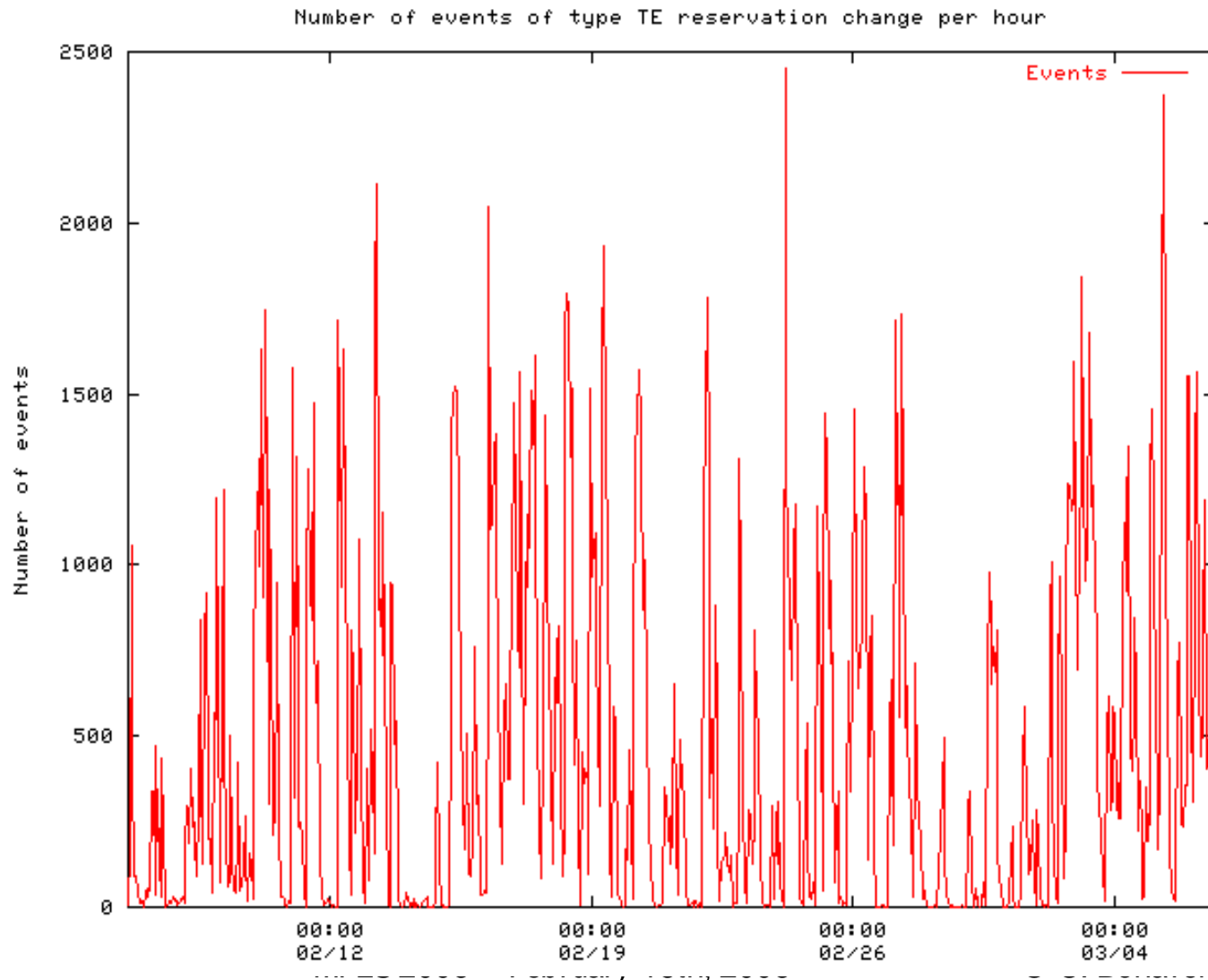
Prefix changes per hour

- Almost no metric changes for prefixes
Prefix up
Prefix down



The LSPs with TE changes

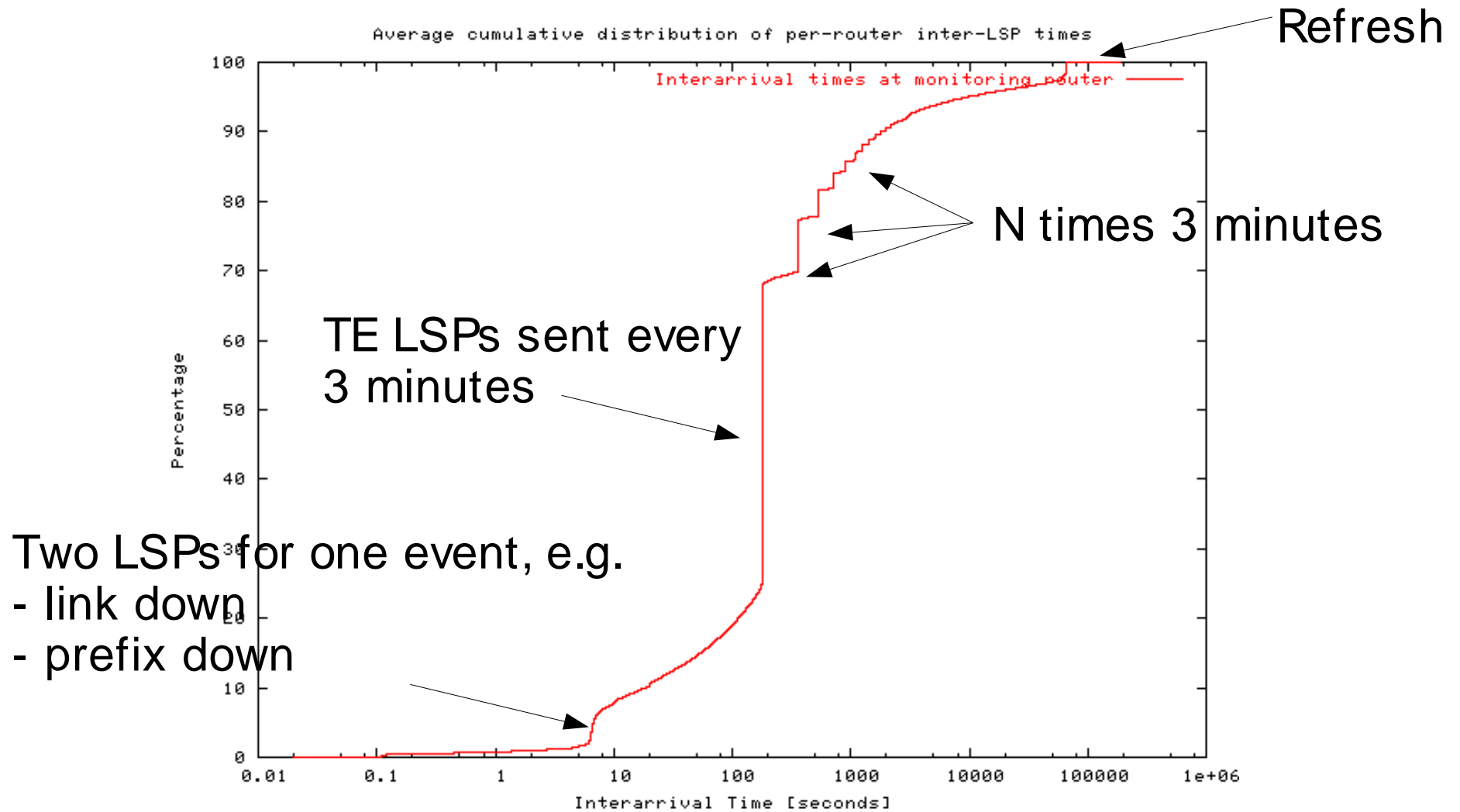
- 88% of all LSPs !



Why so many TE LSPs ?

- Routers are supposed to advertise TE info when crossing percentile thresholds :
 - Up : 15 30 45 60 75 80 85 90 95 96 97 98 99 100
 - Down: 100 99 98 97 96 95 90 85 80 75 60 45 30 15
 - ◆ Such changes are infrequent
 - Most bandwidth changes are only for 10 Kbps
 - ◆ Common value for reserved bandwidth for TE-tunnels with unknown demand
- But, unfortunately those routers also
 - advertise minor TE changes after some delay
 - ◆ default value : 3 minutes

Average per router inter-LSP transmission times



Agenda

- Behaviour of IS-IS in ISP networks

- ● **Simulation study**
 - **Simulation Model**
 - Analysis of link failures
 - Analysis of router failures

- Towards sub 50 msec failure recovery

Simulation model

- Router model follows measurements presented by Clarence Filsfils
- LSP generation
 - ◆ Time to produce a new LSP : 2 milliseconds
 - ◆ upon failure detection, new LSP is produced and placed in LSDB to be flooded
 - ◆ No dampening on the LSP generation
- Failure detection
 - ◆ random delay between [10,15] msec
 - ◆ corresponds to low carrier delay or low BFD timer
 - ◆ larger delay for transoceanic links

Simulation model (2)

- SPF computation time
 - Based on Clarence's Filsfils measurements with some randomness
 - ◆ 2-4 msec for a 22-nodes network
 - ◆ 20-30 msec for a 200-nodes network
- FIB update time
 - Incremental or full FIB update
 - 100-110 microseconds per prefix
 - model uses real prefixes from ISP
- SPF and FIB have exclusive access to CPU
 - No LSP arrival/flooding occurs during SPF+FIB
- Exponential backoff for SPF computation
 - Initial wait : 10, 25, 50, 100 msec
 - Exponential increment : 25, 50, 100 msec

Simulation model (3)

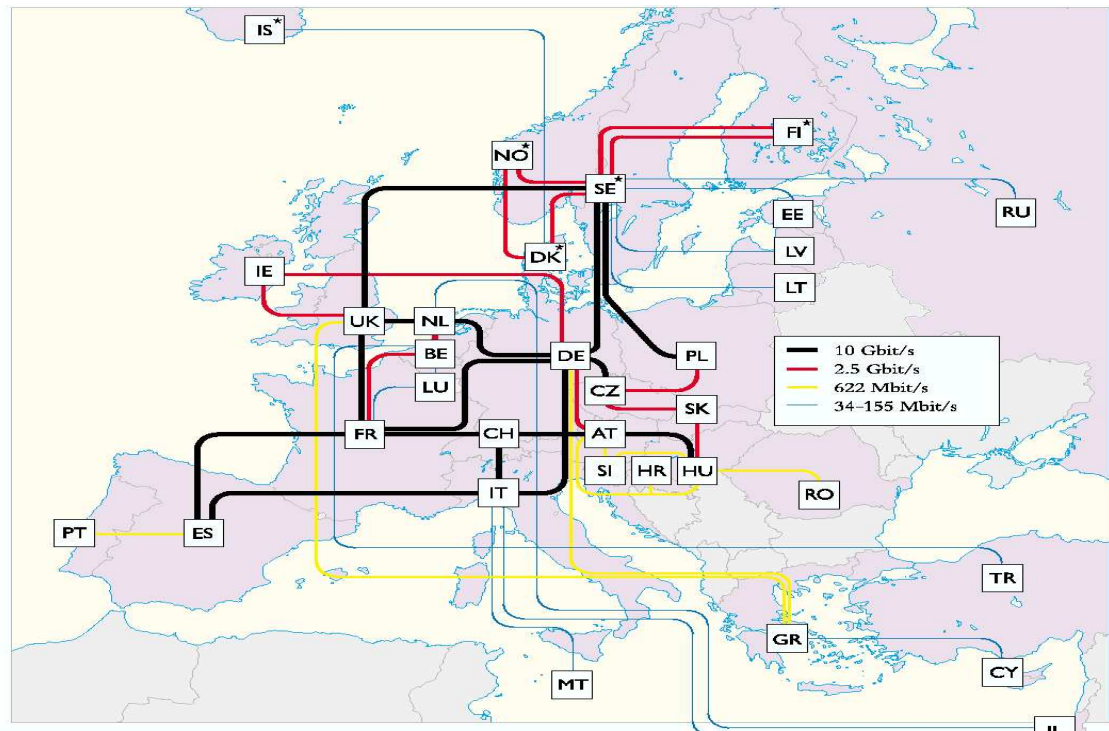
- Normal flooding
 - Timer-based
 - ◆ When timer expires, LSDB is parsed to determine whether a LSP needs to be flooded
 - ◆ Default pacing timer
 - ◆ 33 msec on Cisco
 - Flooding does not run during SPF or FIB
 - ◆ If timer expires during SPF/FIB, flooding will run after FIB
 - Timer expiration
 - ◆ LSDB is parsed
 - ◆ If one LSP is found, it is flooded and timer restarted
 - ◆ Otherwise, timer is cancelled
 - ◆ Arrival or generation of LSP
 - ◆ If pacing timer is running
 - ◆ place LSP in LSDB
 - ◆ If pacing timer is not running
 - ◆ Flood LSP and start pacing timer

Simulation model (4)

- Fast flooding
 - Enhanced flooding mechanism
 - Bypasses pacing timer
 - LSP arrival
 - ◆ If LSP causes SPF
 - ◆ place LSP in LSDB
 - ◆ Flood LSP
 - ◆ maximum number of fast flooded LSPs is configurable, but simulations currently use infinite value
 - ◆ Otherwise
 - ◆ LSP is placed in LSDB and will be flooded by pacing

Simulated networks

- GEANT



- Core backbone of tier-1 ISP
 - 200+ routers in Europe, USA, Asia and South America

Agenda

- Behaviour of IS-IS in ISP networks

- **Simulation study**
 - Simulation Model
 - ● **Analysis of link failures**
 - Analysis of router failures

- Towards sub 50 msec failure recovery

How to evaluate IGP convergence ?

- Packet-based approach
 - Often used to perform measurements
 - Principle
 - ◆ Starting shortly before the failure, send a constant stream of packets from each router to any router in the network
 - ◆ Count the number of packets that
 - ◆ arrive in sequence at their destination
 - ◆ are sent over failed links
 - ◆ loop in the network due to transient loops
 - ◆ are dropped inside routers due to unreachable destination
 - ◆ Derive convergence time for each source/destination pair affected by the failure
 - Drawbacks
 - ◆ Huge simulation cost as most packets are useless
 - ◆ Each packet takes a sample of the routing table of the routers that it passes through

How to evaluate convergence ? (2)

- The Nettester approach

- After each “physical” failure, detection of a failure of FIB update, check consistency of routing tables for each router-router pair

- Definition

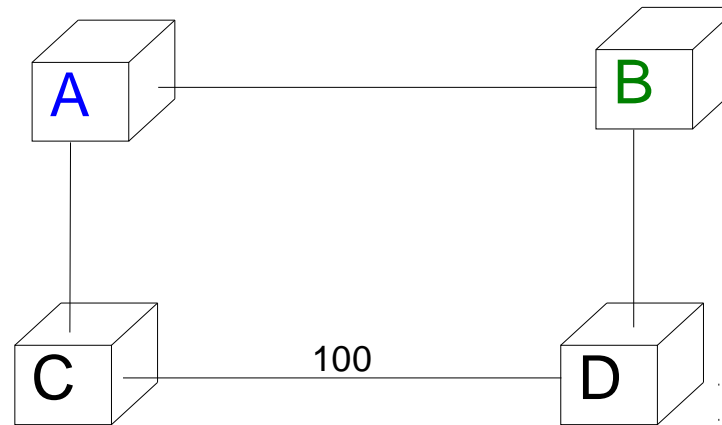
- ◆ Routing is consistent for a pair S-D at time t if all the paths that packets would follow, from S to D, based on the FIB of the routers at time t, are loop-free and finish with D, without passing through a failed link.

- Principle

- ◆ Before the failure, routing is consistent
- ◆ Convergence time is the time when routing becomes and remains consistent for all router-router pairs
 - ◆ Consistency is checked by using the loopback addresses of the routers as source and destination
 - ◆ Note that a packet-based definition could find a lower convergence time than the consistency time

Sample network for link failure

A's FIB
B : East
C : South
D : East



B's FIB
A : West
C : West
D : South

C's FIB
A : North
B : North
D : North

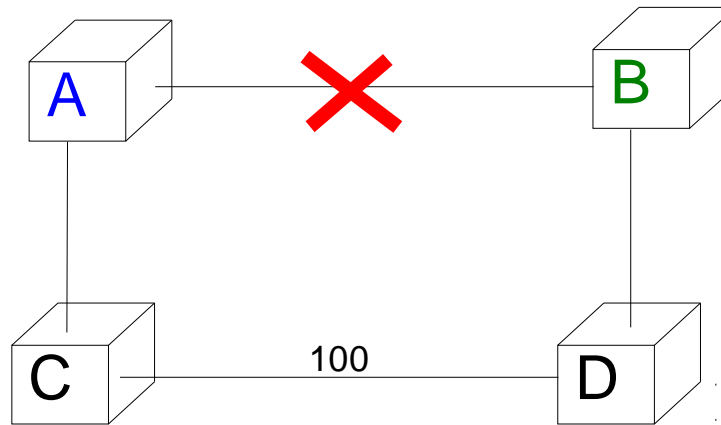
D's FIB
A : North
B : North
C : North

● Parameters

- ◆ 5 msec delay on each link
- ◆ Pacing timer : 33 msec
 - ◆ No fast flooding
- ◆ Initial Wait : 50 msec
- ◆ SPF+ FIB: 0 msec

Example link failure

A's FIB
B : East
C : South
D : East



B's FIB
A : West
C : West
D : South

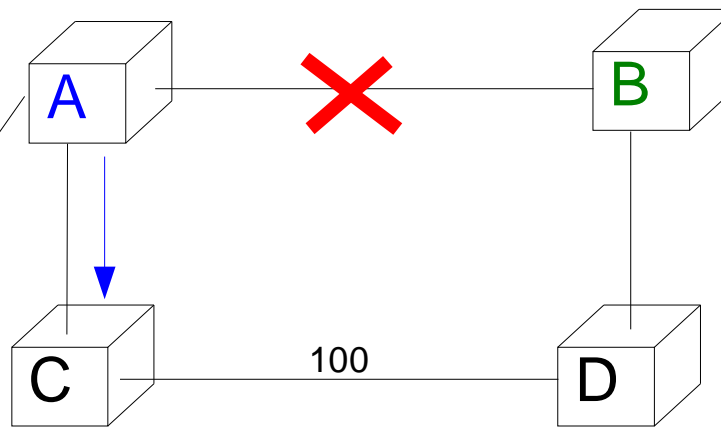
C's FIB
A : North
B : North
D : North

D's FIB
A : North
B : North
C : North

- Nettester at $t=0$ msec : 4 paths out of 12
 - ◆ Link A-B failed but A and B are not yet aware
 - ◆ A can reach C, but not B and D
 - ◆ B can reach D, but not A and C
 - ◆ C can reach A but not B and D
 - ◆ D can reach B but not A and C

Example link failure (2)

A's new FIB
B : Unreach
C : South
D : Unreach



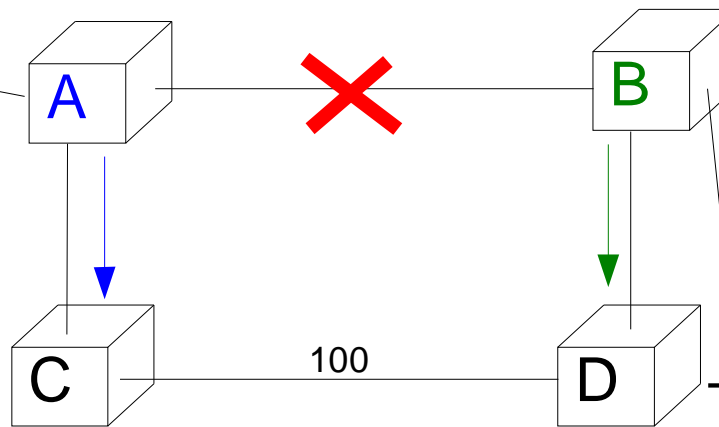
T= 10 msec
Failure detected
T= 12 msec
New LSP(A) produced
SPF will start at 62 msec
New LSP(A) flooded
Pacing expires at 45msec

- **Nettester at t=12 msec : 4 paths out of 12**
 - ◆ A can reach C, but not B and D
 - ◆ B can reach D, but not A and C
 - ◆ C can reach A, but not B and D
 - ◆ D can reach B, but not A and C

Example link failure (3)

Timers

- Pacing at 45 msec
- SPF at 62 msec



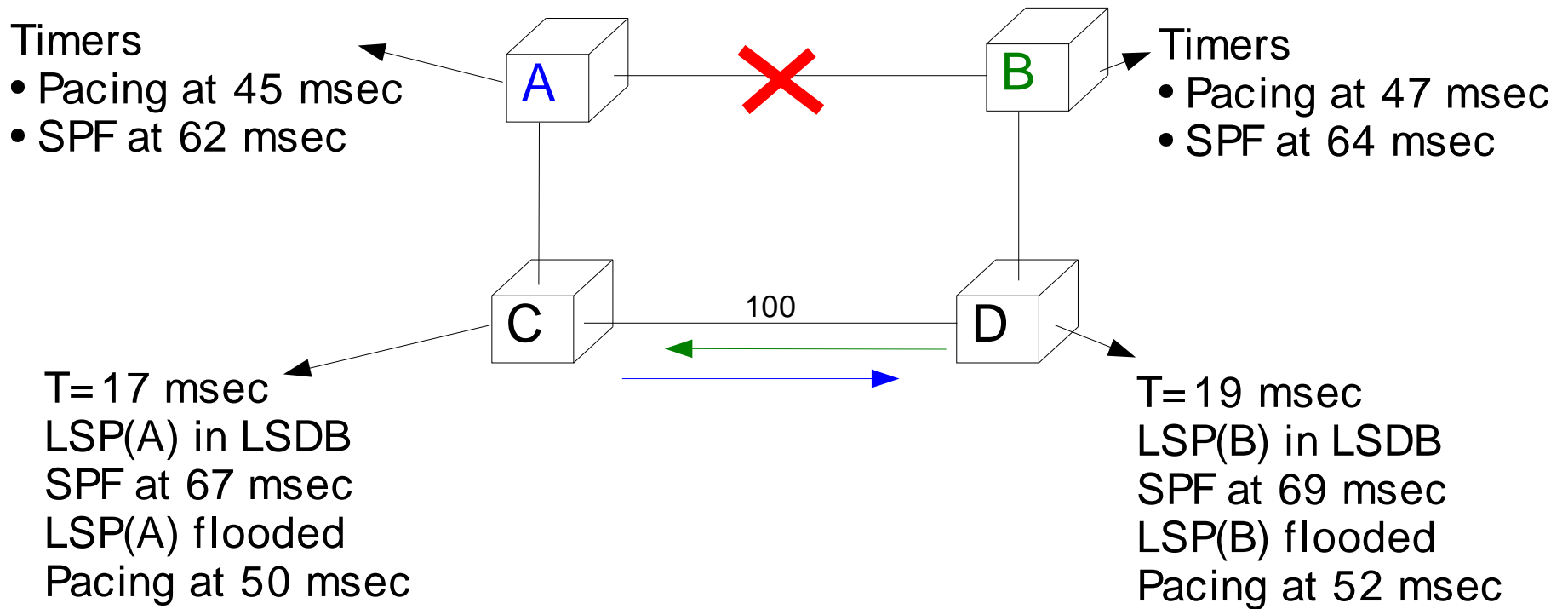
B's FIB
A : Unreach
C : Unreach
D : South

T= 12 msec
Failure detected
T= 14 msec
New LSP(B) produced
SPF start at 64 msec
New LSP(B) flooded
Pacing expires 47 msec

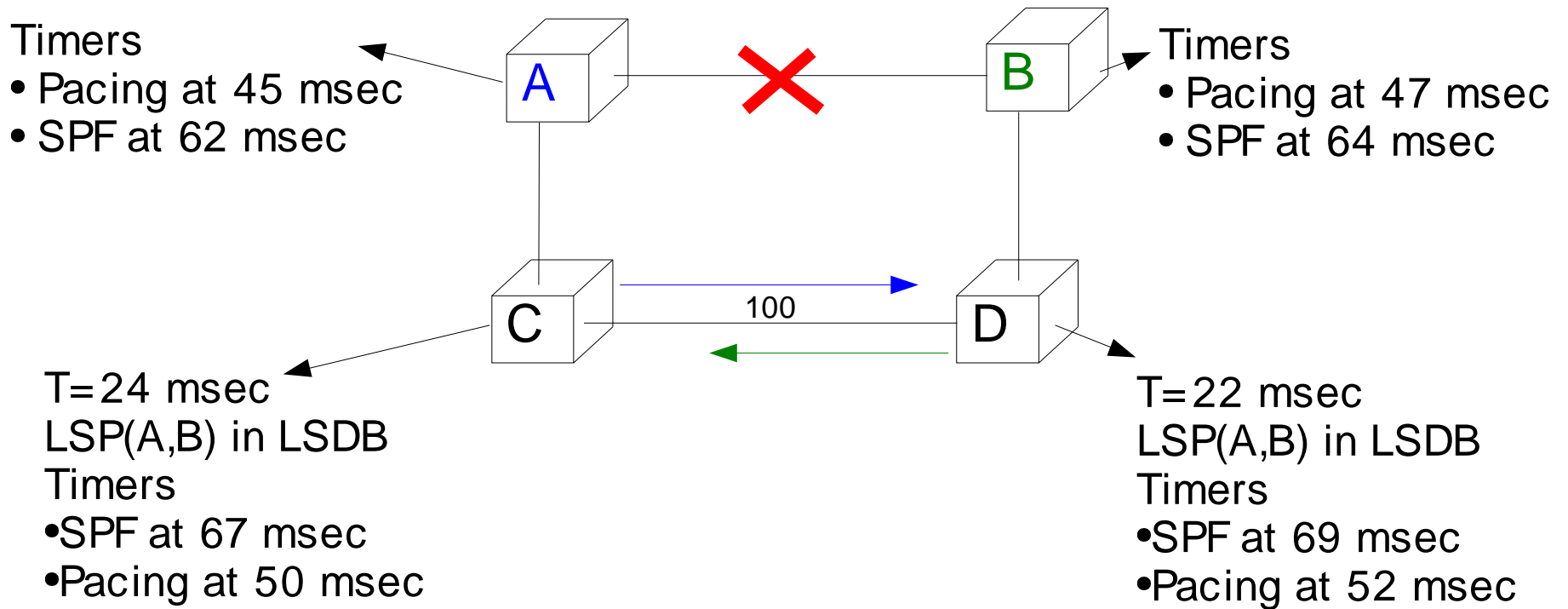
● Nettester at t=14 msec : 4 paths out of 12

- ◆ A can reach C, but not B and D
- ◆ B can reach D, but not A and C
- ◆ C can reach A, but not B and D
- ◆ D can reach B, but not A and C

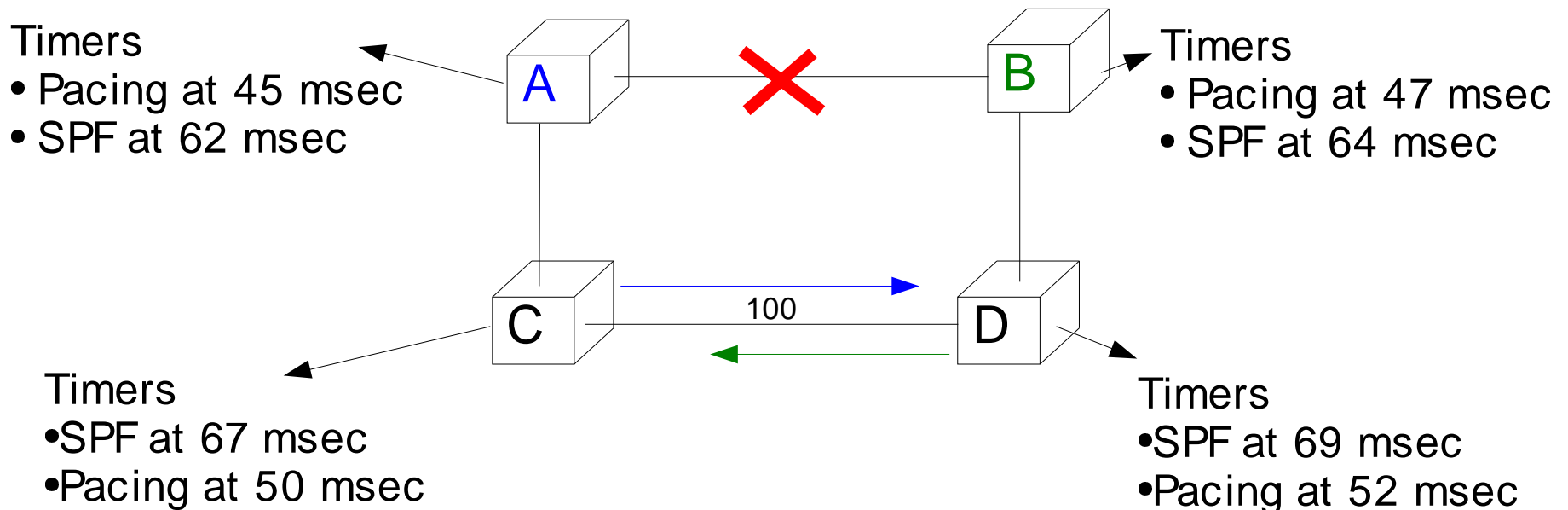
Example link failure (4)



Example link failure (5)



Example link failure (6)



At t=45 msec

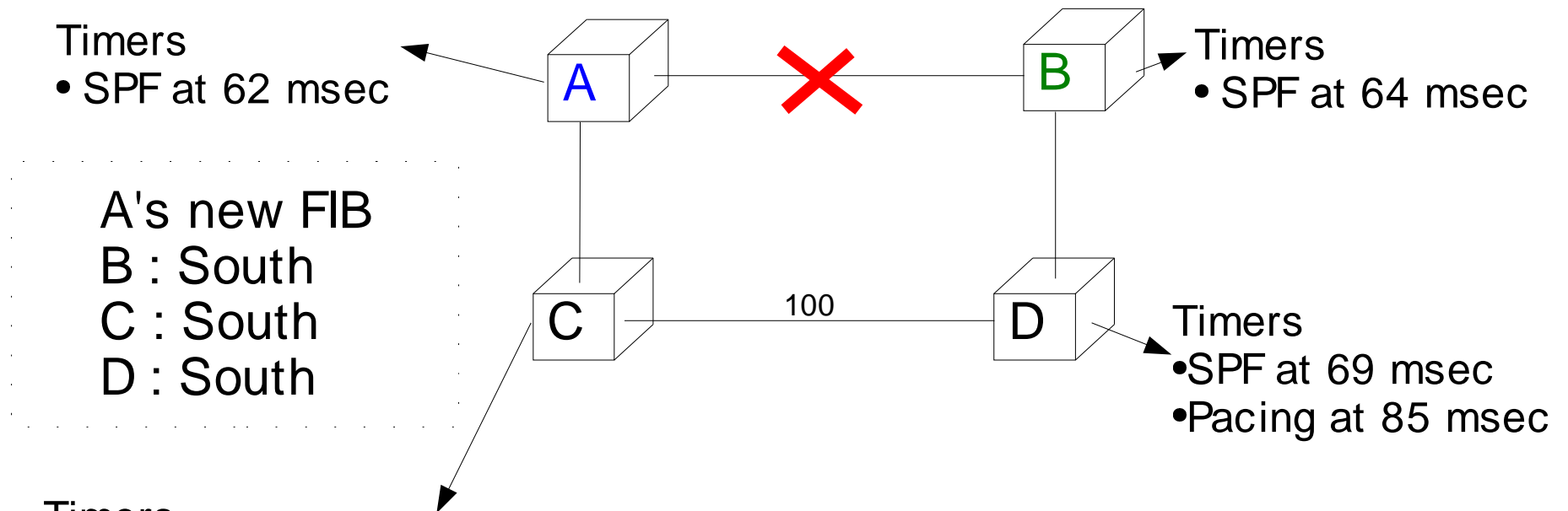
- pacing timer at A cancelled, nothing to flood

At t=47 msec

- pacing timer at B cancelled, nothing to flood

LSP(A) and LSP(B) will eventually reach B and A respectively

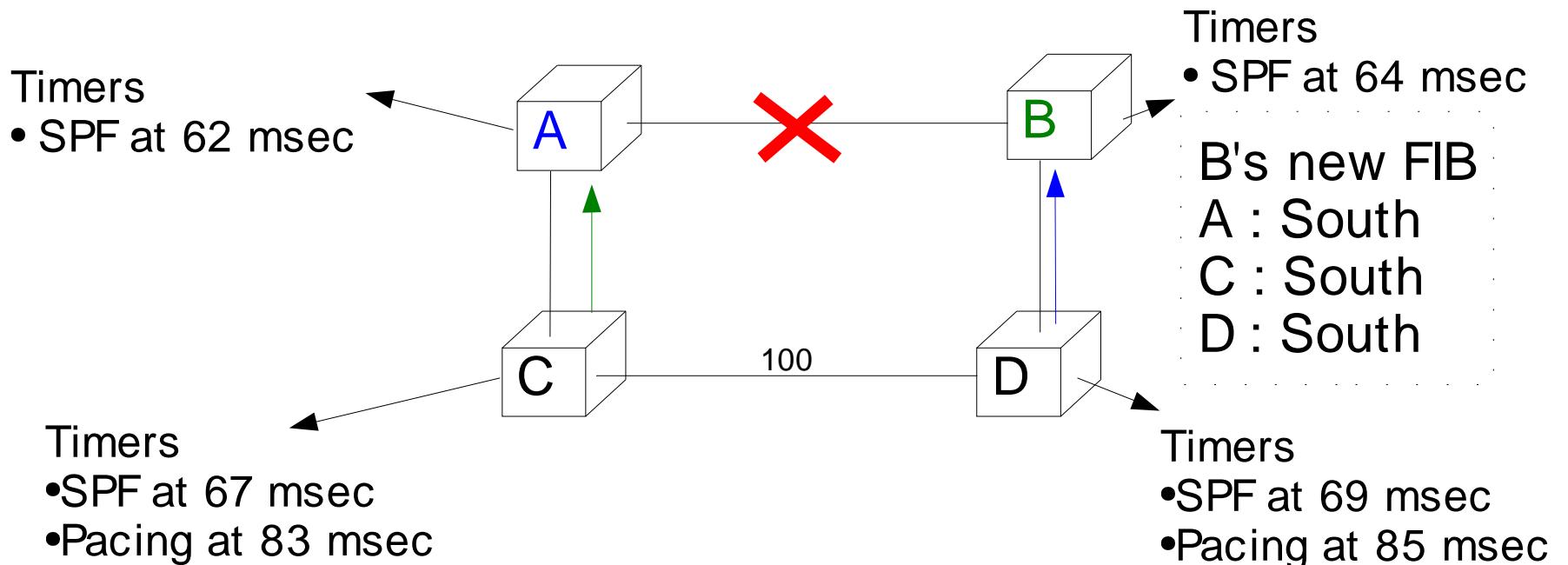
Example link failure (7)



At t=62 msec after FIB update at A

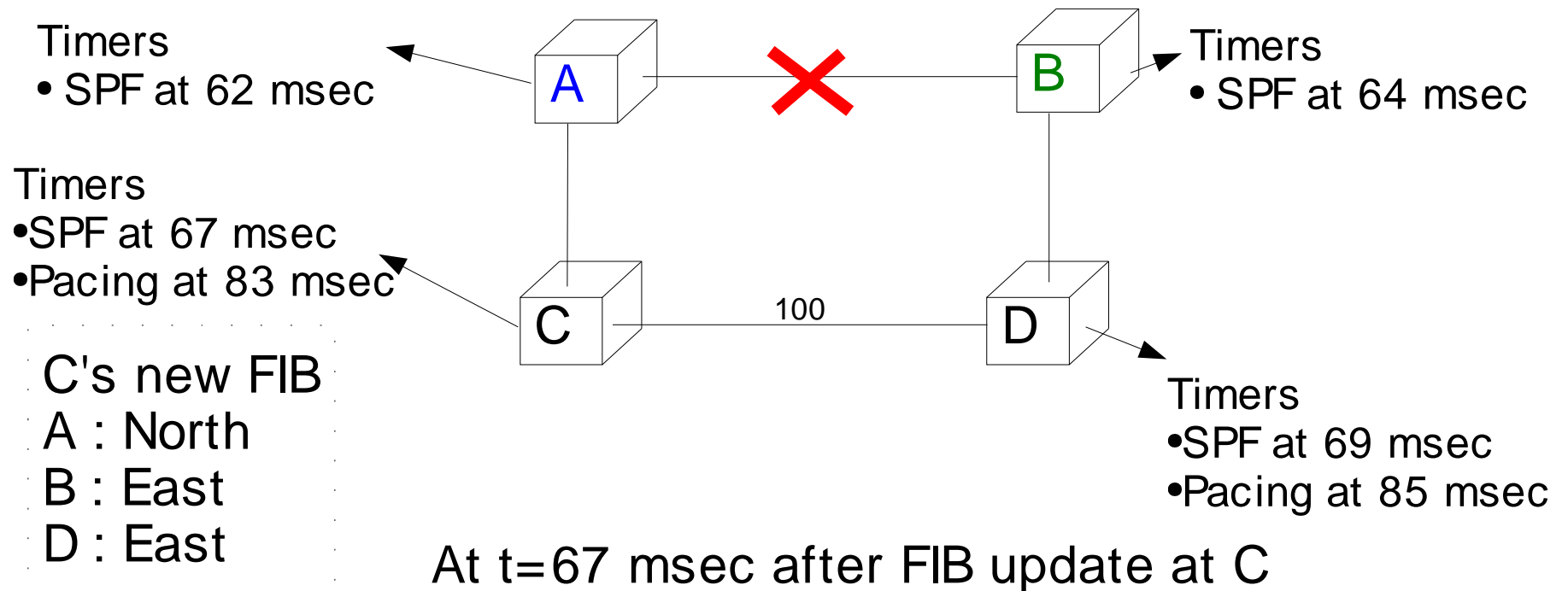
- Nettester at t=62 msec : 4 paths out of 12
 - A can reach C but not B and D loop (A-C)
 - B can reach D, but not A and C
 - C can reach A but not B and D loop (C-A)
 - D can reach B, but not A and C

Example link failure (8)



- At $t=64$ msec after FIB update at B
- Nettester at $t=64$ msec : 4 paths out of 12
 - A can reach C but not B and D (loop A-C)
 - B can reach D, but not A and C (loop B-D)
 - C can reach A but not B and D (loop C-A)
 - D can reach B, but not A and C (loop D-B)

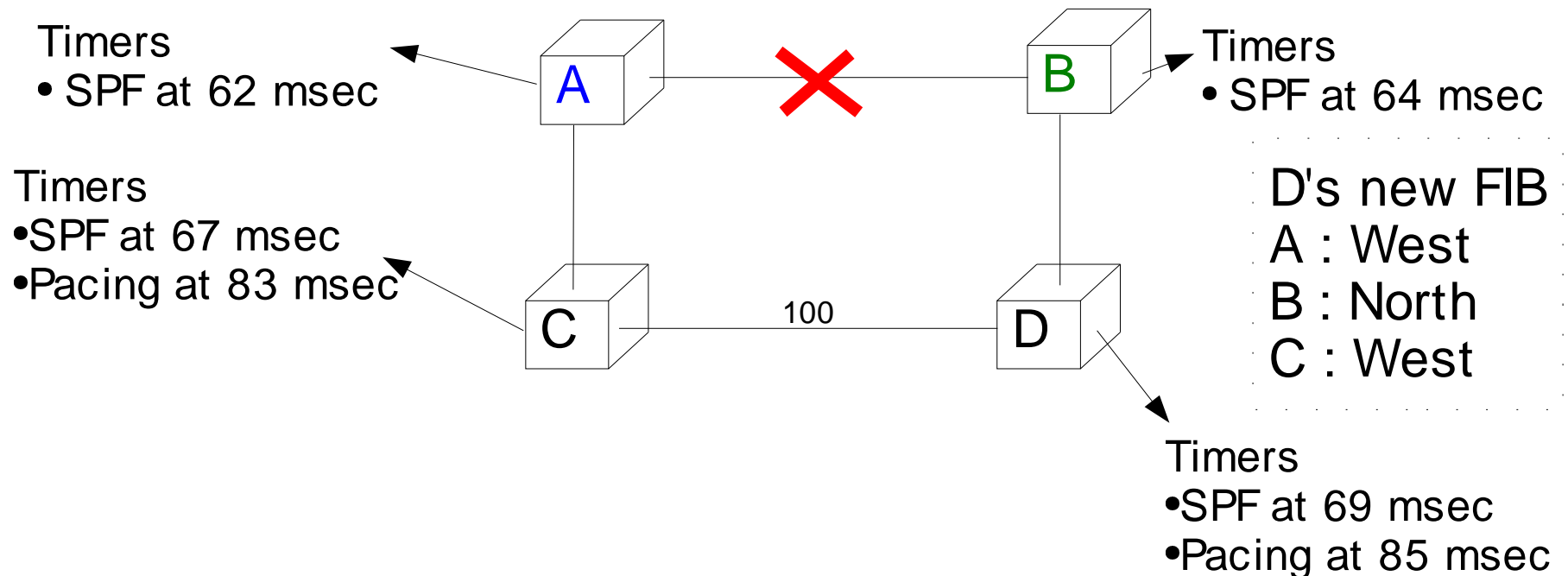
Example link failure (9)



At $t=67$ msec after FIB update at C

- Nettester at $t=67$ msec : 8 paths out of 12
 - A can reach B, C and D
 - B can reach D, but not A and C (loop B-D)
 - C can reach A, B and D
 - D can reach B, but not A and C (loop D-B)

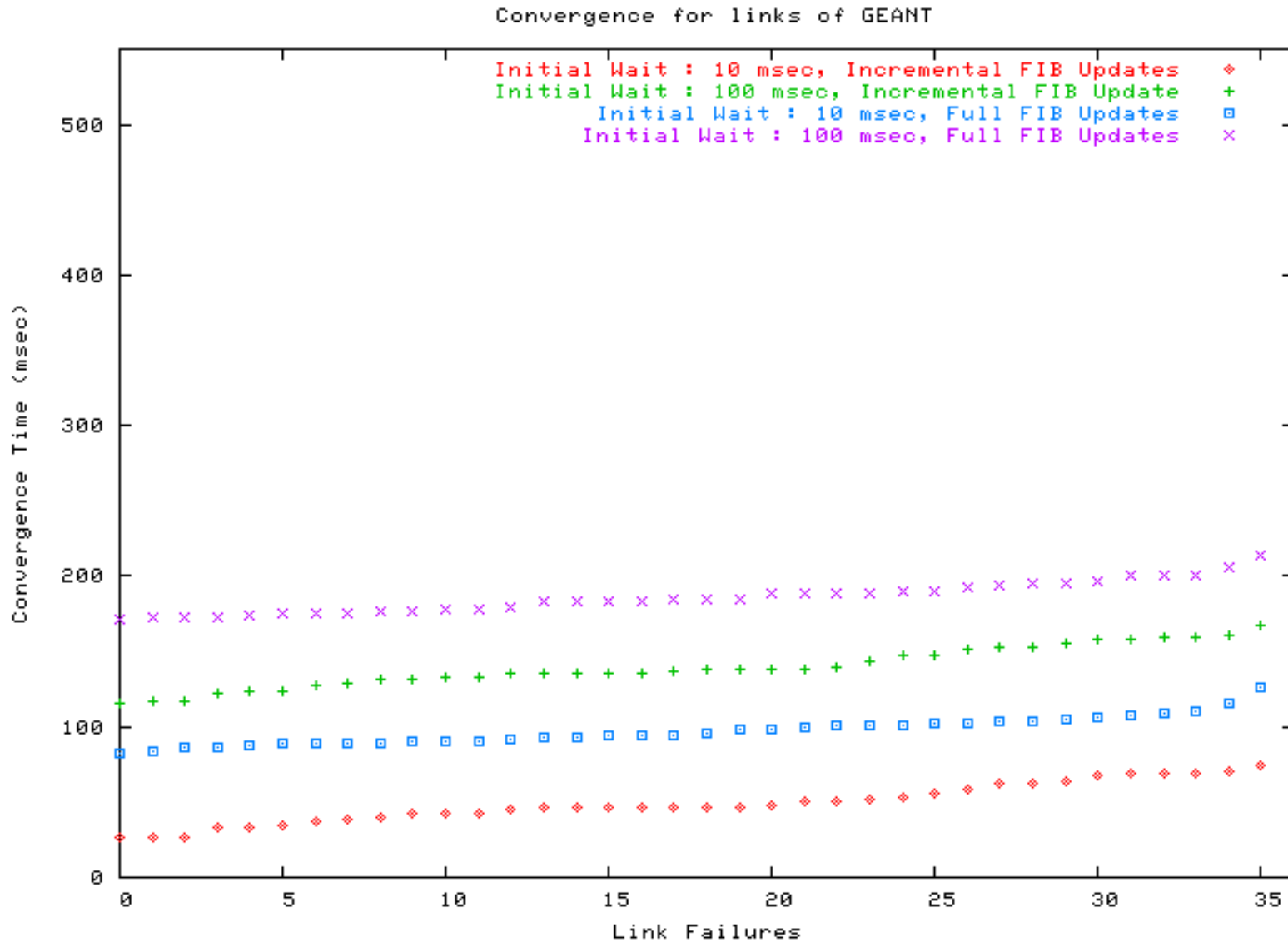
Example link failure (10)



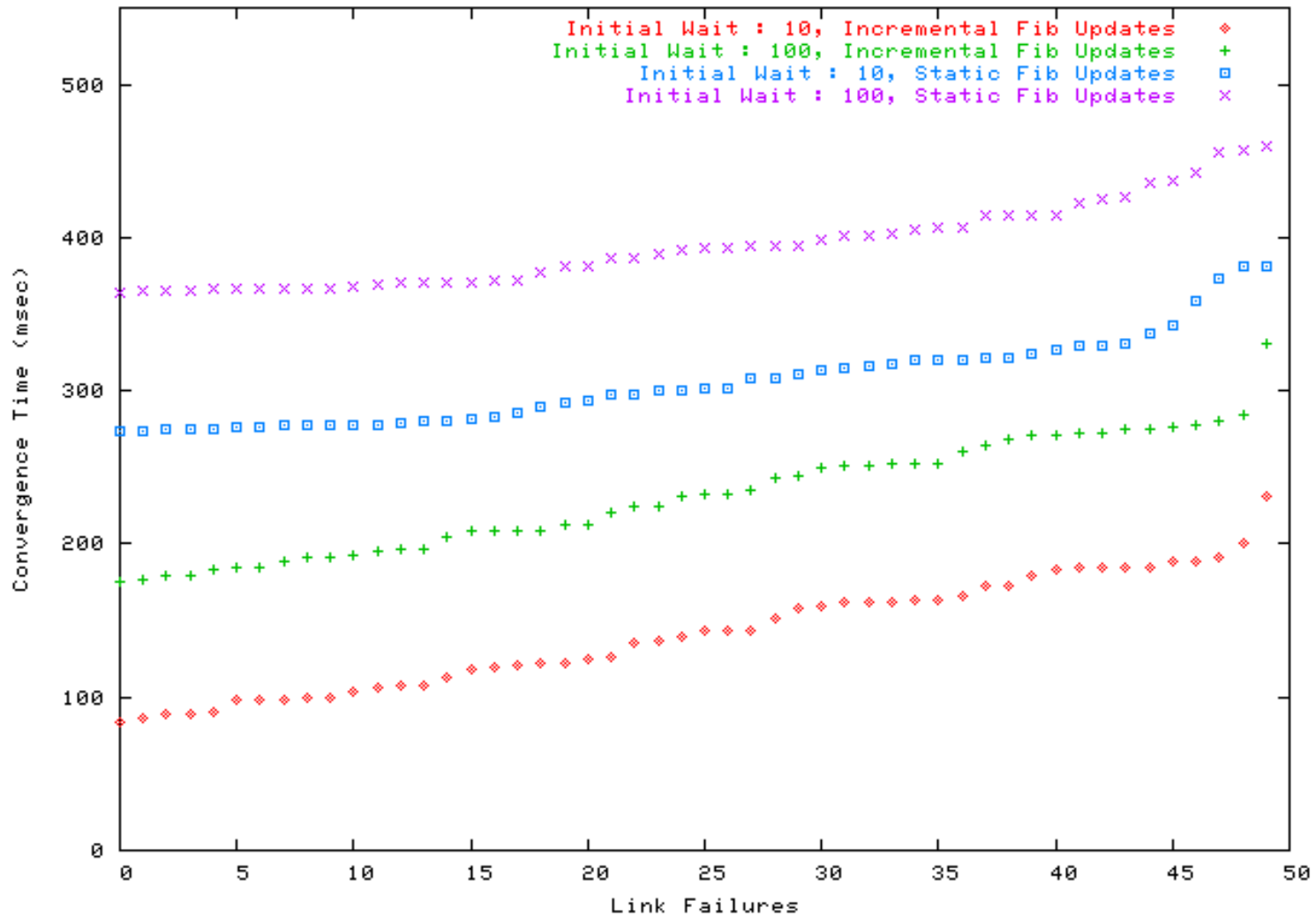
At t=69 msec after FIB update at C

- Nettester at t=69 msec : 12 paths out of 12
 - A can reach B, C and D
 - B can reach A, C and D
 - C can reach A, B and D
 - D can reach A, B and C

All link failures in GEANT

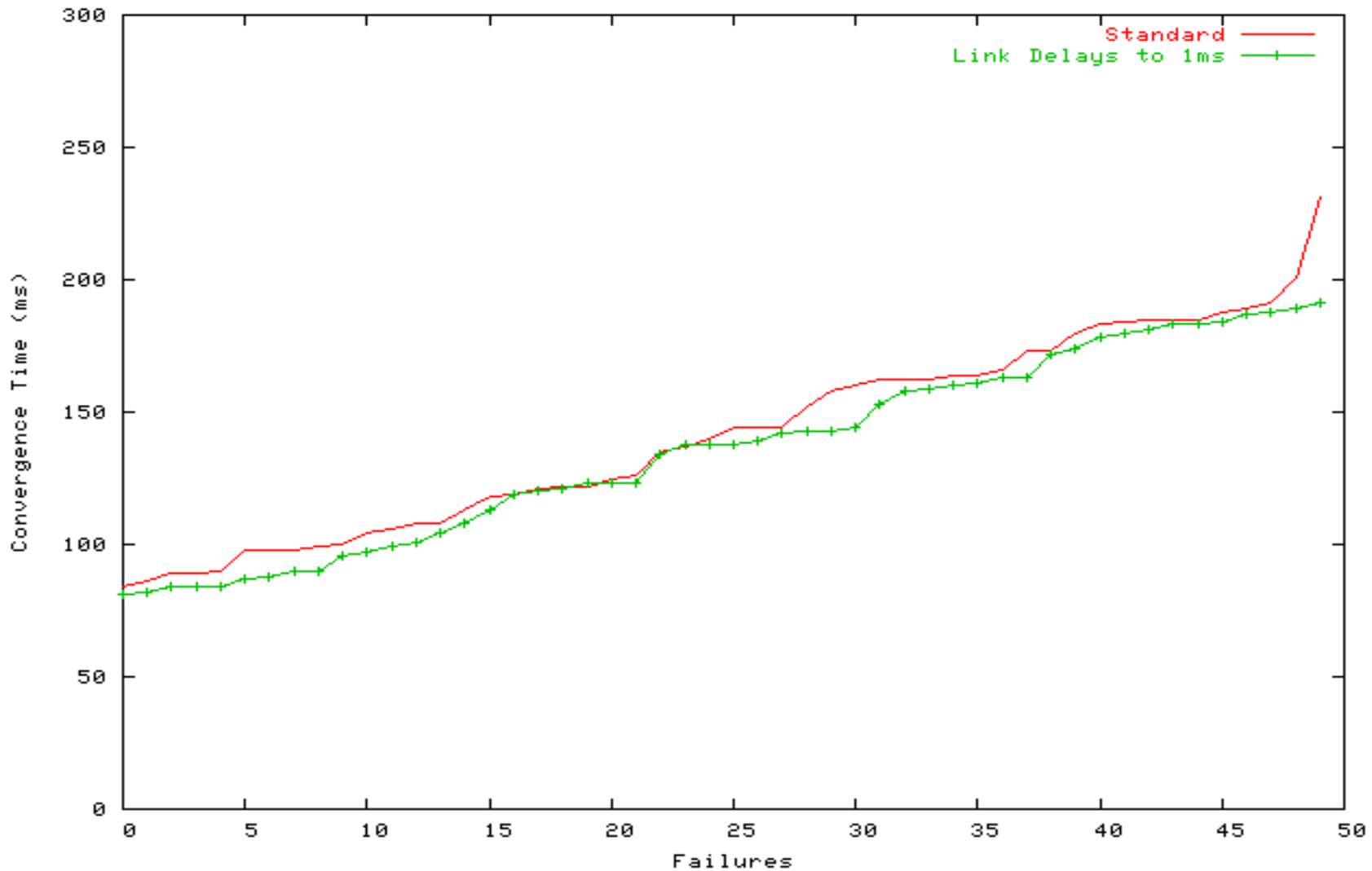


50 link failures in Tier-1 ISP



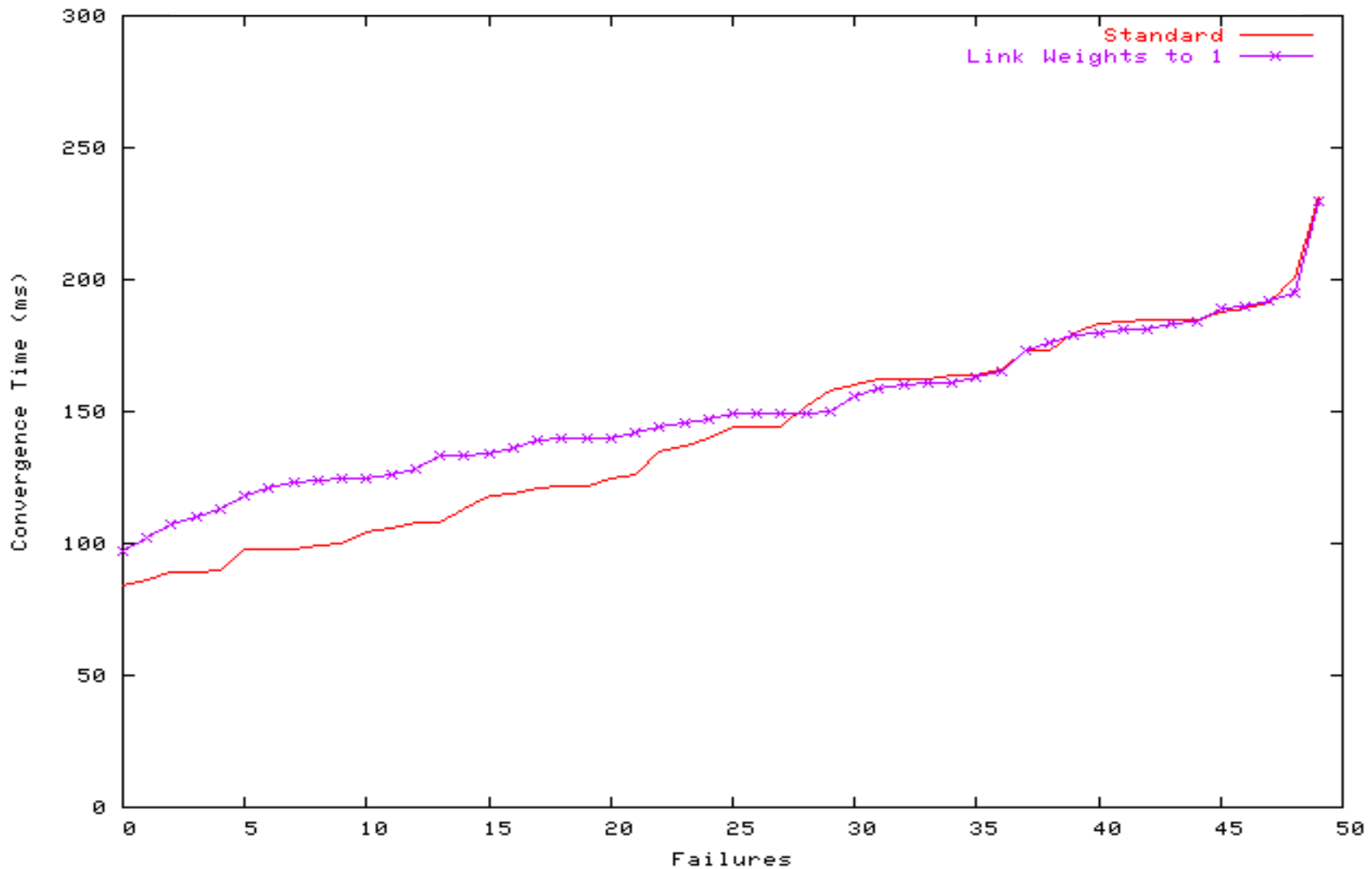
50 link failures in Tier-1 ISP

Impact of link delays



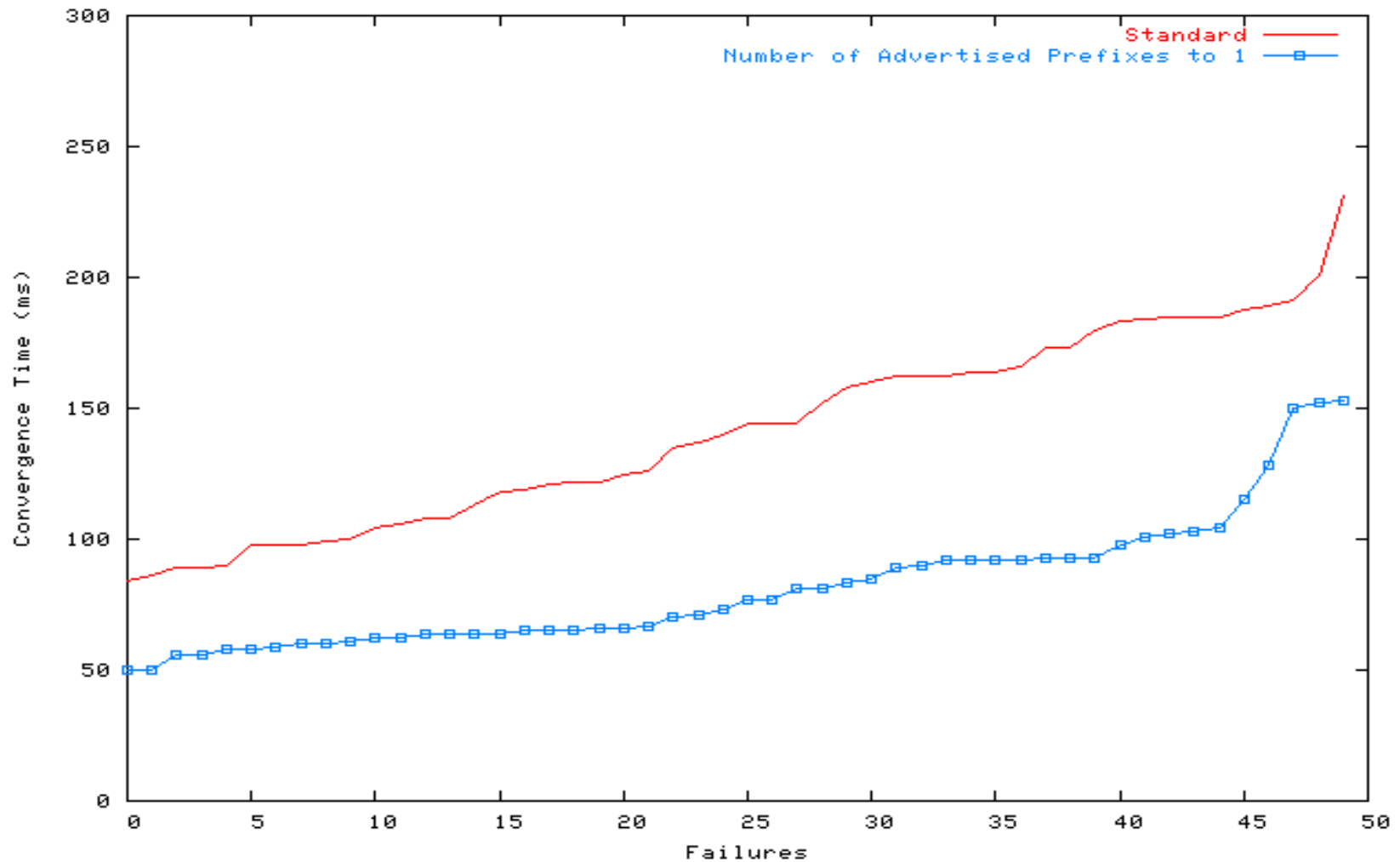
50 link failures in Tier-1 ISP

Impact of ISIS weights



50 link failures in Tier-1 ISP

Impact of number of prefixes



Recommendations for link failures

- Initial wait
 - Should be as small as possible to improve convergence in case of link failures
 - ◆ 70% of the failures are link failures in Sprint
- FIB size
 - A small FIB size is important to ensure fast convergence
 - Reducing the number of prefixes advertised by the IGP reduces convergence time
- IGP weights
 - Should be set to reroute as locally as possible

Agenda

- Behaviour of IS-IS in ISP networks

- **Simulation study**
 - Simulation Model
 - Analysis of link failures
 - ● **Analysis of router failures**

- Towards sub 50 msec failure recovery

Router failures

- Used router failures as a way to model SRLG failures
 - Few SRLG information is available for the GBLX and GEANT topologies
 - Detecting SRLG information from IGP traces is difficult
- What happens when a router fails ?
 - all its links fail and its neighbours detect the link failure within 10-15 msec
 - All neighbours flood their new LSP

Convergence time for router failures

- Modification to Nettester

- Definition

- ◆ Routing is consistent for a pair S-D at time t if all the paths that packets would follow, from S to D, following the FIB of the routers at time t, are loop-free and end at D, without passing through the failed node

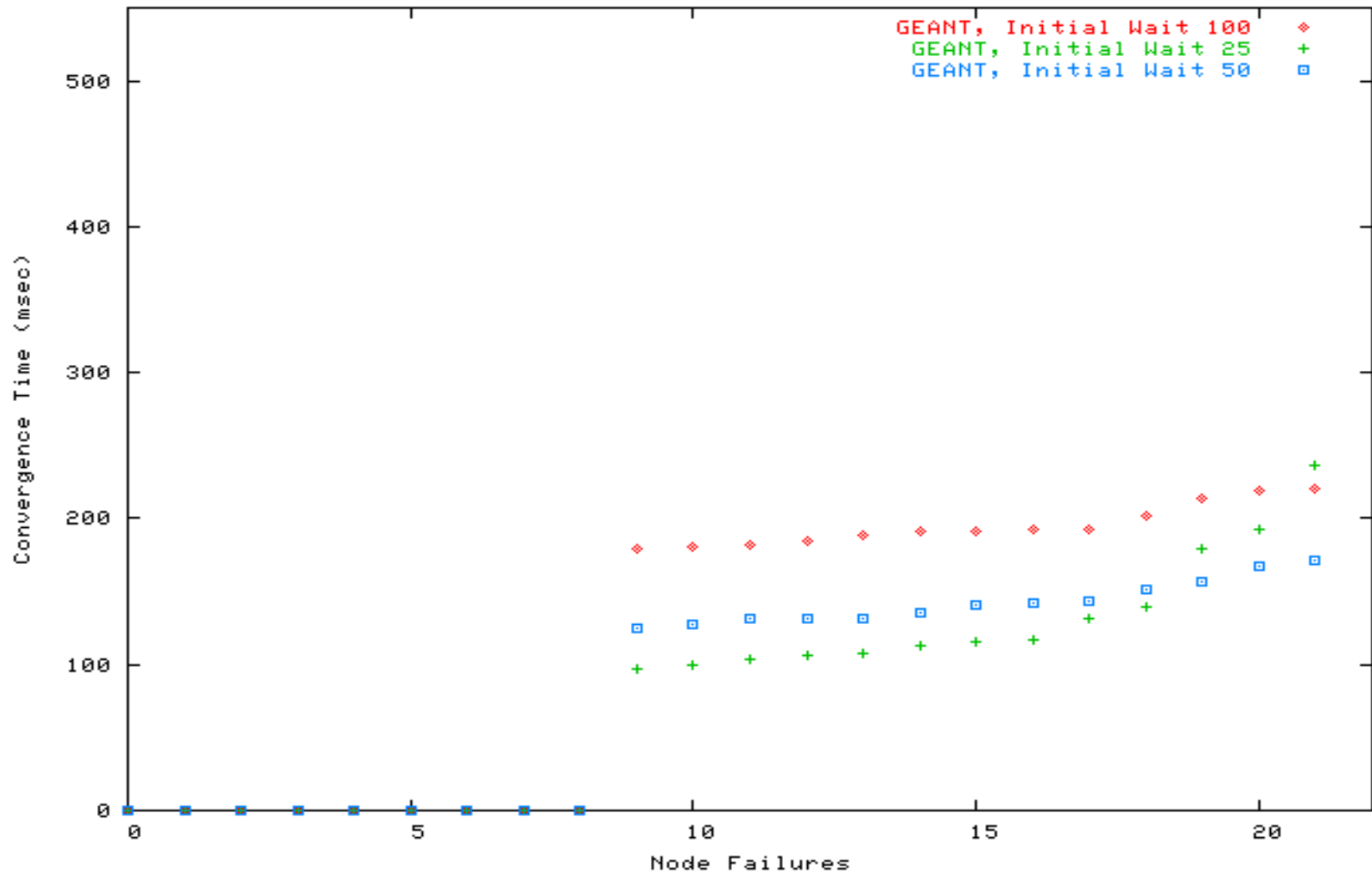
- Principle

- ◆ Before the failure, routing is consistent
 - ◆ Convergence time is the time when routing becomes and remains consistent for all router-router pairs (excluding the failed router)
 - ◆ Consistency is checked by using the loopback addresses of the routers as sources and destinations

All router failures in GEANT

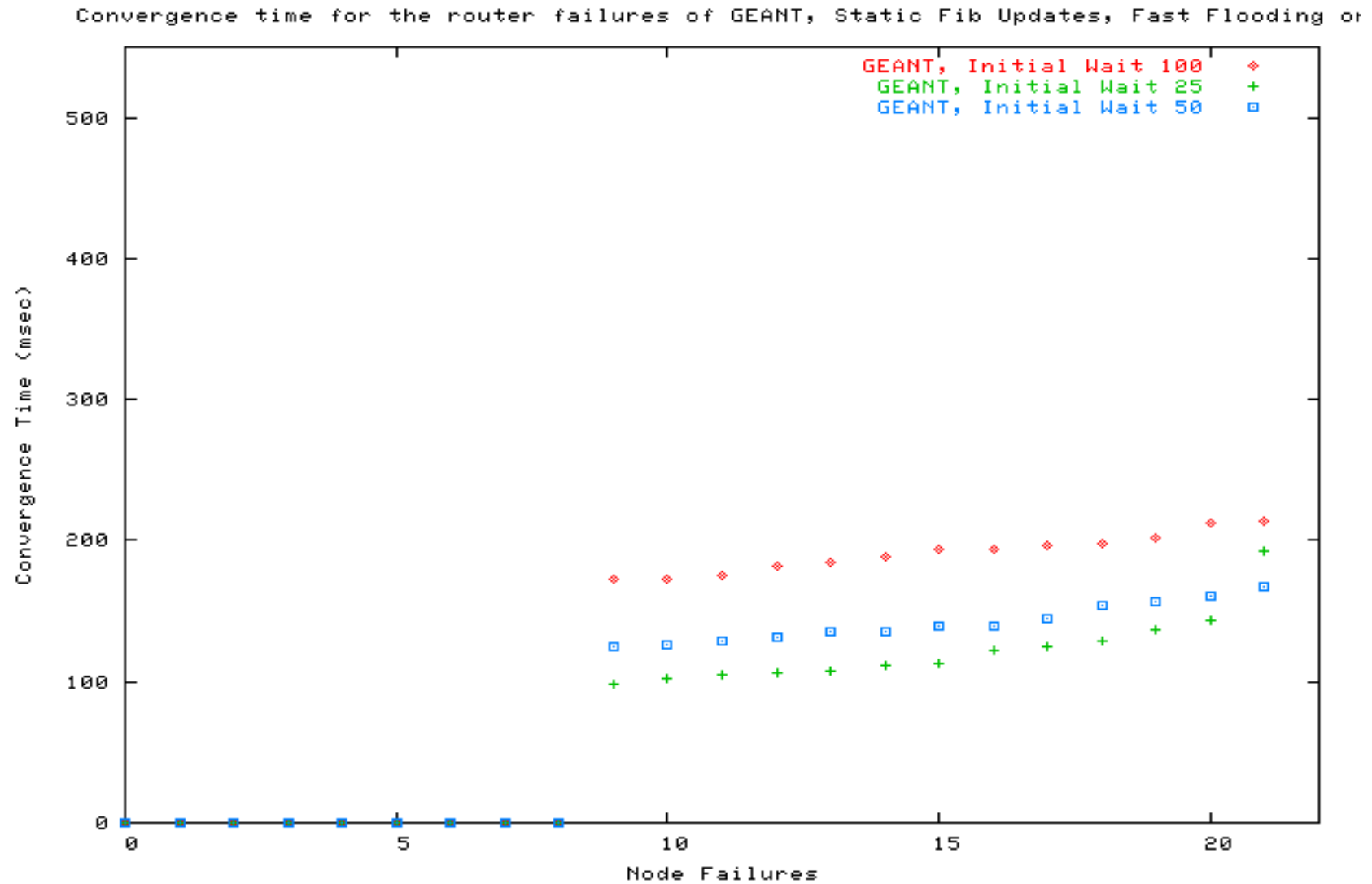
Static FIB updates, 33 msec pacing

Convergence time for the router failures of GEANT, Static Fib Updates, Fast Flooding off, Pacin



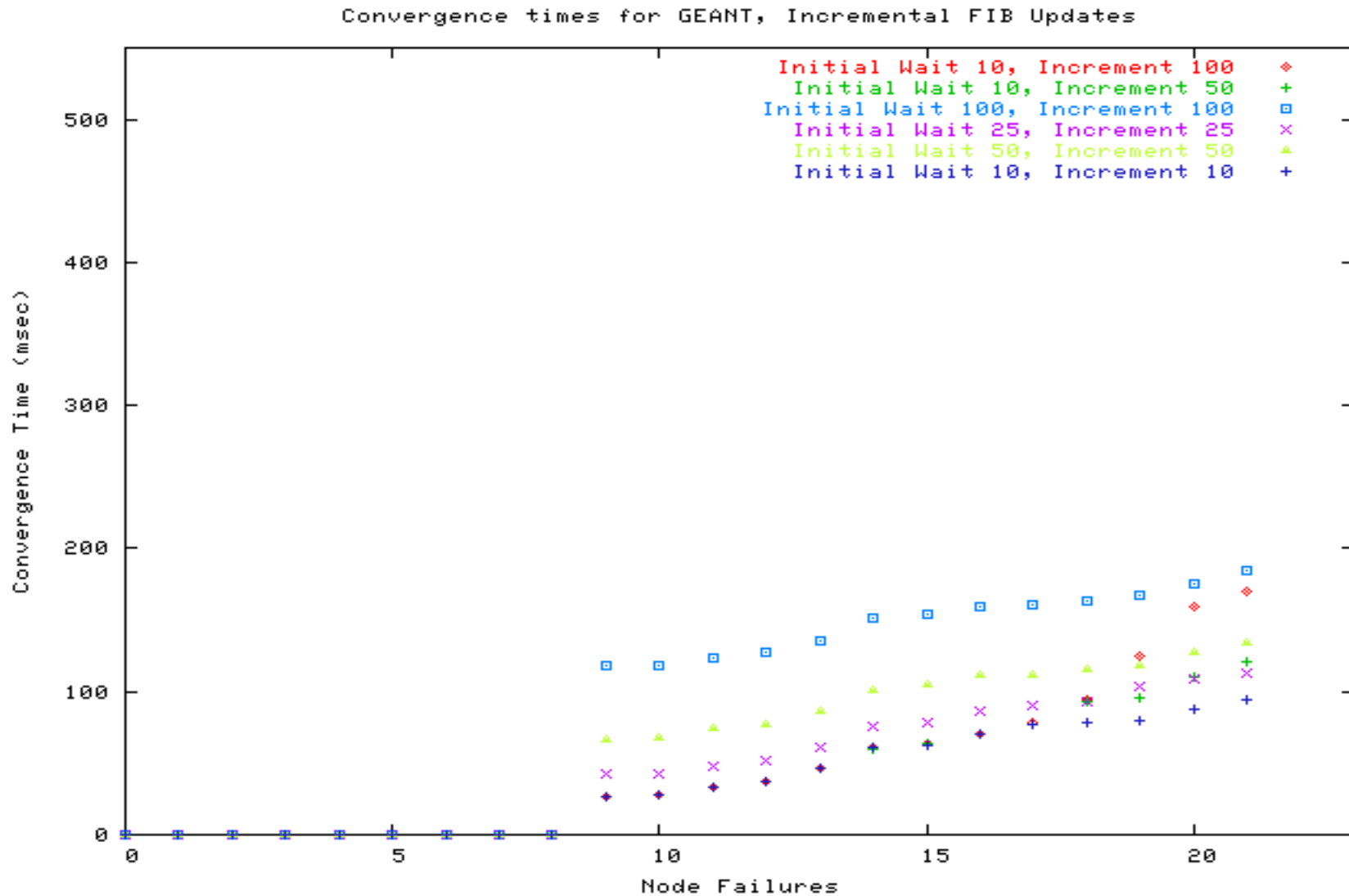
All router failures in GEANT

Static FIB updates, fast flooding



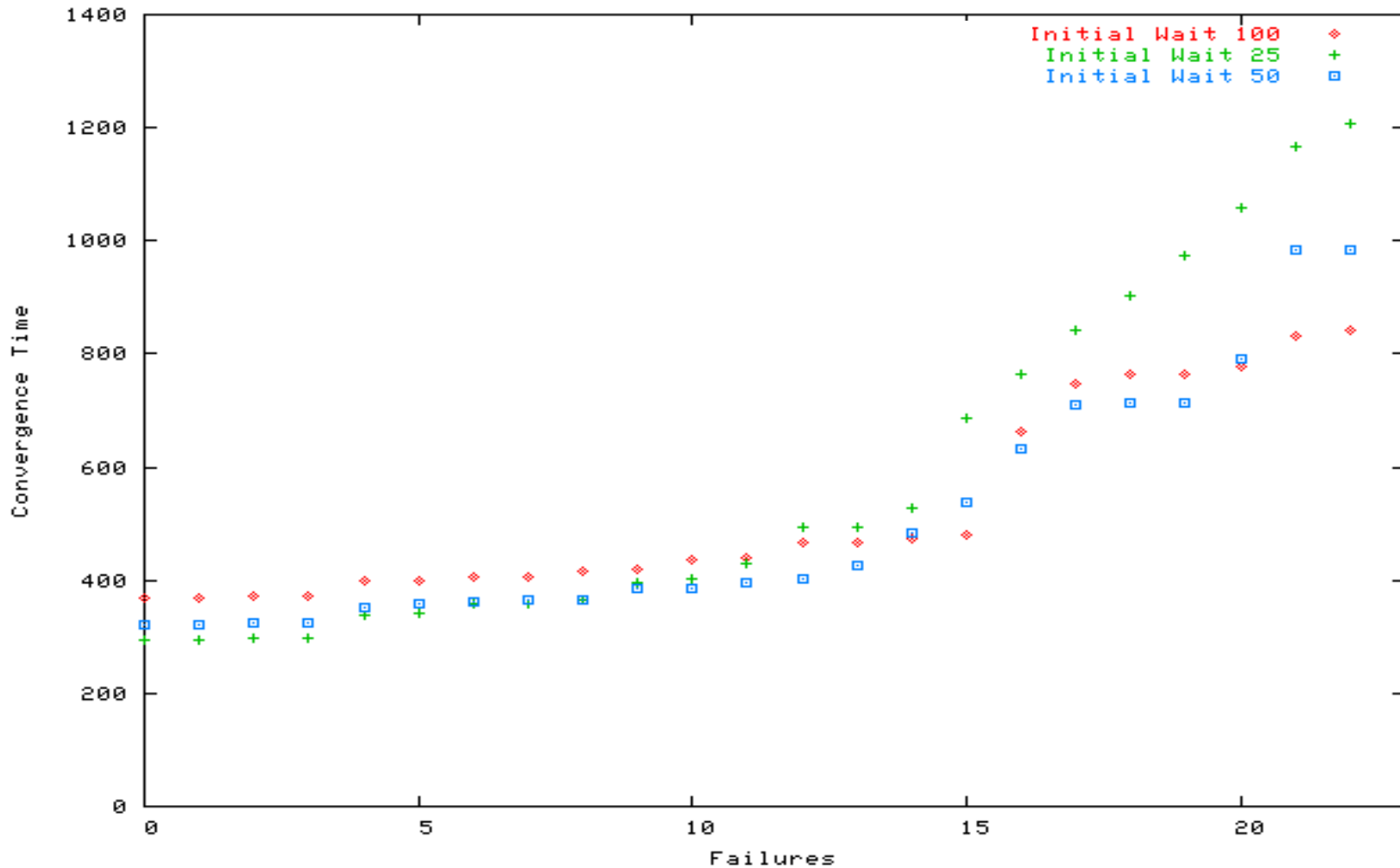
All router failures in GEANT

Incremental FIB updates, fast flooding



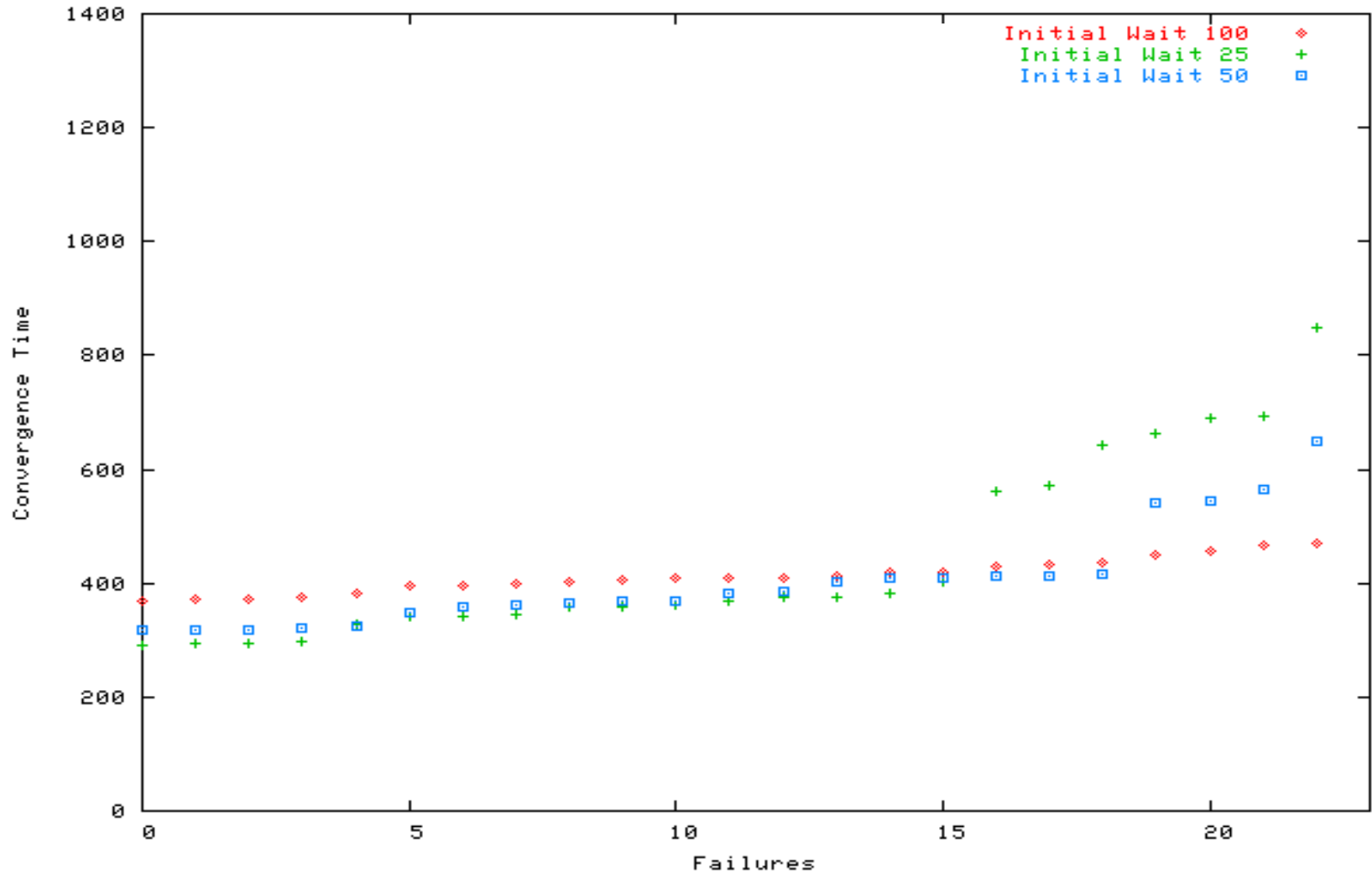
23 router failures in Tier-1 ISP

Static FIB updates, 33 msec pacing



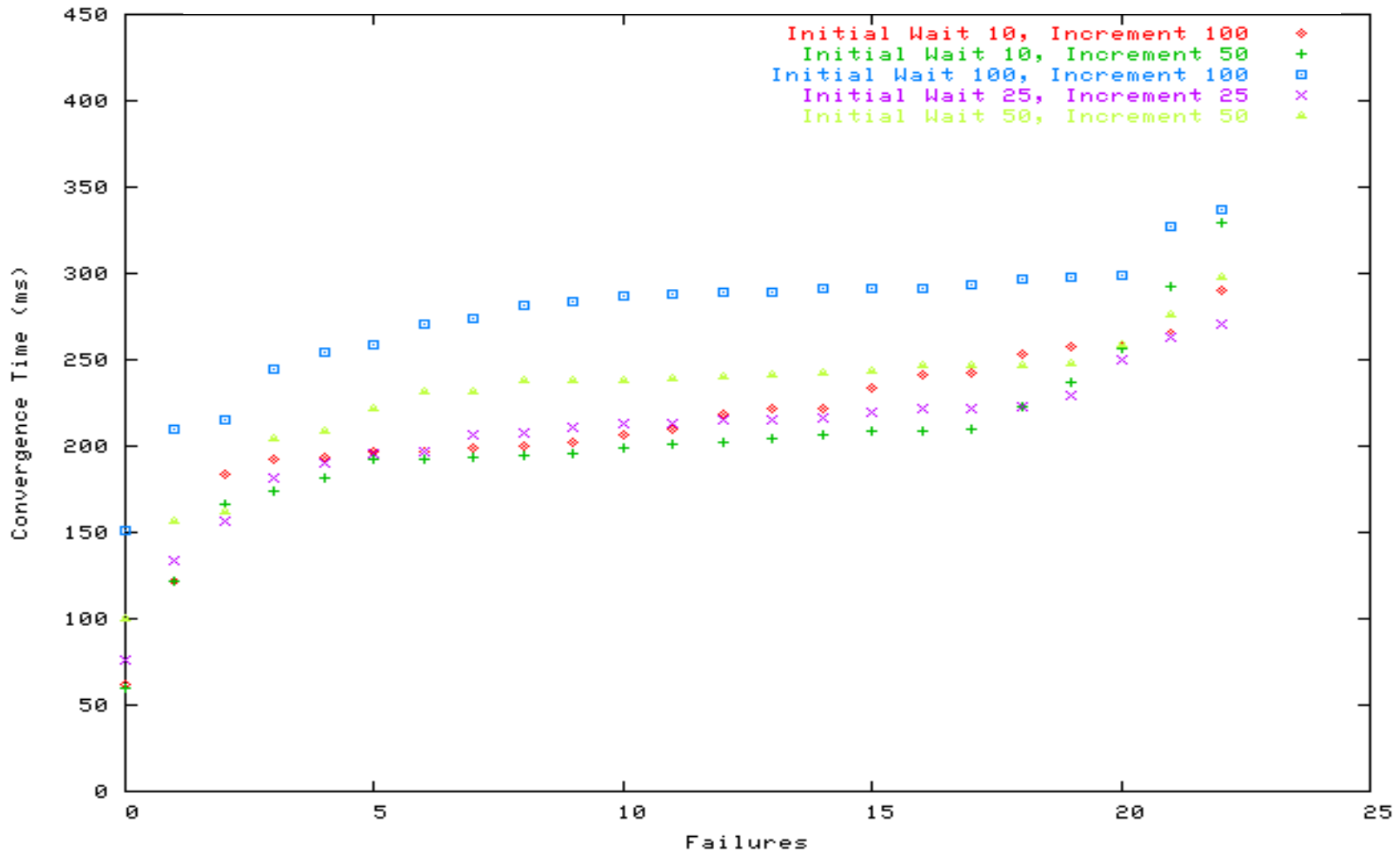
23 router failures in Tier-1 ISP

Static FIB updates, fast flooding



23 router failures in Tier-1 ISP

Incremental FIB updates, fast-flooding



Recommendations for router failures

- Fast flooding
 - Required for fast convergence
 - ◆ allows most LSPs to be flooded before running SPF+FIB
 - ◆ Isolates flooding of urgent LSPs from ISIS noise
- SPF Initial wait
 - Should be large enough to ensure that all important LSPs have been received before running SPF+FIB
- FIB size
 - Reducing the number of prefixes advertised by the IGP reduces convergence time

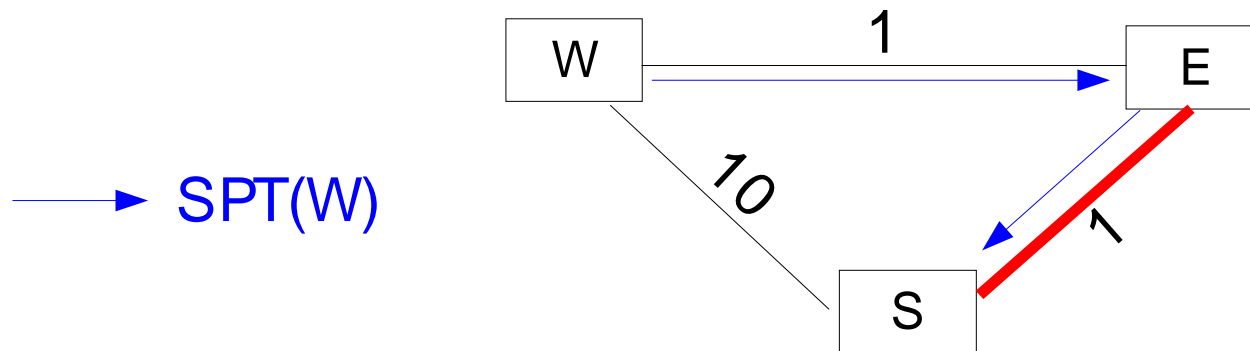
Agenda

- Behaviour of IS-IS in ISP networks
- Simulation study
- ● Towards sub 50 msec failure recovery

How to provide sub 50 msec recovery in pure IP networks ?

- First step
 - When a (directed) link fails, immediately reroute the packets *at the router that detects the failure to a loop-free alternate router*
 - ◆ *This loop-free alternate router is precomputed*
- What is a loop-free alternate router ?
 - For the failure of link S->E and destination D, this is a neighbour N, whose shortest path to reach D does not contain S->E

Loop-free neighbor



- Loop-free neighbour detection algorithm for protected link S->E
 - For each direct neighbour (S-> N_i)
 - ◆ Compute SPT(N_i)
 - ◆ if (S->E) SPT(N_i)
 - ◆ then N_i is a candidate loop-free neighbour for all destinations
 - ◆ otherwise not

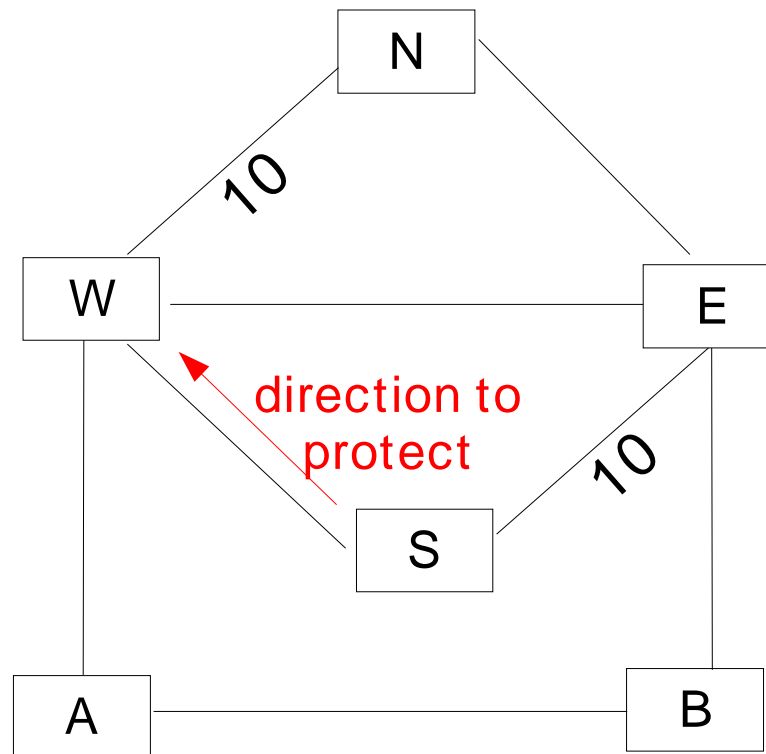
Loop-free neighbours

- Example

- all links have weight=1 except NW and SE

W's routing table

N : East via E
S : SouthEast
E : East
A : South
B : South via A
East via E



E's routing table

N : NorthWest
S : SouthWest
W : West
B : South
A : South via B
West via W

S's routing table

All : via W

- If S->W fails, E is a loop-free neighbor
 - ◆ all S->W's packets sent to E will not loop

Loop-free neighbours (2)

W's routing table

N : East via E

S : SouthEast

E : East

A : South

B : South via A

East via E

A's routing table

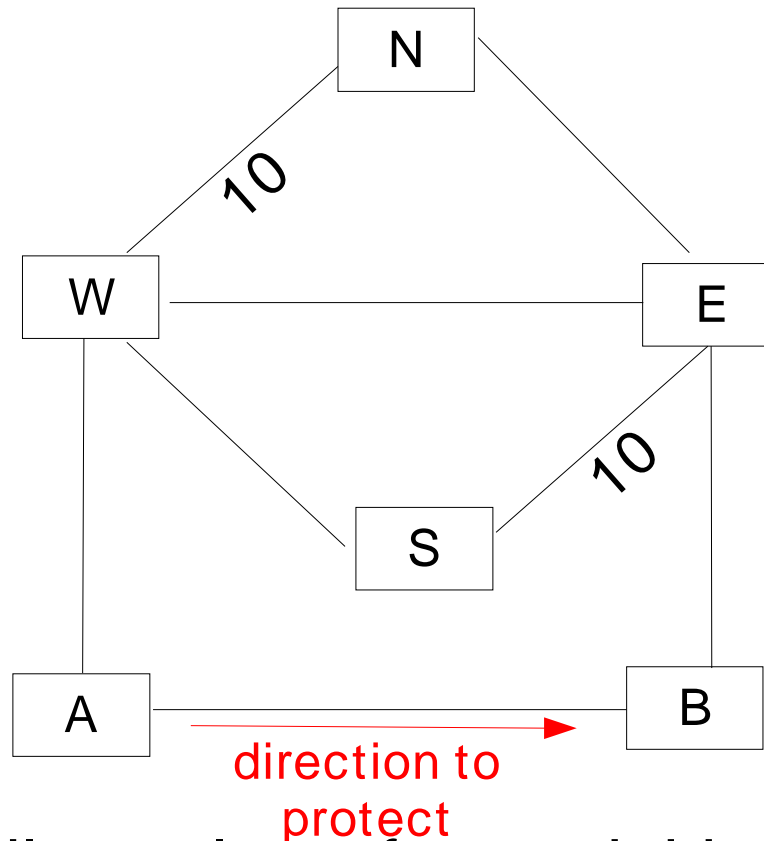
N,S : North via W

W : North

B : East

E : North via W

East via B



B's routing table

E : North

A : West

N,S : North via E

W : West via A

North via E

- If A->B fails, no loop-free neighbour for destination B

- ◆ W would return to A the packets towards node B

Evaluation of loop-free neighbours

- Question

- How many links can be quickly protected by immediately switching over all the traffic that they carry to a loop-free neighbour ?

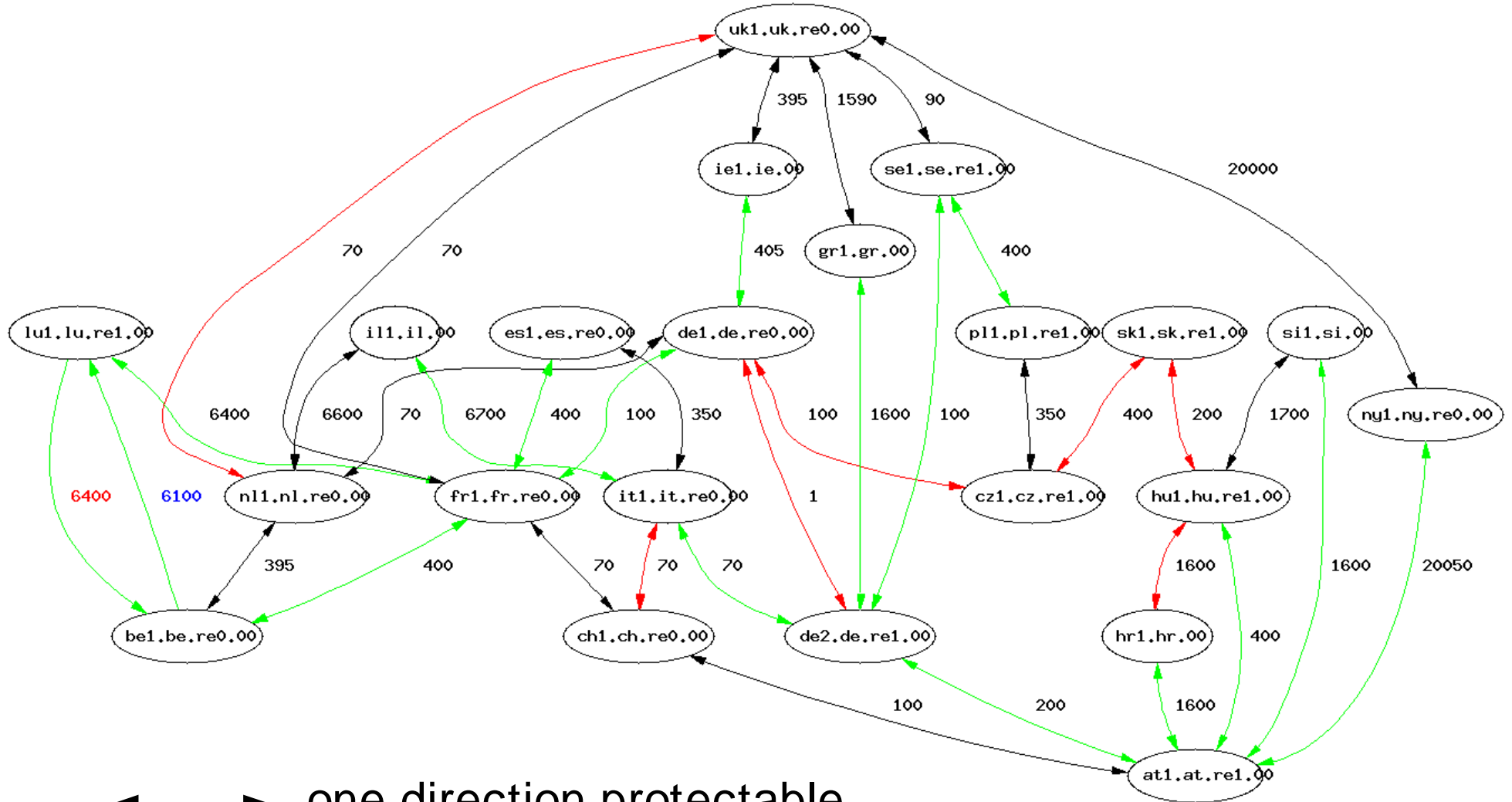
- Algorithm

- for each directed link A->B carrying packets
 - ◆ compute *dlist*, the list of destinations reachable via this directed link
 - ◆ compute the amount of traffic carried on this link
 - ◆ for all neighbours N of A except B
 - ◆ check whether N can reach all destinations inside *dlist* without using link A->B
 - ◆ if yes, N can be used to protect directed link A->B
 - ◆ if no, N is not a valid candidate loop-free neighbour

Loop-free neighbours in GEANT

- Total traffic : 4024 units
 - based on real traffic matrix
- Protectable traffic with loop-free neighbours
 - 1859 (46%)
- Number of directed links carrying traffic
 - 72
- Number of protectable directed links carrying traffic
 - 48

The protectable links with loop-free neighbours in GEANT



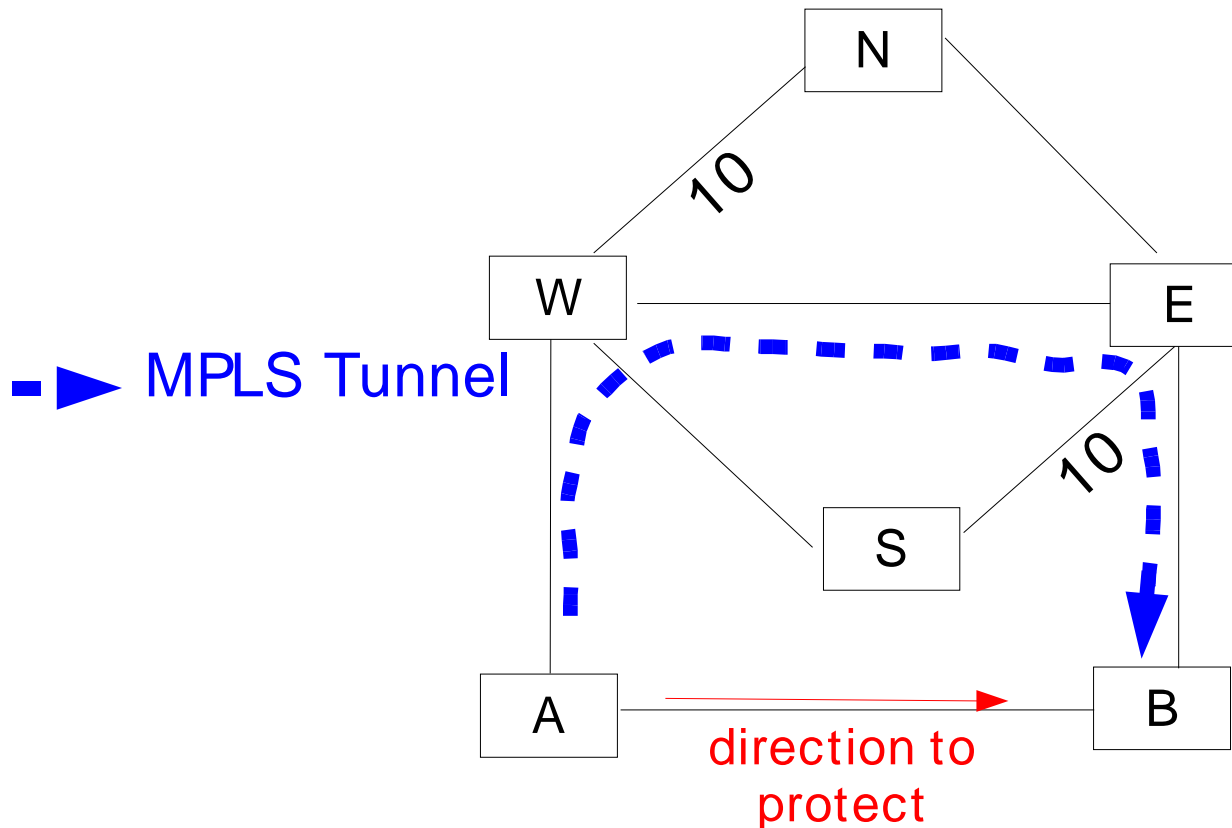
- ↔ one direction protectable
- ↔ No direction protectable
- ↔ Both directions protectable

Loop-free neighbours in a Tier-1 ISP

- Total traffic : 216459 units
 - based on real traffic matrix
- Protectable traffic : 166482 (76.9 %)
 - 84.9% of the intrapop traffic is protectable
 - 70.9% of the interpop traffic is protectable
- Directed links carrying traffic : 756
 - 358 intrapop links (out of 486) are protectable
 - 187 interpop links (out of 270) are protectable

Loop-free alternate routers

- How to improve the coverage ?
 - Use non-neighbours as alternate routers
 - Simple solution
 - ◆ MPLS tunnel to protect failed link



U-turns

- Principle

- If there is no loop-free neighbour, a neighbour of our neighbours might be loop-free...
- When failure occurs, return the packets to sender who will send them to its loop-free neighbour
 - ◆ Assumes hardware support on interfaces

W's routing table

N : East via E

S : SouthEast

E : East

A : South

B : South via A

East via E

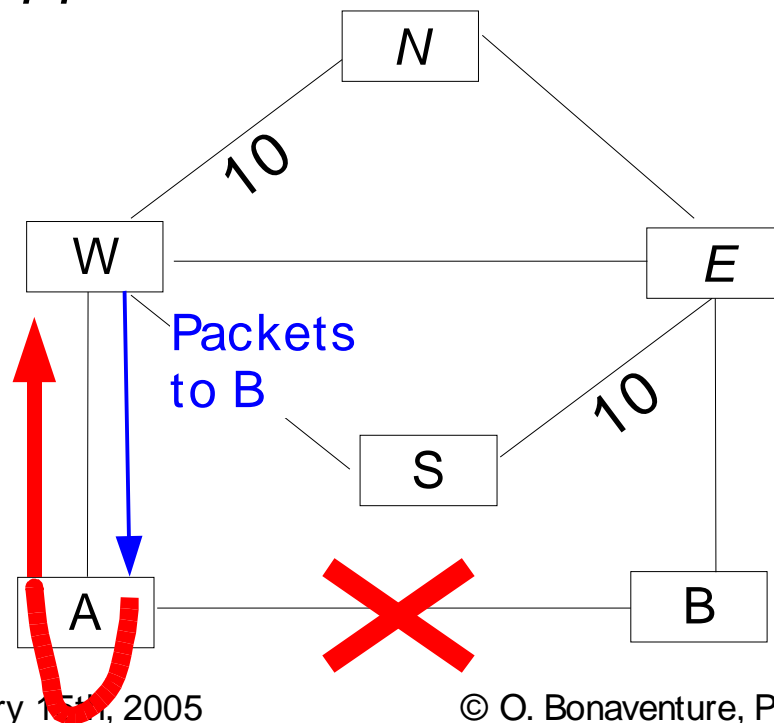
Loop-free neighbours

If W->A fails

Use E to reach B

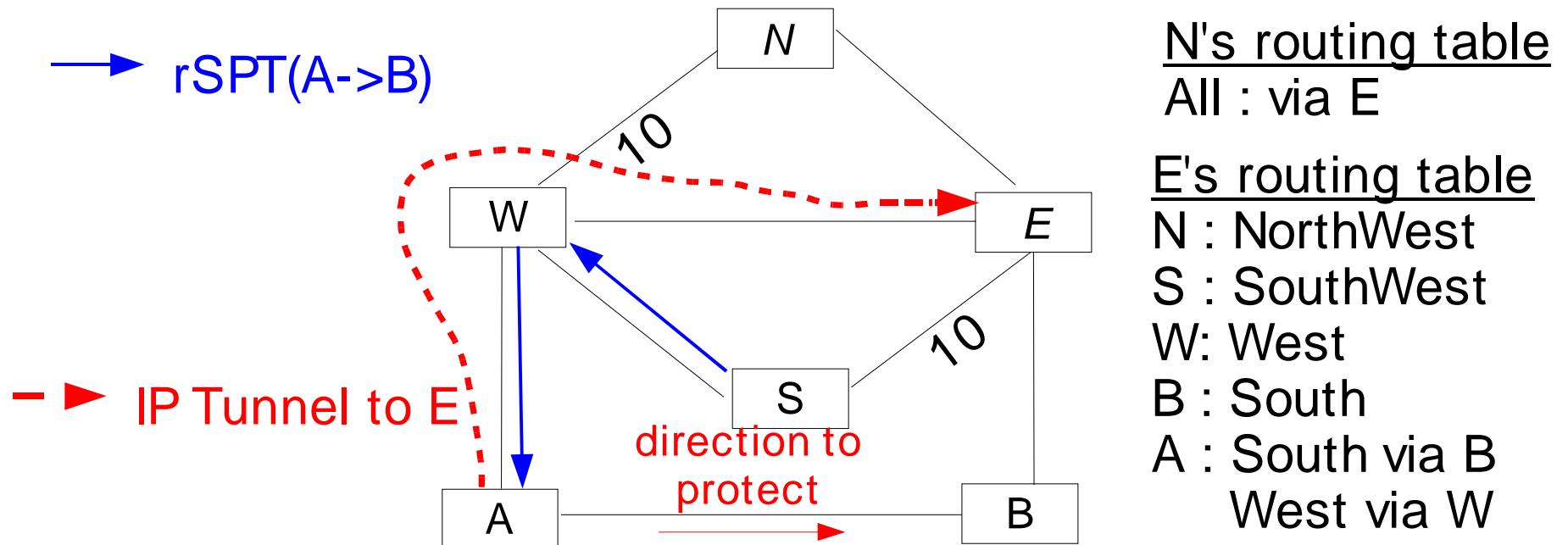
U-turned
to B

Packets
to B



Loop-free alternate routers (2)

- Another solution
 - Use as loop-free alternate a router that does not use the (directed) link to be protected



- Precompute a tunnel towards loop-free alternate router to protect link from failure

Are loop-free alternates sufficient ?

- Consider the failure of link A->B
 - A immediately updates its FIB to use tunnel

A's routing table

N,S : North via W

W : North

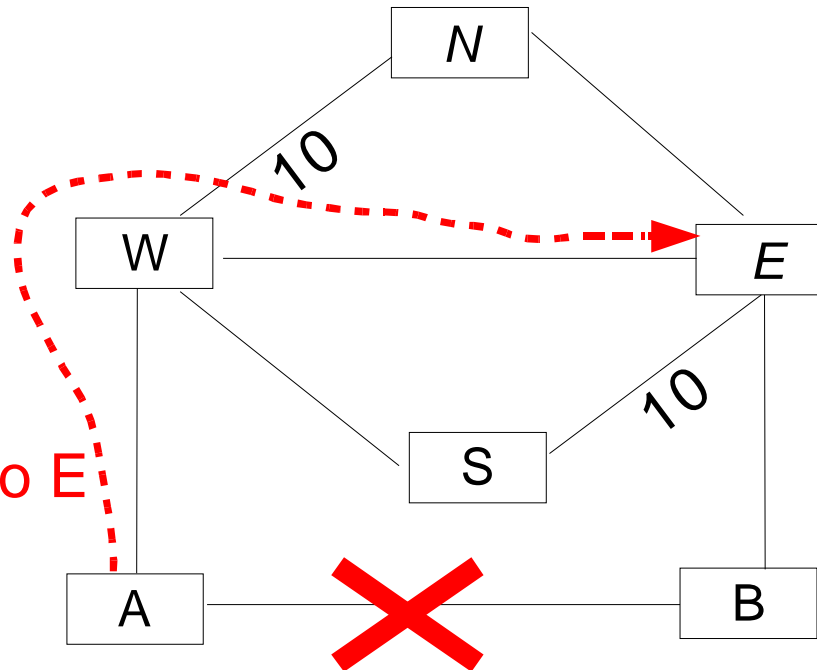
~~B : East~~

E : North via W

~~East via B~~

B : via tunnel to E

- ► IP Tunnel to E



- Is this sufficient to avoid all packet losses ?

Are loop-free alternates sufficient ? (2)

- Unfortunately, the protection tunnel is not optimal
- A will flood its new link-state packet and all routers will eventually update their FIB

A's routing table

N,S : North via W

W : North

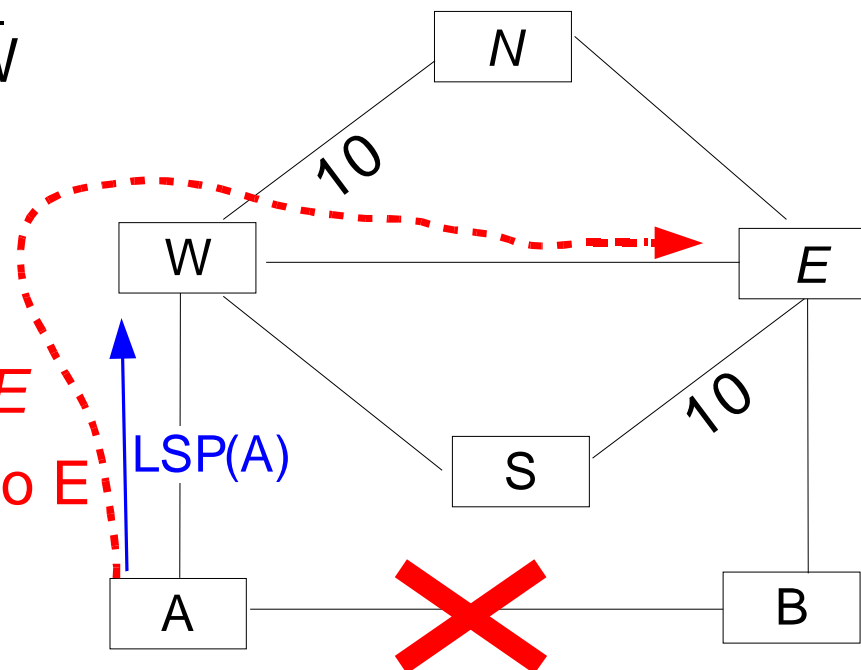
~~B : East~~

E : North via W

~~East via B~~

B : via tunnel to E

- ► IP Tunnel to E



Are loop-free alternates sufficient ? (3)

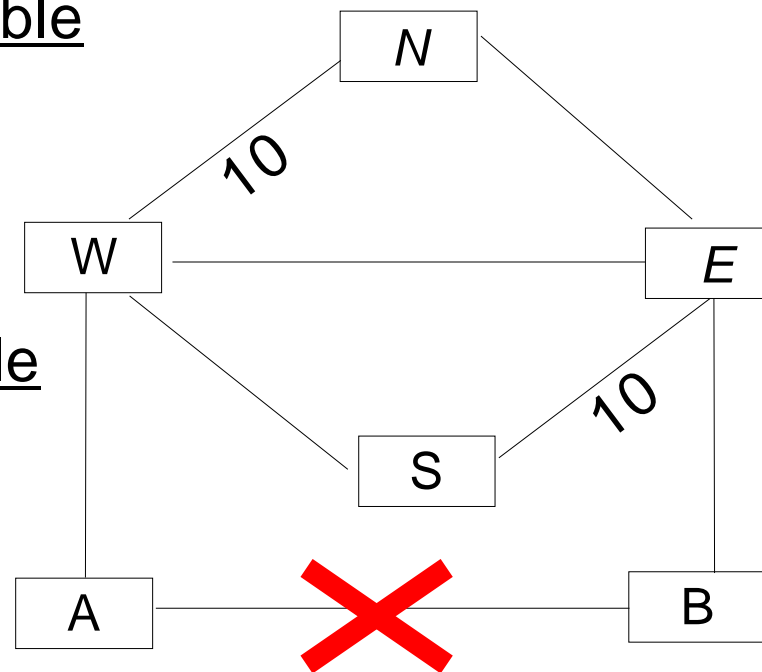
- W updates its FIB before A
 - Everything is fine, no packets are lost
- A updates its FIB before W
 - Packets towards B loop between A and W

A's new routing table

...
E : North via W
B : North via W

A's old routing table

...
E : North via W
B : via tunnel to E



W's old routing table

...
E : East
B : South via A
East via W

W's new routing table

...
E : East
B : East via E

How to avoid transient loops during FIB updates ?

- Three solutions are discussed within IETF
 - Synchronised update of all the FIBs
 - Timer-based ordering the updates of the FIBs
 - Distributed ordering of the updates of the FIBs
- Solution developed could also handle all non-urgent topology changes
 - link brought up/down for maintenance
 - router reboot
 - change in link weights

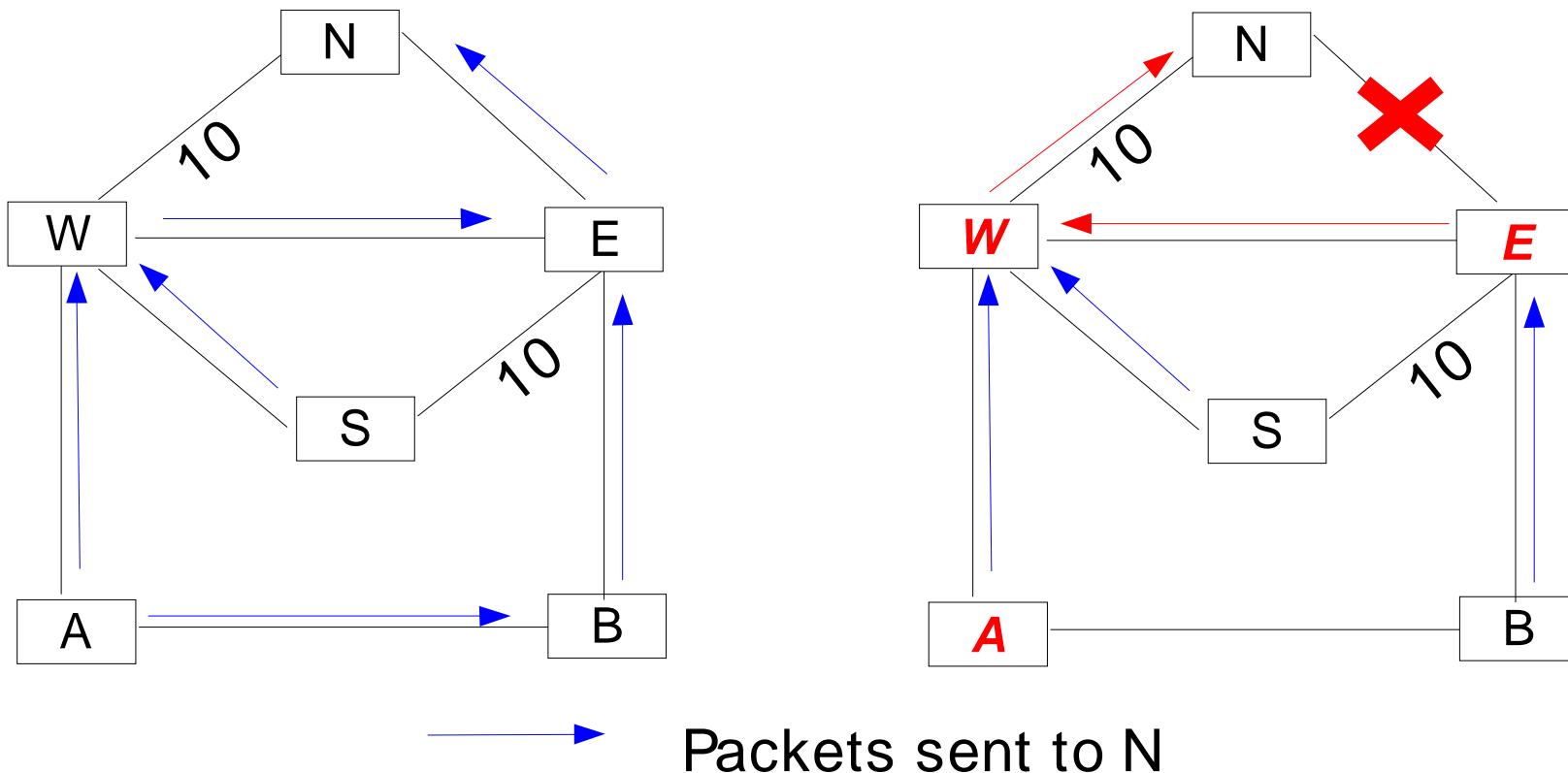
Synchronised FIB updates

- Principle
 - To avoid transient forwarding loops during the updates of the FIBs, ensure that all FIBs are updated **exactly at the same time**
 - ◆ update time can be included inside LSP
- Difficulties
 - Routers need to be synchronised
 - ◆ GPS clocks, NTP
 - Router must be able to update their FIB quickly
 - ◆ A possibility is to have two FIB copies on the line card and switch FIB, but FIBs use expensive memory
 - How to make sure that this technique works on low-end as well as high-end routers ?

Timer-based ordering of FIB updates

- Principle

- When a link fails, routers far away from the failure must update their FIB before routers close to the failure



Timer-based ordering of FIB updates (2)

- Algorithm used by router R receiving a LSP indicating a non-urgent failure of link X->Y
- Check if X->Y belongs to router's SPT
 - ◆ if not, FIB is already up-to-date
 - ◆ Because router R is not using link X->Y
 - ◆ if yes, R's FIB must be updated
 - ◆ Compute RSPT centered on Y
 - ◆ The RSPT is computed by considering the network topology before the failure of link X->Y
 - ◆ Find N, farthest (in hops) node upstream of R inside RSPT(Y)
 - ◆ Timer at router R is $\text{Flood} + T * \text{distance}_{\text{hops}}(R, N)$
 - ◆ Flood is the expected flooding time inside network
 - ◆ T should be larger than SPF + FIB computation time
 - ◆ Timer expiration
 - ◆ Update FIB

Timer-based ordering of FIB updates (3)

- Example computation of timers

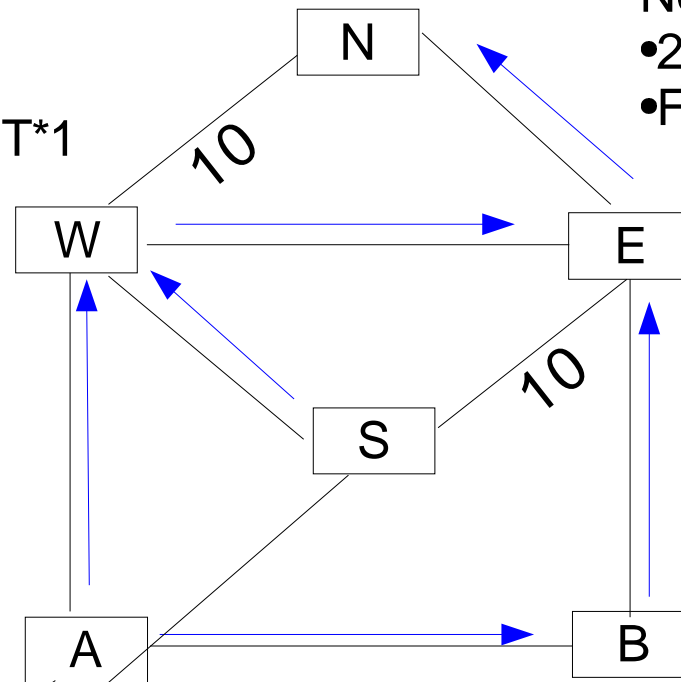
→ RSPT centered on N before failure

Node W :

- 1 hop from A
- FIB updated at : $\text{Now} + T * 1$

Node E :

- 2 hops from A
- FIB updated at : $\text{Now} + T * 2$



Node B :

- 1 hop from A
- FIB updated at : $\text{Now} + T * 1$

Farthest node from failure, FIB updated at : $\text{Now} + T * 0$

Protocol-based ordering of FIB updates

- To avoid transient loops during IGP convergence
- Order the updates of the FIBs on the distant routers
 - ◆ a non urgent failure can be handled in a few seconds if required, fast convergence in this case is not required
- ensure that a router will only update its FIB when it knows that it will not create transient loops
 - ◆ ordering of the FIB updates is built by exchanging HELLO PDUs containing special TLVs between routers

HELLO extension for link changes

- Link-event TLV contains
 - FIB bit
 - LSPid of first router attached to link
 - LSPid of second router attached to link
 - old ISIS metric
 - new ISIS metric
- Role of the FIB bit for failure of link X->Y
 - Router A sends **FIB=1** to router B
 - ◆ Router A is not (anymore) using router B to reach X->Y
 - Router A sends **FIB=0** to router B
 - ◆ Router A is currently using router B to reach X->Y
 - ◆ This implies that router B receiving this message should wait before updating its FIB

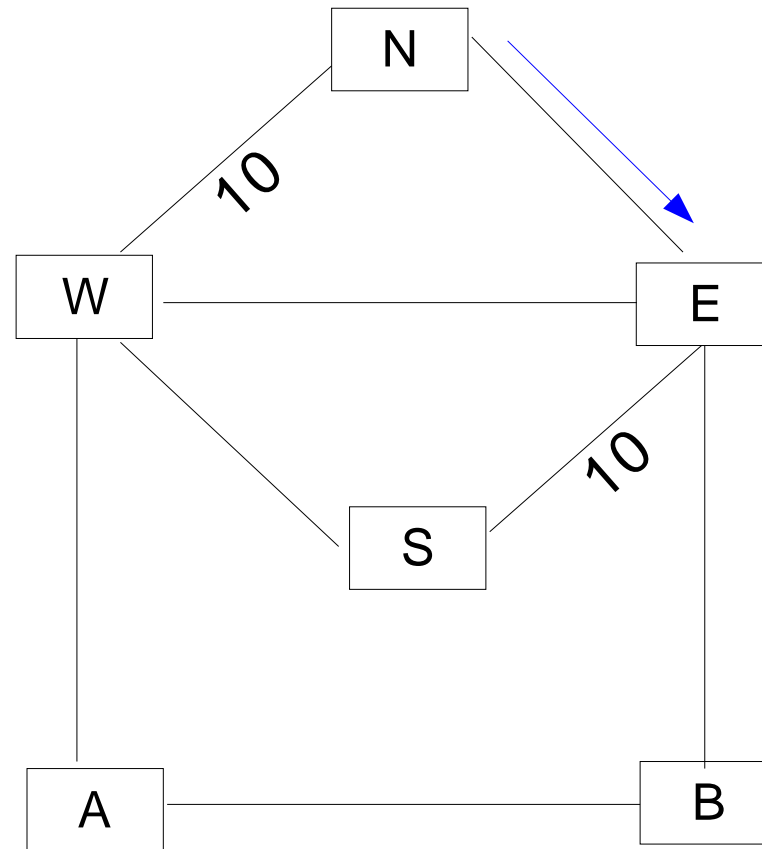
Router behaviour when link X->Y fails

- Arrival of LSP indicating failure of X->Y
 - LSP is flooded as usual
- Behaviour of Router R
 - if X->Y does **not** belong to SPT(R)
 - ◆ R is not using the failed link and will not update its FIB
 - ◆ If R receives a HELLO(X->Y) from a neighbour, it will reply by sending **HELLO(X->Y,FIB=1)**

Router behaviour when link X->Y fails (2)

- if X->Y **belongs** to SPT(R)
 - R is currently using the failed link and will update its FIB in an appropriate order
 - **W=neighbours(R)**
 - ◆ R must wait for a confirmation for all routers in W before updating its FIB
 - For all neighbours that R **uses** as nexthop to reach X
 - ◆ R sends **HELLO(X->Y,FIB=0)**
 - For all neighbours that R **does not use** as nexthop to reach X
 - ◆ R sends **HELLO(X->Y,FIB=1)**
 - R will only update its FIB once it has received
 - ◆ **HELLO(X->Y,FIB=1)** from all its neighbours
 - ◆ after the FIB update, R will send **HELLO(X->Y,FIB=1)** to all neighbours that it used to reach X before the failure of X->Y

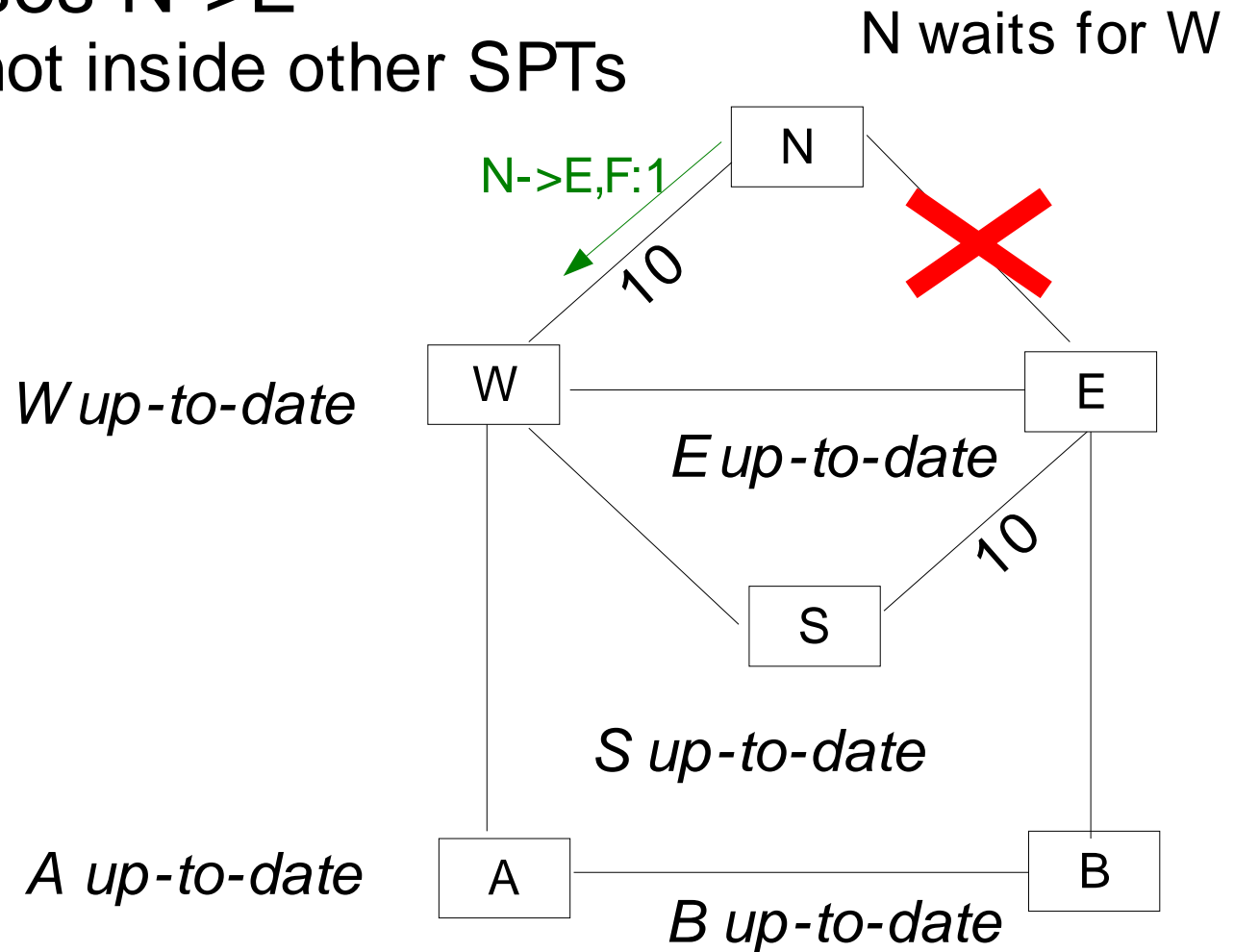
Users of link N->E



Packets sent to E

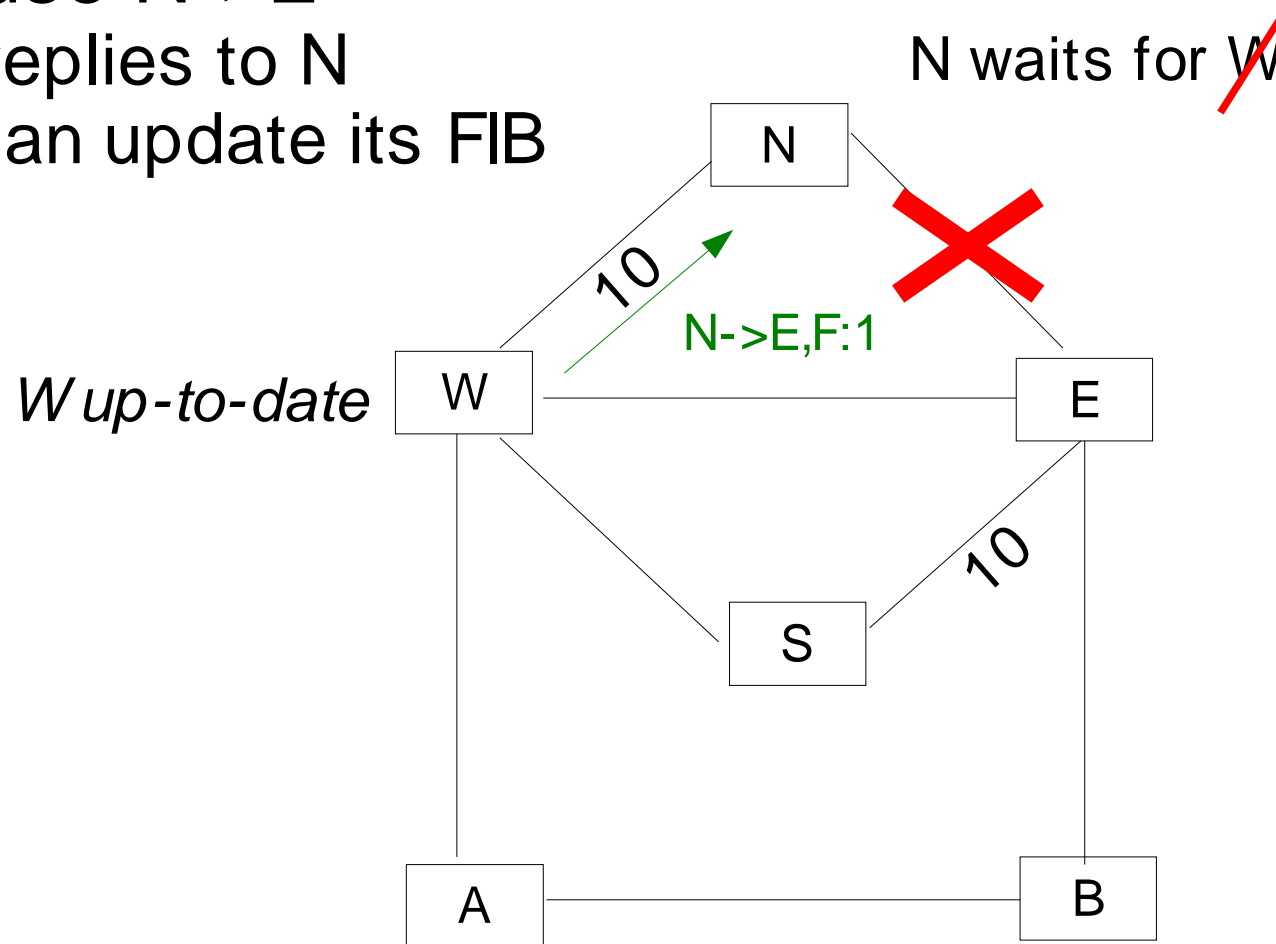
Graceful failure of link N->E

- LSP is flooded
- Only N uses N->E
 - N->E is not inside other SPTs

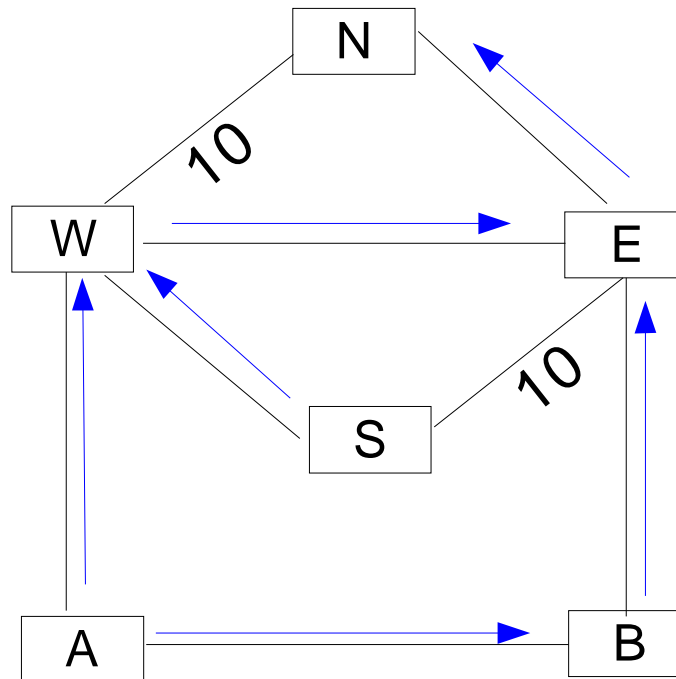


Graceful failure of link N->E (2)

- W's FIB is already up to date since it does not use N->E
 - W replies to N
 - N can update its FIB



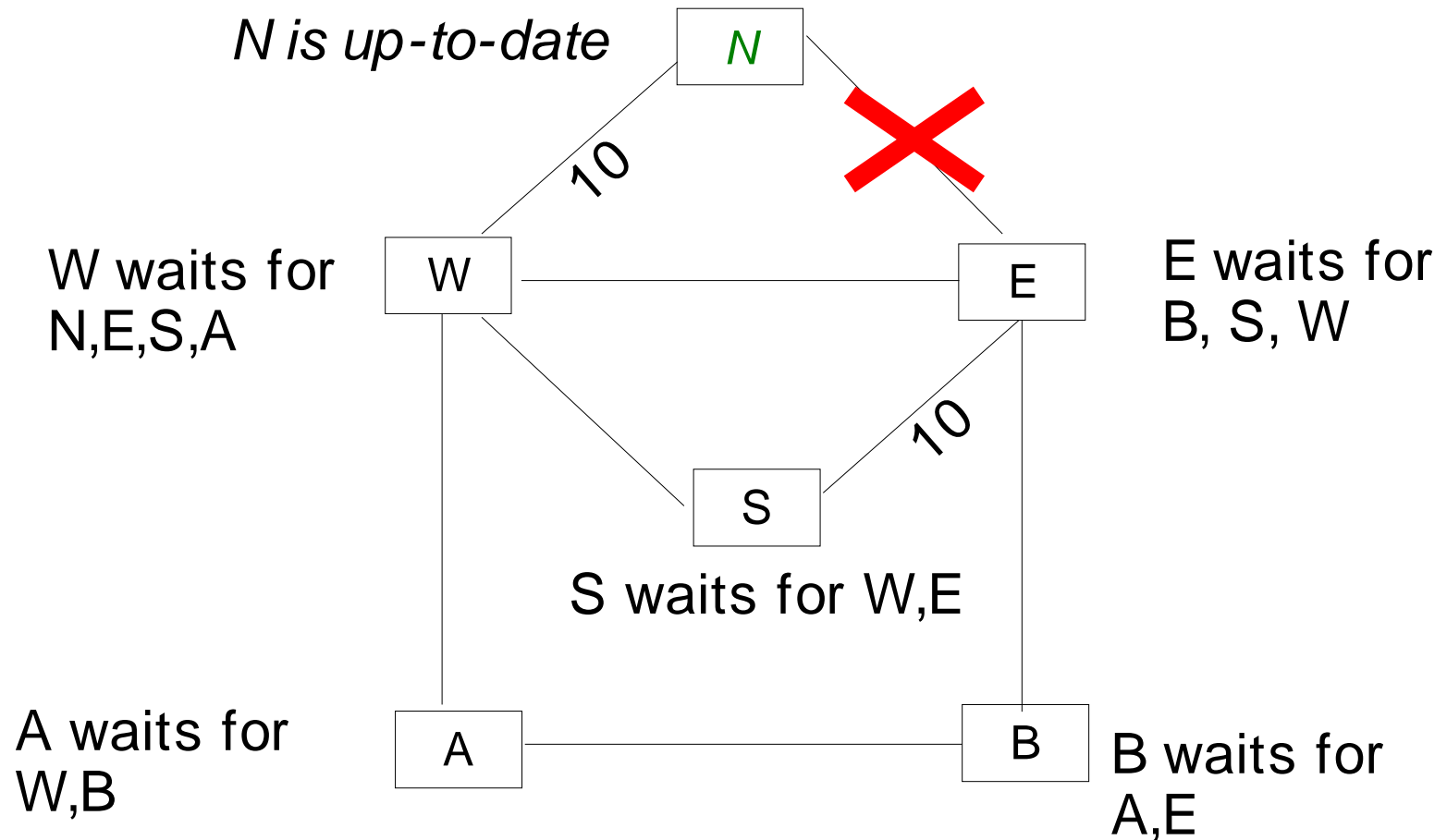
Users of link E->N



→ Packets sent to N (RSPT(N))

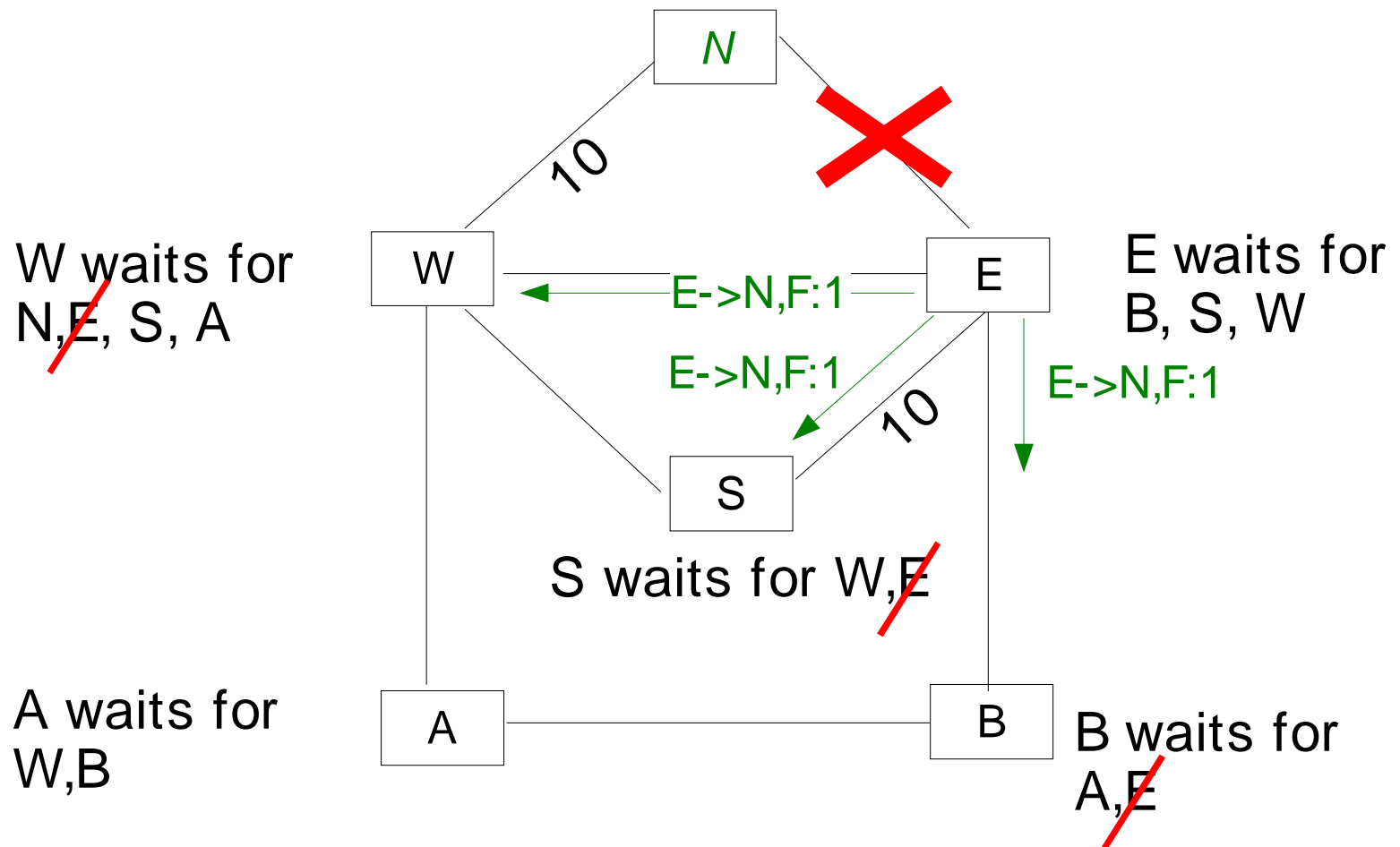
Graceful failure of link E->N

- The waiting lists



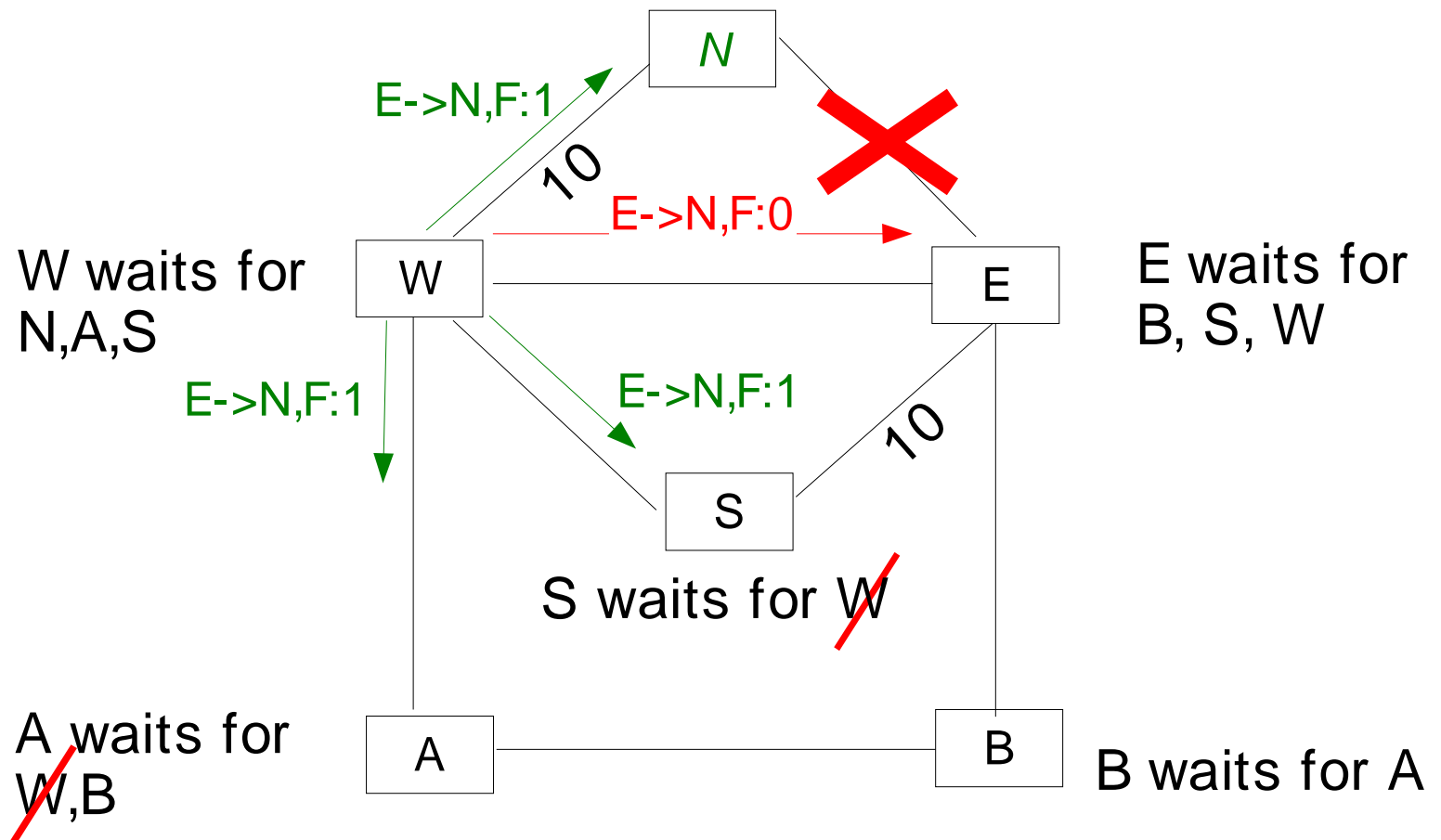
Graceful failure of link E->N (2)

- Exchange of the HELLO(E->N) PDUs



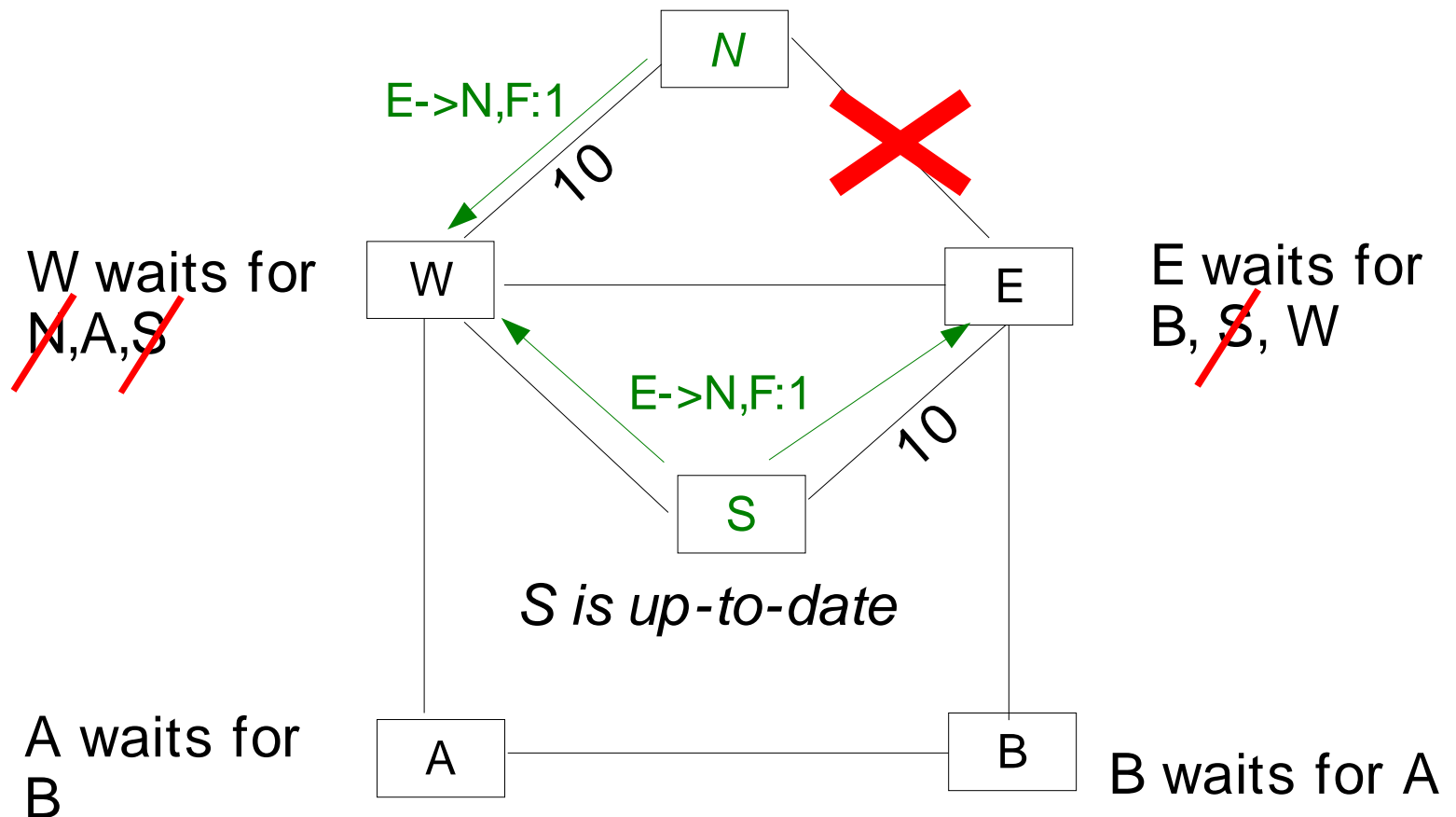
Graceful failure of link E->N (3)

- Exchange of the HELLO(E->N) PDUs



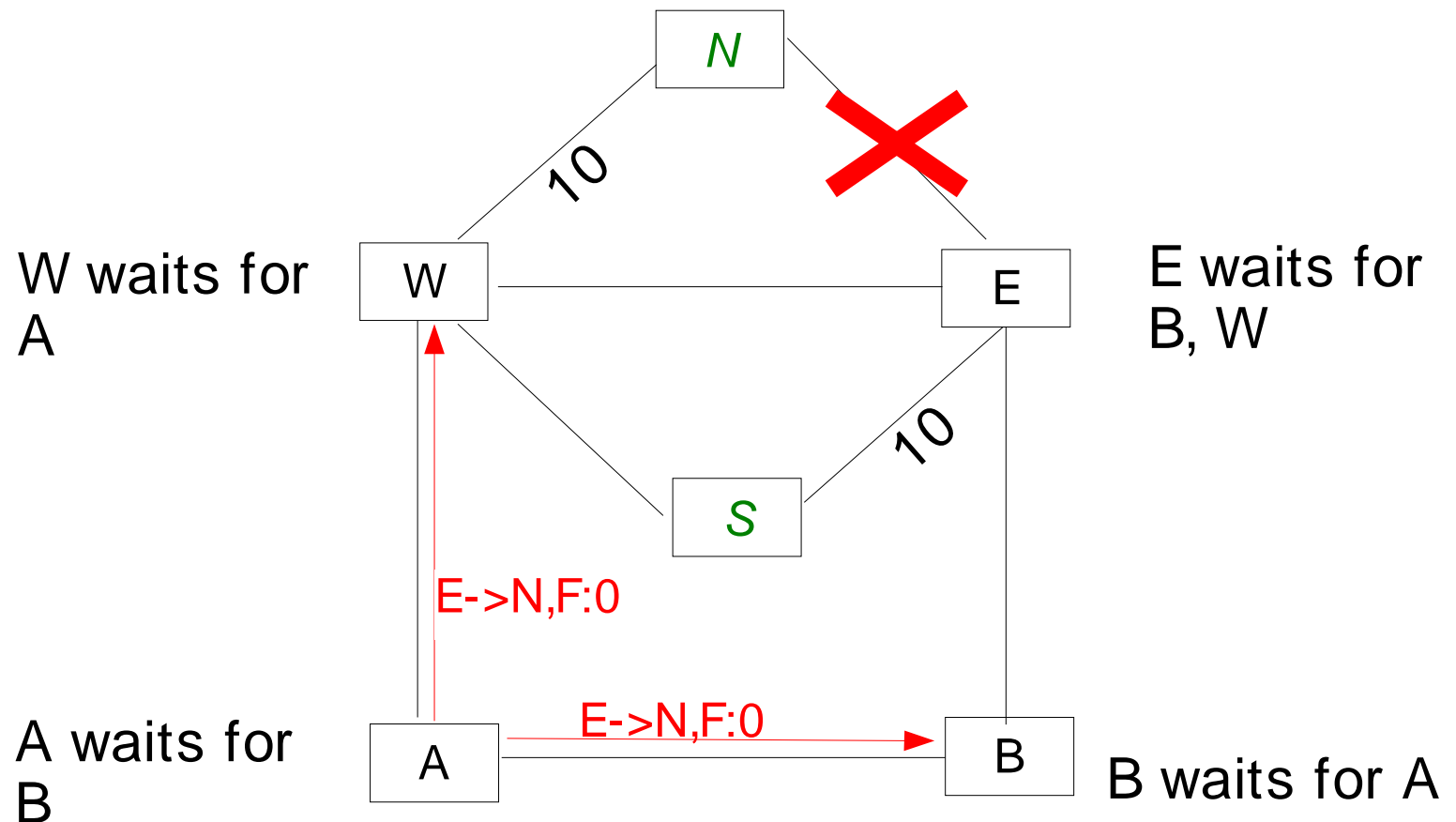
Graceful failure of link E->N (4)

- S has updated its FIB



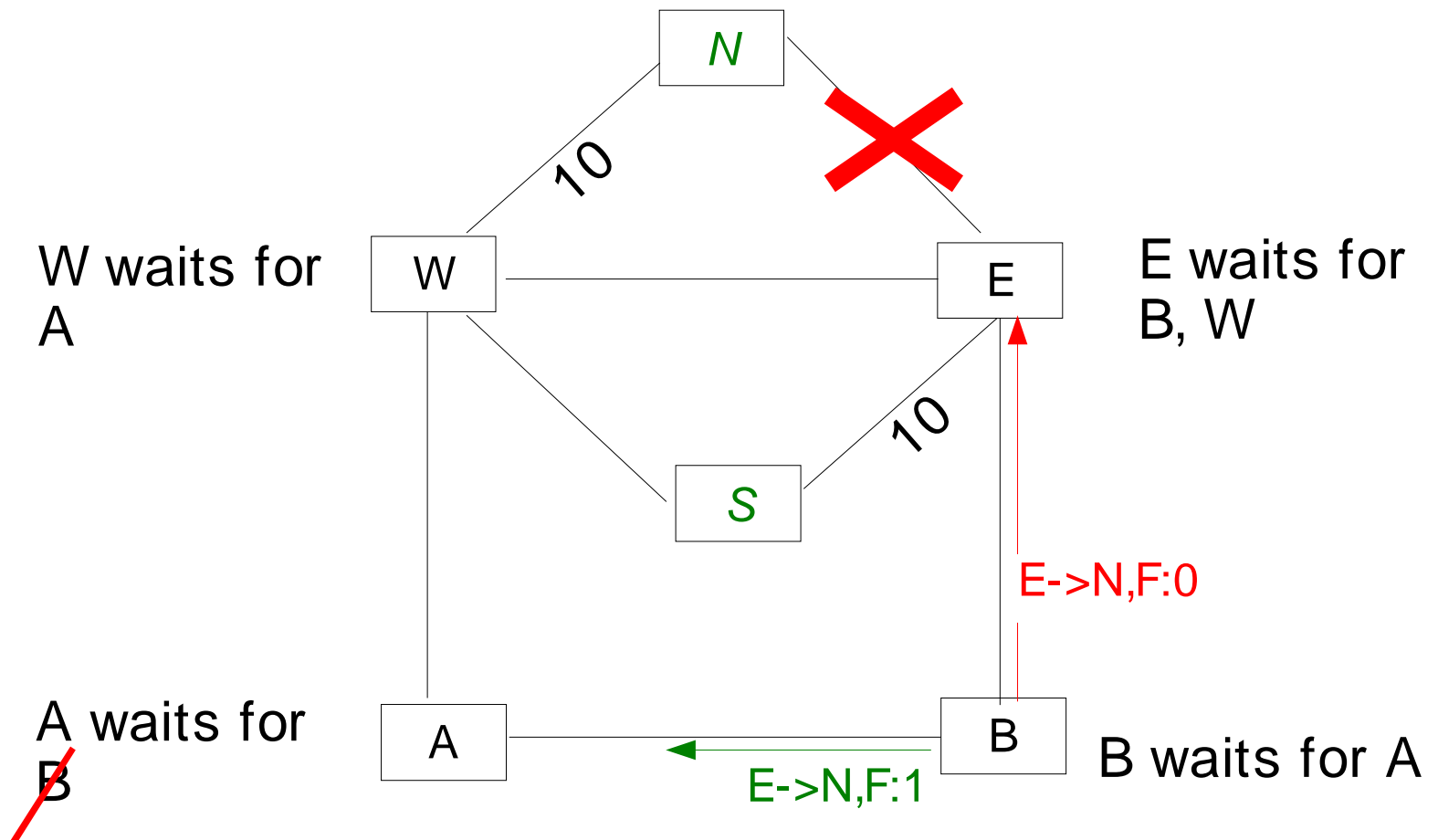
Graceful failure of link E->N (5)

- A sends its initial HELLO(E->N)



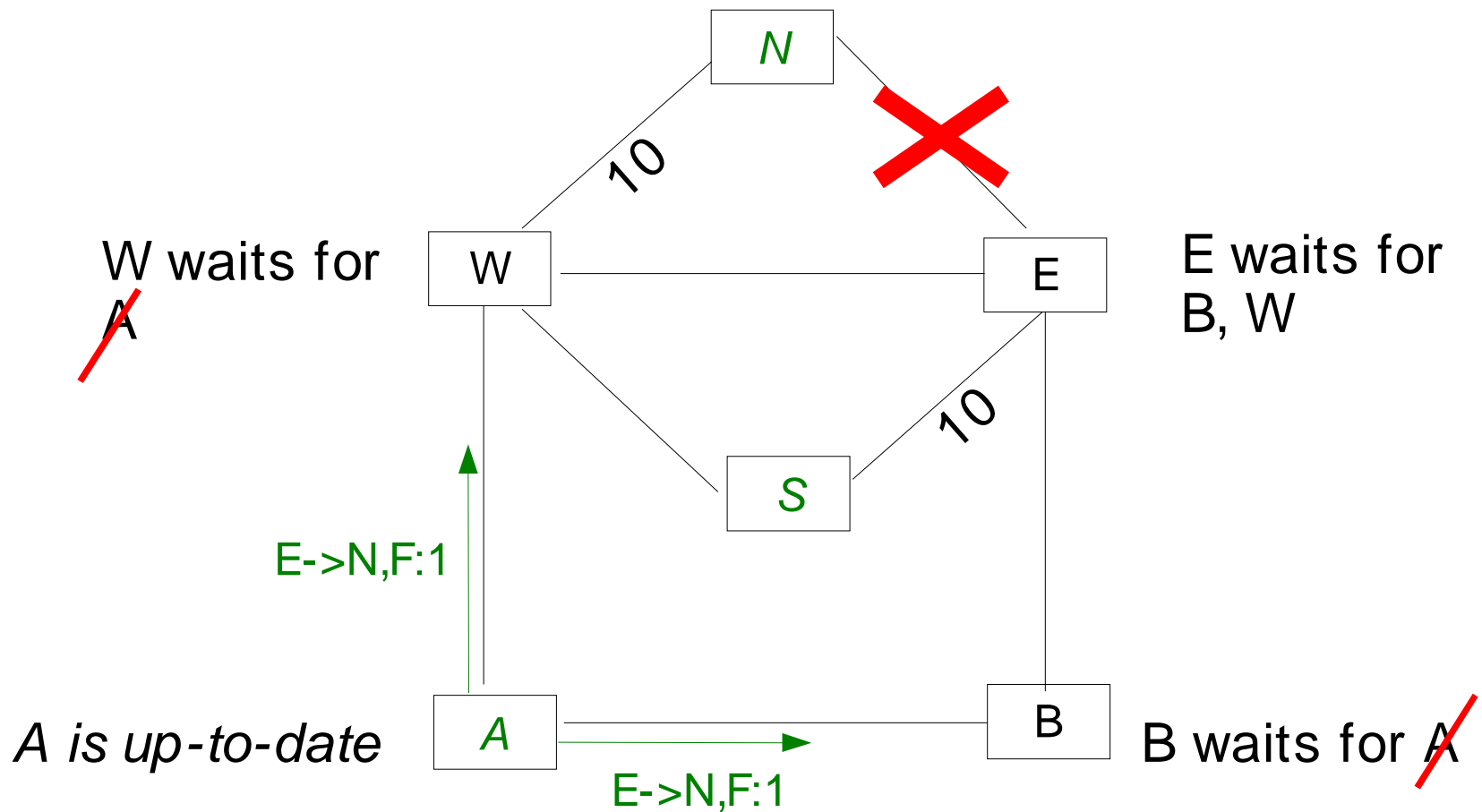
Graceful failure of link E->N (6)

- B sends its initial HELLO(E->N)



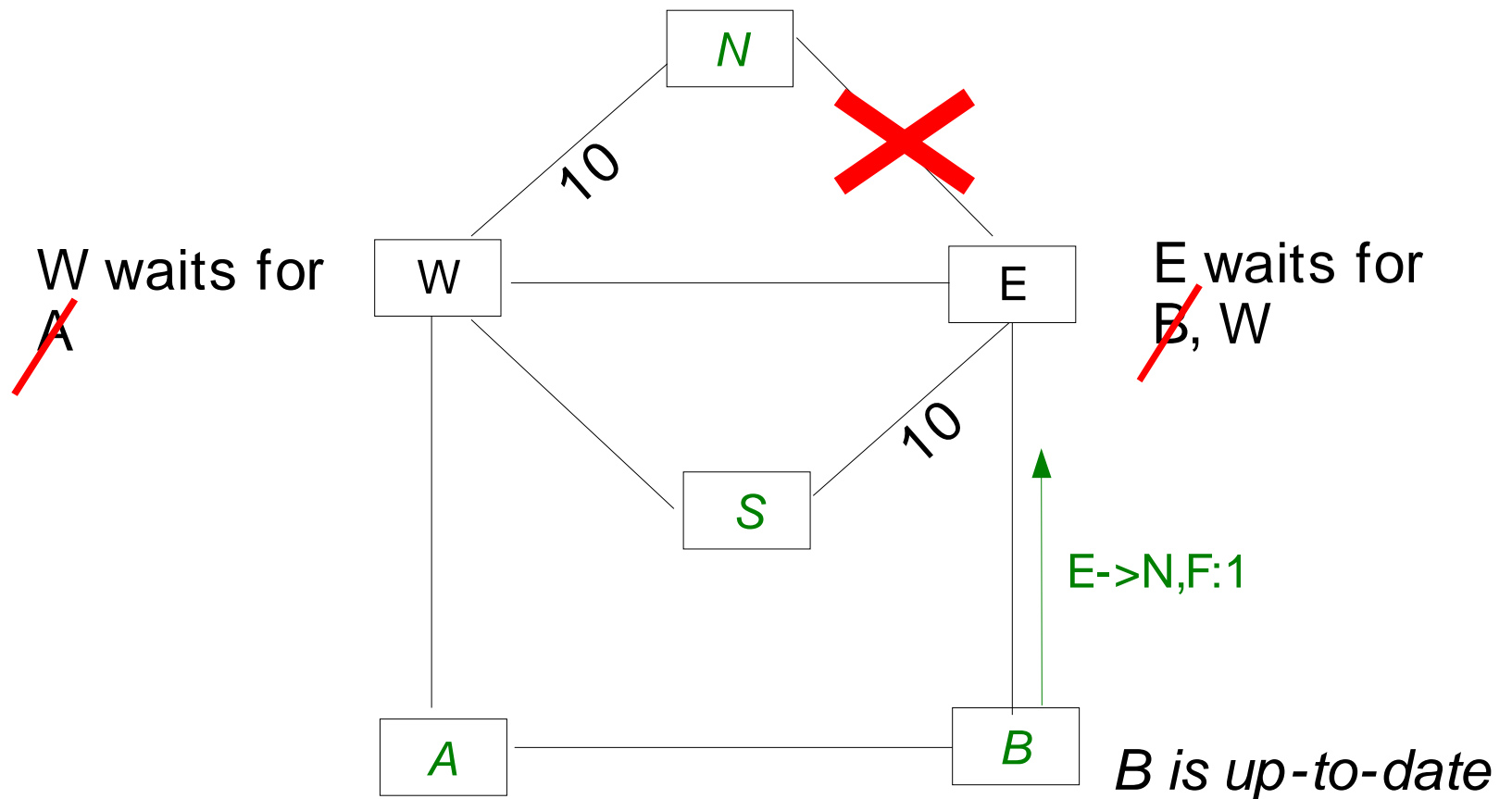
Graceful failure of link E->N (7)

- A updates its FIB



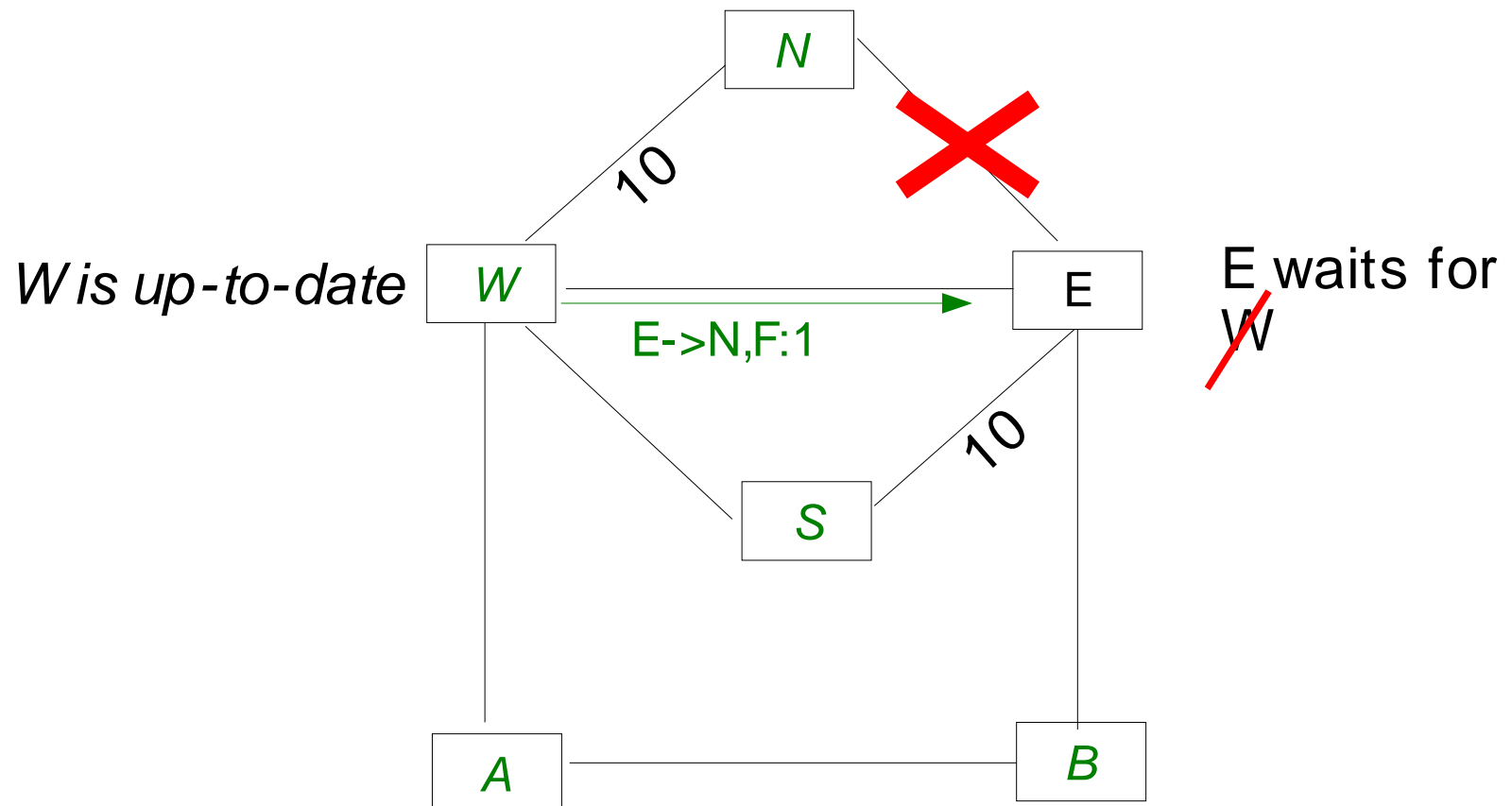
Graceful failure of link E->N (8)

- B updates its FIB



Graceful failure of link E->N (9)

- W updates its FIB



- E can safely update its FIB
 - No transient loops during IGP convergence

Other types of non-urgent IS-IS events

- The proposed protocol supports all single link changes
 - link up as well as link down
 - link metric increase
 - link metric decrease

- and also
 - Non-urgent router failures
 - ◆ Transitions of overload bit from Unset to Set
 - Non-urgent router arrivals
 - ◆ Transitions of overload bit from Set to Unset

Ongoing work with IP fast reroute

- IGP areas
 - Current solutions were designed by considering a single OSPF/ISIS area
 - Extensions to support areas are necessary
- SRLG failures
 - Several links can fail at the same time
 - ◆ links through the same fibre or same interface
 - Issues with SRLG failures
 - ◆ IGP must know the SRLG of each link
 - ◆ Accurately documenting SRLG may be a difficult operational issue in a network where the physical topology is managed by one team and the IP routers are managed by another
 - ◆ If two links share the same SRLG, they do not necessarily fail at the same time

Conclusion

- IGP behaviour in large ISP networks
 - configuration tuning can reduce IGP load
- IGP convergence after a failure
 - sub second convergence
 - ◆ possible with some IGP tuning in worldwide network
 - ◆ easy in small or MAN network
- IP fast reroute
 - Several techniques being developed to provide sub-50 millisecond recovery for intradomain link failures
 - Providing sub-50 millisecond recovery in global Internet is a difficult research challenge

For more information

- The need for fast IGP convergence
 - ◆ C. Alaettinoglu, V. Jacobson, and H. Yu. Towards millisecond IGP convergence. Internet draft, draft-alaettinoglu-ISIS-convergence-00.ps, work in progress, November 2000.
 - ◆ N. Dubois, B. Fondeviole, and N. Michel. Fast convergence project. Presented at RIPE47, <http://www.ripe.net/ripe/meetings/ripe-47/presentations/ripe47-routing-fcp.pdf>, January 2004.
 - ◆ G. Iannaccone, C. Chuah, S. Bhattacharyya, and C. Diot. Feasibility of IP restoration in a tier-1 backbone. *IEEE Network Magazine*, January-February 2004.
- ISIS and OSPF behaviour in real networks
 - ◆ A. Markopoulou, G. Iannaccone, S. Bhattacharyya, C.-N. Chuah, and Ch. Diot. Characterization of failures in an IP backbone. In *IEEE Infocom2004*, Hong Kong, March 2004.
 - ◆ A. Shaikh, C. Isett, A. Greenberg, M. Roughan and J. Gottlieb, A case study of OSPF behavior in a large enterprise network, Proceedings of the 2nd ACM SIGCOMM Workshop on Internet measurement

For more information (2)

● IP Fast reroute and related

- A. Atlas, R. Torvi, G. Choudhury, C. Martin, B. Imhoff, and D. Fedyk. Ip/Idp local protection. Internet draft, draft-atlas-ip-local-protect-00.txt, work in progress, February 2004.
- S. Bryant, C. Filsfils, S. Previdi, and M. Shand. IP Fast Reroute using tunnels. Internet draft, draft-bryant-ipfrr-tunnels-00.txt, work in progress, May 2004.
- M. Shand. "IP fast reroute framework". Internet draft, draft-ietf-rtgwg-ipfrr-framework-01.txt, June 2004.
- N. Shen and P. Pan. Nexthop Fast ReRoute for IP and MPLS. Internet draft, draft-shen-nhop-fastreroute-00.txt, work in progress, December 2003.
- P. François and O. Bonaventure, Avoiding transient loops during IGP convergence, IEEE INFOCOM 2005, March 2005
- IETF work : <http://psg.com/zinin/ietf/rtgwg/>