# Improving Load Balancing with Multipath Routing

*Pascal Mérindol*
*Jean-Jacques Pansiot, Stéphane Cateloin*

## ICCCN 2008

Improving Load Balancing
with Multipath Routing

**Pascal Mérindol**

# *Table of Contents*

- Multipath Routing Paradigm - Related Works
  - Goals and general context
  - Source routing
  - Hop by Hop routing

- Dijkstra Transverse at depth p
  - Dijkstra Transverse computation DT
  - Validation process DT(p)

- Load Balancing
  Related Works & Traffic Engineering module

- Evaluation
  - Path diversity
  - TE results

August 4, 2008

**ICCCN2008**
Improving Load Balancing
with Multipath Routing
*Pascal Mérindol*

2/17

# Multipath Routing
## Outlines

- Purposes and objectives
  - *Load balancing – circumvent congestions*
  - *Protection and restoration - circumvent failures*
  - ➢ Increase throughput and reduces delays
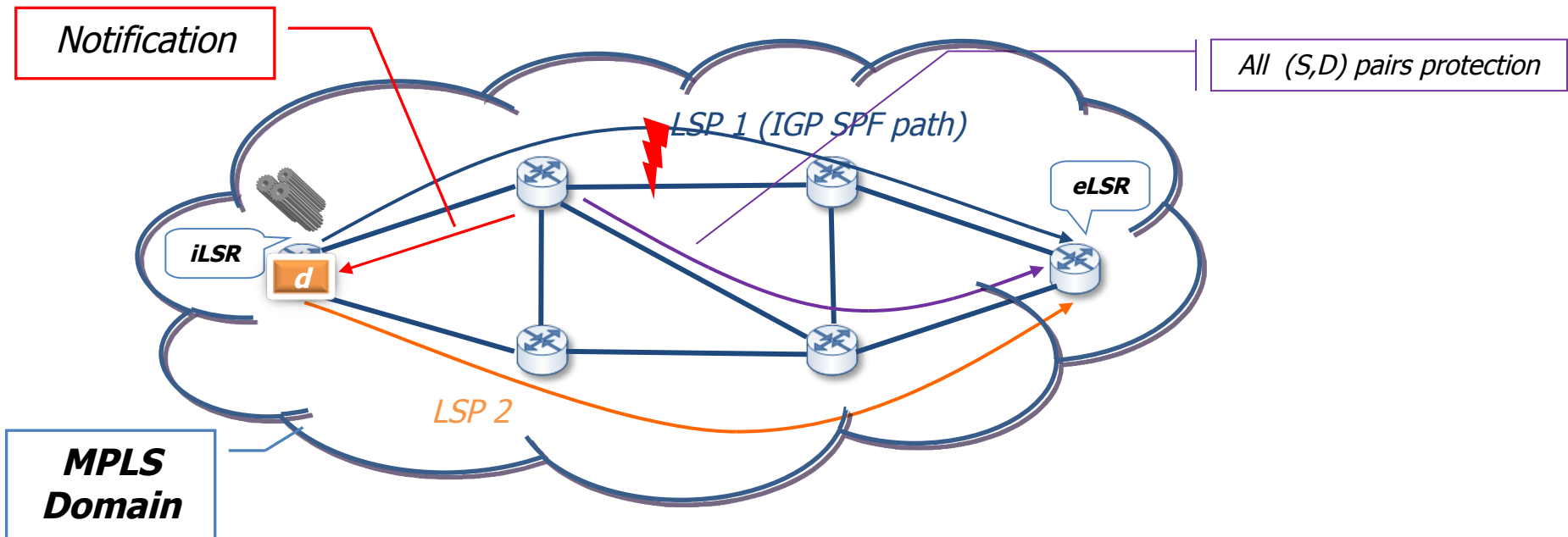  - ➔ Two ways : **source** or **hop by hop** routing

- General scheme
  1. *Path computation & signalization or validation (loop free paths)*
  2. *Path or link (global or local) traffic analyse*
  3. *Load balancing policy*
  4. *Traffic splitting*

August 4, 2008

**ICCCN2008**
Improving Load Balancing
with Multipath Routing
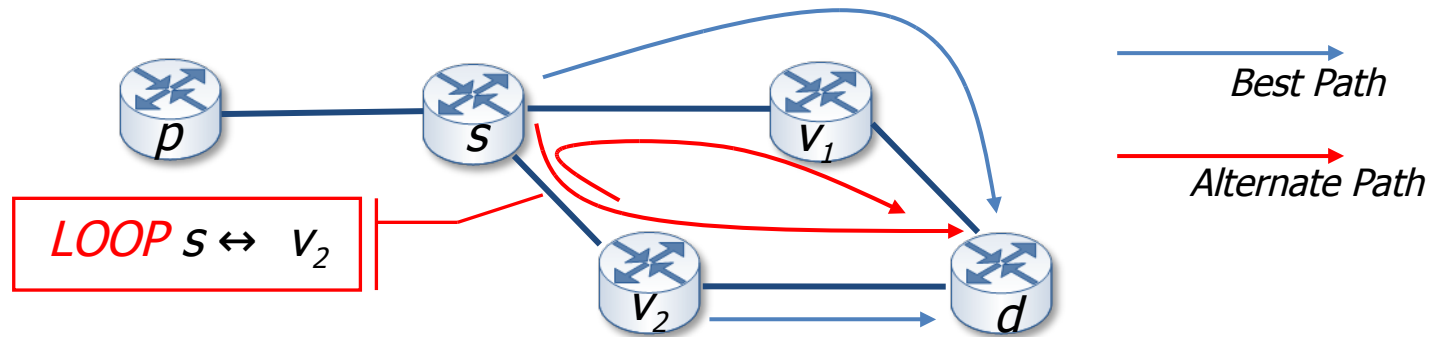**Pascal Mérindol**

3/17

# MPLS based scheme
## Pseudo-Source routing

- Multi protocol label switching (**MPLS**)
*Explicit path signalling mechanisms such as **RSVP-TE** or **CR-LDP***
- Additional label Switch Paths (**LSP**)
*Since the ingress Label Switch Router (**LSR**) towards the egress LSR*

Notification

All (S,D) pairs protection

LSP 1 (IGP SPF path)

eLSR

iLSR

d

LSP 2

**MPLS Domain**

*August 4, 2008*

**ICCCN2008**
Improving Load Balancing
with Multipath Routing
**Pascal Mérindol**

*4/17*

# Multipath hop by hop routing
## Related works



**LOOP FREE PROPERTY**

Loops can be condidered at two levels :

- ## At node level
  - *Equal Cost Multipath Routing : ECMP*
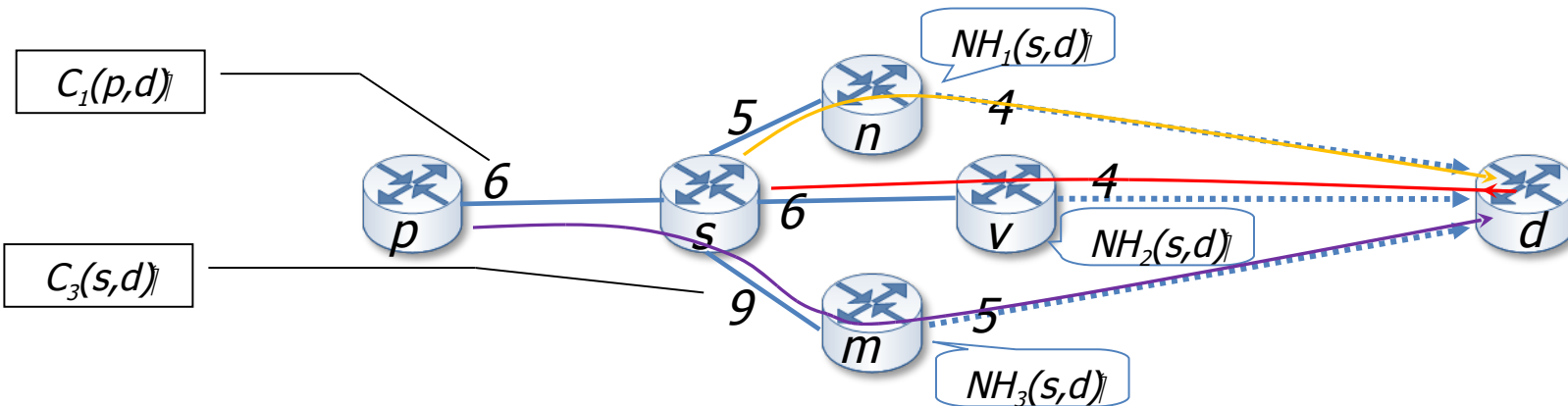  - *Downstream Criteria (One hop vision) : OMP-OSPF, LFI , etc.*
- ## At link level
  - *Two Hop vision depending to the incoming interface*

*August 4, 2008*

**ICCCN2008**
Improving Load Balancing
with Multipath Routing
**Pascal Mérindol**

*5/17*

## Loop free conditions

An alternate next hop v is viable if :

$$ECMP : \quad C_j(s,d) = C_1(s,d) \; / \; v = NH_j(s,d)$$

$$Downstream\ Criteria\ (LFI) : \quad C_1(v,d) < C_1(s,d)$$

$$Two\ Hop\ vision : \quad C_1(v,d) < C_1(p,d)$$

$C_1(p,d)$

$C_3(s,d)$

$NH_1(s,d)$

$NH_2(s,d)$

$NH_3(s,d)$

$C_j(s,d)$ : $j^{th}$ best cost computed on s towards d

$NH_j(s,d)$ : $j^{th}$ next hop computed on s towards d

**ICCCN2008**
Improving Load Balancing
with Multipath Routing
**Pascal Mérindol**

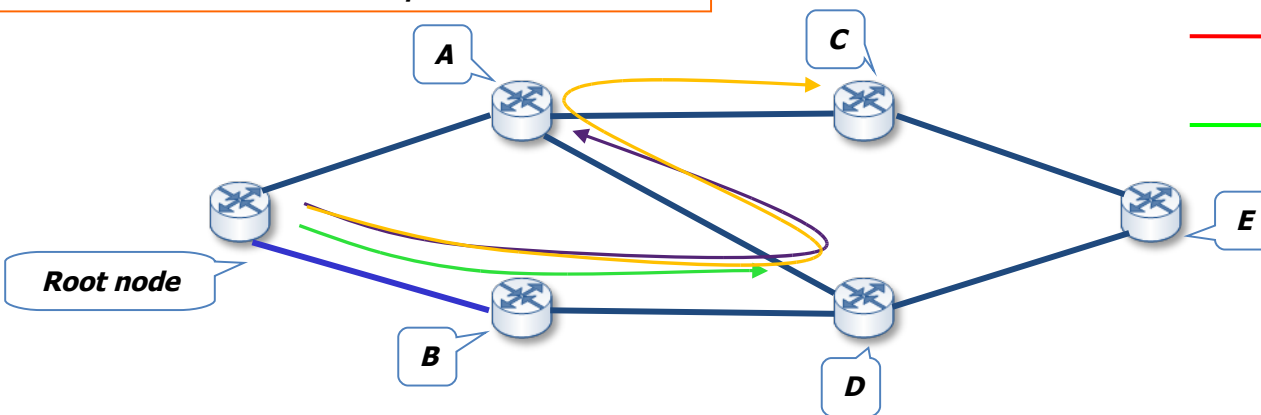# Dijkstra Transverse
## Path computation

- Dijkstra Transverse (DT) is a enhanced SPT algorithm

- DT computes at least one alternate next hop to every destination

### DT computes four sets of paths

- *Best cost path*

- *Simple transverse path*

- *Backward transverse path*

- *Forward transverse path*

### Next Hop candidates are store in a cost matrix

- *2 dimensions :*
  - cardinal of the successor set of the root
  - cardinal of the destination set



— Shortest path tree
*(with a lexicographical order)*

— Transverse edge

|   | A | B | C | D | E |
|---|---|---|---|---|---|
| A | 1 | 3 | 2 | 2 | 3 |
| B | 3 | 1 | 4 | 2 | 3 |

August 4,  2008

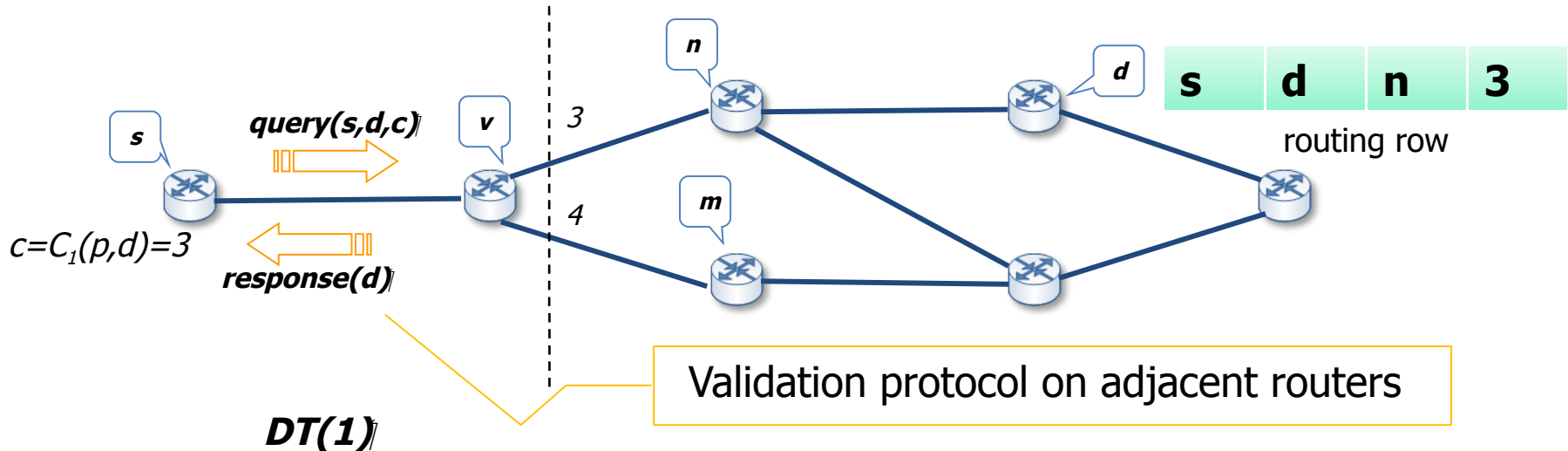**ICCCN2008**
Improving Load Balancing
with Multipath Routing
**Pascal Mérindol**

7/17

- Validation at the granularity of the **incoming interface**

- The incoming interface <u>loop free criteria</u>

$$C_j(v,d) \leq C_1(s,d)$$

candidates on v

| d | n | 3 |
|---|---|---|
| d | m | 4 |

| s | d | n | 3 |
|---|---|---|---|

routing row



$c = C_1(p,d) = 3$

**query(s,d,c)**

**response(d)**

**DT(1)**

Validation protocol on adjacent routers

*August 4, 2008*

**ICCCN2008**
Improving Load Balancing
with Multipath Routing
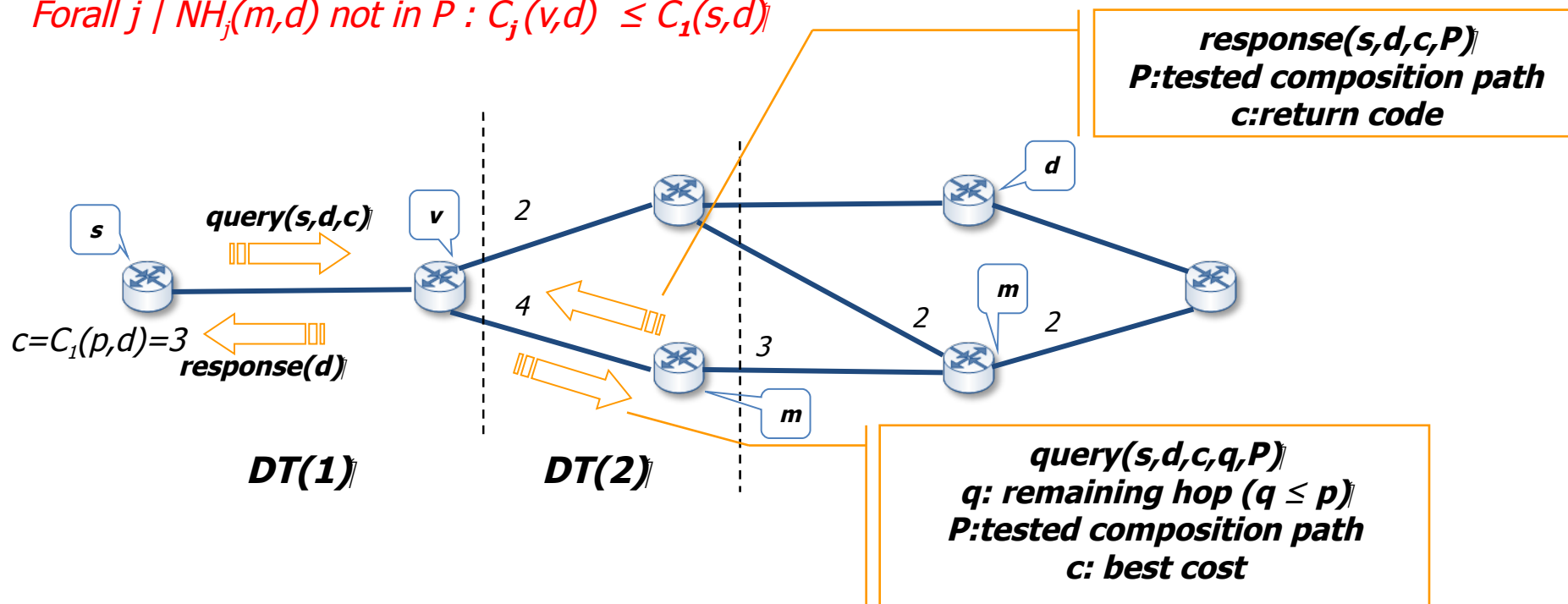**Pascal Mérindol**

*8/17*

# DT(p) Validation process

## Routing row validation at depth p

- Validation process extension

*Each path composition P is evaluated with a wave of request (On the exemple, p=2 and P={v})*

- The loop free criteria becomes

Forall $j \mid NH_j(m,d)$ not in $P : C_j(v,d) \leq C_1(s,d)$



**response(s,d,c,P)**
**P:tested composition path**
**c:return code**

query(s,d,c)

c=C₁(p,d)=3

$c=C_1(p,d)=3$
**response(d)**

**query(s,d,c,q,P)**
**q: remaining hop (q ≤ p)**
**P:tested composition path**
**c: best cost**

**DT(1)**          **DT(2)**

August 4, 2008

**ICCCN2008**
Improving Load Balancing
with Multipath Routing
**Pascal Mérindol**

9/17

## Related works and TCP incidence

- ● **Load balancing**

  - ●Modifying links weights *– [Fortz & Thorup, Wang & al., Sridharan & al., etc]*

  - ●Optimization statement with logical paths *(e.g MPLS) :*

    **OFFLINE** *(with traffic matrix) oblivious or/and normal case routing –*
    *[Applegate & Cohen, Zhang & al., COPE, etc]*

    **ONLINE** *(with probe protocols) – [MATE, TeXCP, etc]*

  - ●Incremental heuristics for hop by hop routing *– [ECMP, Vutkury &*
    *Garcia-Luna-Aceves, OSPF-OMP, Gojmerac &al., etc]*

- ● **Traffic splitting** (without disordering TCP packets?)

  - ●packet level : *round robin, probabilistic, etc*

  - ●flow level : *<src,dst,port...> hash function, tag, hybrid, etc.*

  - ●burst level *- [FLARE]*

*August 4, 2008*

**ICCCN2008**
Improving Load Balancing
with Multipath Routing
**Pascal Mérindol**

*10/17*

- **Constraints : proportions integrity**

For each destination d,
a proportions vector :  $\{x_1^d(p),\ldots,x_j^d(p),\ldots,x_n^d(p)\}$

$$\forall p \sum x_j^d(p) = 1$$

- **Objective : minimize the maximum link utilization**

Utilization ratio of an outgoing
link  l on a router s :

$$U(l) = \sum_{\forall d \in N, \forall p \in pred(s)} \frac{x_j^d(p) \times V_d(p)}{c_l}$$

*If  $x_j^d(p)$ corresponds to outgoing link l*

*traffic coming from p to d*

$$\min \ \max U(l)$$

- **Local load shifting incremental process**

Two parameters for up and
down reaction thresold :

*transient*

$$\underbrace{0,\ldots,\beta}_{non\ stressed}, \ \ldots \ , \underbrace{\alpha,\ldots,1}_{stressed} \quad \forall p,d \ \ x_j^d(p) \leftarrow x_j^d(p) \times \frac{\alpha}{U(l)}$$

*August 4,  2008*

**ICCCN2008**
Improving Load Balancing
with Multipath Routing
**Pascal Mérindol**

*11/17*

**Path Diversity**

Alternet

Open Transit

Traffic engineering

Geant

| | #node | #link | Diameter |
|---|---|---|---|
| Alternet | 83 | 334 | 8 |
| Open Transit | 76 | 206 | 11 |
| *Geant* | *23* | *74* | *6* |

**ICCCN2008**
*Improving Load Balancing
with Multipath Routing*
**Pascal Mérindol**

*August 4, 2008*

*12/17*

# Path diversity
# Number of routes and rerouting capacities



DT(3)



Open Transit & Alternet  (routes number)



|  | LFI | DT(1) | DT(3) |
|---|---|---|---|
| Open transit | 18 | 98 | 99 |
| Alternet | 16 | 60 | 78 |
| Geant | 37 | 37 | 75 |

Local re-routing capacities

*August 4,  2008*

**ICCCN2008**
Improving Load Balancing
with Multipath Routing
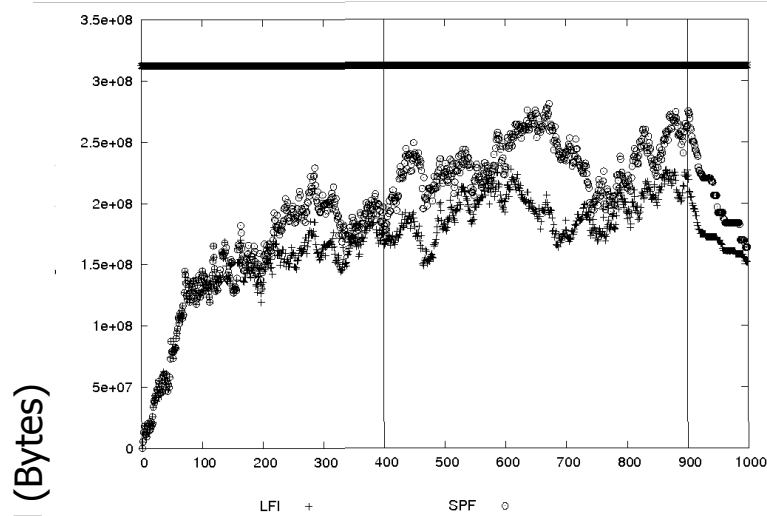*Pascal Mérindol*

*13/17*

# TE results
## Simulations setup on Geant network

- Based on realistic traces : **Totem** traffic matrix (900 s)
  - *Each entry is decomposed in TCP flows (Reno → Sack)*
  - *GEANT is over-provisioned → artificial congestions*

*August 4,  2008*

**ICCCN2008**
Improving Load Balancing
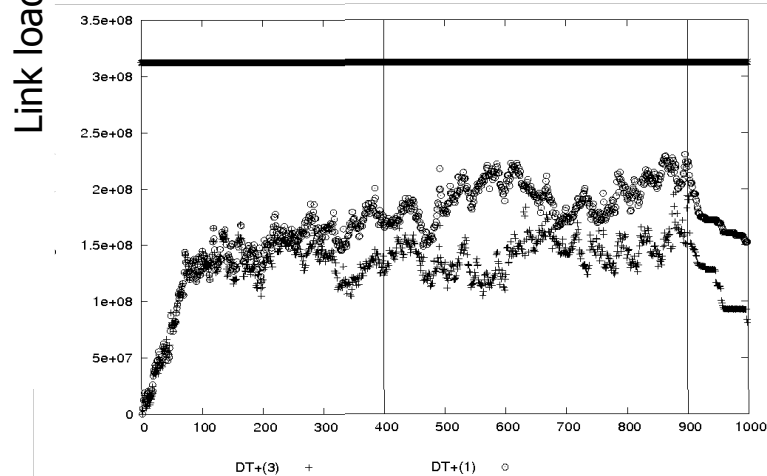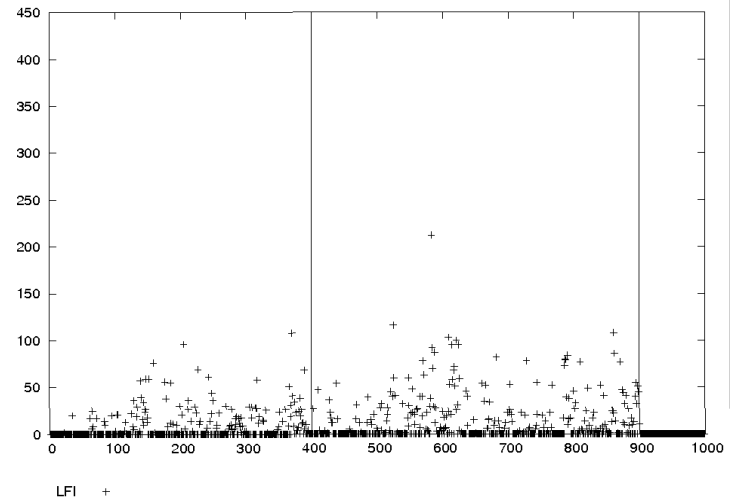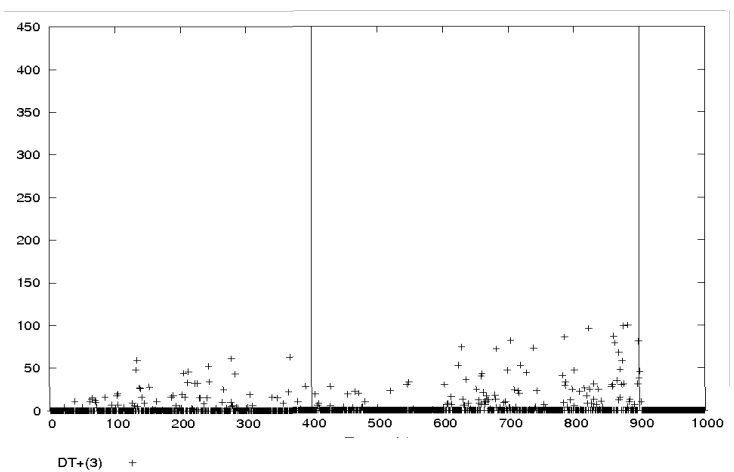with Multipath Routing
**Pascal Mérindol**

*14/17*

# TE results
## A single case as an example



LFI vs
SPF

DT(3) vs
LFI

Link load (Bytes)

Loss (Dropped packets)

August 4, 2008

**ICCCN2008**
*Improving Load Balancing
with Multipath Routing*
**Pascal Mérindol**

15/17

- ## Configuration parameters and measured indicators
  - ### *α=0.5, β=0.25  and t=1s* (sender windows bounded by 65 packets)
  - ### *link load and dropped packets*



### Results compilation

|  | LFI | DT(1) | DT(3) |
|---|---|---|---|
| Average loss reduction ratio (compared as SPF) | 3.8 | 4.2 | 6.5 |
| Average load of most loaded link (SPF : 76%) | 61.4 | 61.4 | 51.8 |

### Averages calculated on 12 simulations :
- congestions are triggered on the most natural loaded link (1→ 1, 1→ n, n→ 1)
- for each run, the link load average confidence intervall (95%) is below 0.1% of the link capacity

*August 4,  2008*

**ICCCN2008**
Improving Load Balancing
with Multipath Routing
**Pascal Mérindol**

16/17

# Conclusions
## *Perspectives & work in progress*

- ## Multipath Routing and Path Diversity
  - The efficiency of load balancing scheme depends of the path diversity *(routes number, coverage & cumulative bandwidth)*
  - DT(p)-TE allows to reduce congestions impact
  - Global notification can enhance the redirection coverage

- ## Current Work
  - Notification and probing protocols
  - Global load balancing problem statement

- ## Future work
  - Congestions and/or failures scenarios
  - Scalability and inter-domain issues

*August 4, 2008*

**ICCCN2008**
*Improving Load Balancing
with Multipath Routing*
**Pascal Mérindol**

*17/17*