# Revisiting the *same service for all* paradigm

IP packet

**No delivery guarantee**

INTERNET

*Regular mail*

**Enhancing the best-effort service**

**Introduce Service Differentiation**

P r I o r I t a I r e
First Class Letter
**A**

IP packet

DiffServ

1

---

# Service Differentiation

> The real question is to choose which packets shall be dropped. The first definition of differential service is something like "not mine."
> -- Christian Huitema

- ❑ Differentiated services provide a way to specify the relative priority of packets
- ❑ Some data is more important than other
- ❑ People who pay for better service get it!

**SLA**
Service Level Agreement

DiffServ

2

---

1

# Divide traffic into classes

**Differentiated IP Services**

**E-Commerce**

**Application Traffic**

**E-mail, Web Browsing**

**Voice**

**Traffic Classification**

**Voice** — Platinum Class Low Latency

**Gold** — Guaranteed: Latency and Delivery

**Silver** — Guaranteed Delivery

**Bronze** — Best Effort Delivery

Borrowed from Cisco 3

---

# Design Goals/Challenges

❑ Ability to charge differently for different services
❑ No per flow state or per flow signaling
❑ All policy decisions made at network boundaries
  ❑ Boundary routers implement policy decisions by tagging packets with appropriate priority tag
❑ Traffic policing at network boundaries
❑ Deploy incrementally: build simple system at first, expand if needed in future

4

# IP implementation: DiffServ

RFC 2475

## No per flow state in the core

Flow 1
Flow 2
Flow 3
Flow 4
...

IP packet

10Gbps=2.4Mpps

with 512-byte packets

**Stateful approaches scalable at gigabit rates**

6 bits used for Differentiated Service Code Point (DSCP) and determine PHB that the packet will receive

1997

1993

1981

DiffServ

IntServ/ RSVP

IP TOS

| IP header | IP Data Area |
|---|---|

| Ver | Len | Typ.Ser. | Total Length |
|---|---|---|---|
| | | Fl. | Frag.Offset |
| TTL | | ader Checksum |

DSCP    CU

Padding

DiffServ

---

# DiffServ building blocks

DIFFSERV

TRAFFIC CONDITIONING SHAPING → TOKEN BUCKET

MARKING POLICY → INTRA-DOMAIN

→ INTER-DOMAIN

PER HOP BEHAVIOR → SCHEDULING (RR, WRR, FQ, WFQ)

→ AQM (DT, RED,...)

DiffServ

6

# Traffic Conditioning

❑ User declares traffic profile (eg, rate and burst size); traffic is metered and shaped if non-conforming

5Mbps

SLA
2Mbps

Service
Level
Agreement

r tokens per second

b tokens

<= C bps

regulator

meter

forward

packets → classifier → marker → Shaper/dropper

drop

# Token Bucket (1)

## Example

- B = 4000 bits, R = 1 Mbps, C = 10 Mbps
- Packet length = 1000 bits
- Assume the bucket is initially full and a "large" burst of packets arrives

R

B

C = 10 Mbps

istoica@cs.cmu.edu

# Token Bucket (2)

B=4000 bits, R=1Mbps, C=10Mbps

# Token Bucket for traffic characterization

❑ Given b=bucket size, C=link capacity and r=token generation rate

# Differentiated Architecture

Ingress Edge Router

DiffServ Domain

Egress Edge Router

Interior Router

scheduling

**Marking:**

per-flow traffic management

marks packets as in-profile and out-profile

**Per-Hop-Behavior (PHB):**

per class traffic management

buffering and scheduling based on marking at edge

preference given to in-profile packets

Ingress

Egress

DiffServ

---

# Pre-defined PHB

❑ Expedited Forwarding (EF, premium):
  - ❑ departure rate of packets from a class equals or exceeds a specified rate (logical link with a minimum guaranteed rate)
  - ❑ Emulates leased-line behavior

❑ Assured Forwarding (AF):
  - ❑ 4 classes, each guaranteed a minimum amount of bandwidth and buffering; each with three drop preference partitions
  - ❑ Emulates frame-relay behavior

DiffServ

12

# Premium Service Example



Drop always

10Mbps

Fixed Bandwidth

DiffServ

source Gordon Schaffee    13

# Assured Service Example



Drop if congested

10Mbps

Assured Service

Uncongested

Congested

DiffServ

source Gordon Schaffee    14

# Border Router Functionality

Premium Service

Token Bucket

| 0 | | | | | | |
|---|---|---|---|---|---|---|
| DSCP | | | | | | |
| 1 | 0 | 1 | 1 | 1 | 0 | |

Packet Input → Data Queue → **Wait for token** → **Set P-bit** → Packet Output

| 0 | DSCP | Class 1 | Class 2 | Class 3 | Class 4 |
|---|---|---|---|---|---|
| | Low drop probability | **001**010 | **010**010 | **011**010 | **100**010 |
| | Medium drop proba. | **001**100 | **010**100 | **011**100 | **100**100 |
| | High drop proba. | **001**110 | **010**110 | **011**110 | **100**110 |

Assured Service

Token Bucket

No token

Packet Input → **Test if token** → Token → **Set A-bit** → Data Queue → Packet Output

DiffServ

source Gordon Schaffee, modified by C. Pham  15

---

# Internal Router Functionality

Packets In → P-bit set? → Yes → High Priority Queue

No

If A-bit set, a_cnt++ → Low Priority Queue

if congested

If A-bit set, a_cnt--

RED In/Out Queue Management

Packets Out

A DSCP codes aggregates, not individual flows
No state in the core
Should scale to millions of flows

DiffServ

source Gordon Schaffee, modified by C. Pham  16

# Scheduling

❑ DiffServ PHB relies mainly on scheduling
  ❑ choose the next packet for transmission
  ❑ FIFO: in order of arrival to the queue; packets that arrive to a full buffer are either discarded, or a discard policy is defined.
  ❑ More complex policies: FCFS, PRIORITY, EDD…

arrivals → [ queue ] ( link ) → departures

queue (waiting area)    link (server)

---

# Priority Queueing

❑ Priority Queuing: classes have different priorities;
❑ Transmit a packet from the highest priority class with a non-empty queue
❑ Preemptive and non-preemptive versions

high priority queue (waiting area)

arrivals → classify → high priority queue / low priority queue → link (server) → departures

low priority queue (waiting area)

arrivals

packet in service

departures

# Round Robin (RR)

❑ Round Robin: scan class queues serving one from each class that has a non-empty queue

one round

arrivals

time

packet in service

departures

time

19

# Weighted Round Robin, WRR

❑ Assign a weight to each connection and serve a connection in proportion to its weight

Connection A, B and C with same packet size and weight 0.5, 0.75 and 1. How many packets from each connection should a round-robin server serve in each round?

A: Normalize each weight so that they are all integers: we get 2, 3 and 4. Then in each round of service, the server serves 2 packets from A, 3 from B and 4 from C.

$w_1$

$w_2$     one round

$w_i$

20

# (Weighted) Round-Robin Discussion

- ❑ Advantages: protection among flows
  - ❑ Misbehaving flows will not affect the performance of well-behaving flows
  - ❑ FIFO does not have such a property
- ❑ Disadvantages:
  - ❑ More complex than FIFO: per flow queue/state
  - ❑ Biased toward large packets: a flow receives service proportional to the number of packets
- ❑ If packet size are different, we normalize the weight by the packet size
  - ❑ ex: 50, 500 & 1500 bytes with weight 0.5, 0.75 & 1.0

DiffServ

21

# Generalized Processor Sharing (GPS)

- ❑ Assume a fluid model of traffic
  - ❑ Visit each non-empty queue in turn (like RR)
  - ❑ Serve infinitesimal from each
  - ❑ Leads to "max-min" fairness
- ❑ GPS is un-implementable!
  - ❑ We cannot serve infinitesimals, only packets



**max-min fairness**

Let n sources requiring resources $x_1,..,x_n$ with $x_1<x_2..<x_n$ for instance. Server has a capacity of C.

We assign C/n to source 1. If $C/n>x_1$, give $C/n+(C/n-x_1)/(n-1)$ to the (n-1) remaining sources. If this amount is greater than $x_2$, process again.

DiffServ

22

# Packet Approximation of Fluid System

❑ GPS un-implementable

❑ Standard techniques of approximating fluid GPS

   ❑ Select packet that finishes first in GPS assuming that there are no future arrivals (emulate GPS on the side)

❑ Important properties of GPS

   ❑ Finishing order of packets currently in system independent of future arrivals

❑ Implementation based on virtual time

   ❑ Assign virtual finish time to each packet upon arrival

   ❑ Packets served in increasing order of virtual times

---

# Fair Queuing (FQ)

❑ Idea: serve packets in the order in which they would have finished transmission in the fluid flow system

❑ Mapping bit-by-bit schedule onto packet transmission schedule

❑ Transmit packet with the lowest finish time at any given time

# Weighted Fair Queueing

❑ Variation of FQ: Weighted Fair Queuing (WFQ)
❑ Weighted Fair Queuing: is a generalized Round Robin in which an attempt is made to provide a class with a differentiated amount of service over a given period of time

# Implementing WFQ

❑ WFQ needs per-connection (or per-aggregate) scheduler state→implementation complexity.
   ❑ complex iterated deletion algorithm
   ❑ complex sorting at the output queue on the service tag
❑ WFQ needs to know the weight assigned for each queue →manual configuration, signalling.
❑ WFQ is not perfect...
❑ Router manufacturers have implemented as early as 1996 WFQ in their products
   ❑ from CISCO 1600 series
   ❑ Fore System ATM switches

# Putting it together!



DiffServ

Source VTHD

27

# DiffServ for grids



scheduling

marking

Wide-area interactive simulations

FTP

Ingress/Ingress

Egress

Egress

Egress

Egress

Assured Forwarding

Premium

28

14

# DiffServ for grids (con't)

Wide-area interactive simulations

FTP

scheduling

Ingress/Ingress

Egress

Egress

A DSCP codes aggregates, not individual flows
No state in the core
Should scale to millions of flows

Assured Forwarding

Premium

29

# Bandwidth provisioning

N
E
W

C
H
A
P
T
E
R

❑ DWDM-based optical fibers have made bandwidth very cheap in the backbone

❑ On the other hand, dynamic provisioning is difficult because of the complexity of the network control plane:

❑ Distinct technologies

❑ Many protocols layers

❑ Many control software

IP

ATM

SONET/SDH

DWDM

MPLS

30

# Provider's view



Today's setting time is several weeks/months! We want to set dynamic links within hours

# Review of IP routing



| Destination | Next Hop |
|---|---|
| D | R3 |
| E | R3 |
| F | R5 |

# Review of IP routing



| | | | |
|---|---|---|---|
| Ver | HLen | T.Service | Total Packet Length |
| Fragment ID | | Flags | Fragment Offset |
| TTL | | Protocol | Header Checksum |
| Source Address | | | |
| Destination Address | | | |
| Options (if any) | | | |
| Data | | | |

1    4         16         32

20 bytes

33

---

# Review of telephone network



*First automatic Branch Exchange Almond B. Strowger, 1891…*

**Signaling replaces the operator**

Source J. Tiberghien, VUB

34

# The telephone circuit view

Trunk
lines

SW
SW
SW
SW
SW
SW
SW
SW
SW
PABX
PABX

# Advantages of circuits

❑ Provides the same path for information of the same connection: less out-of-order delivery

❑ Easier provisioning/reservation of network's resources: planning and management features

# Time Division Circuits

- ❑ Most trunks time division multiplex voice samples
- ❑ At a central office, trunk is demultiplexed and distributed to active circuits
- ❑ Synchronous multiplexor
  - ❑ N input lines
  - ❑ Output runs N times as fast as input

Simple, efficient, but low flexibility and wastes resources

**1**
**2**
**3**
**. . .**
**N**

**MUX**

**Fixed bandwitdh**

| 1 | 2 | 3 | ... | N |

**De-MUX**

**1**
**2**
**3**
**. . .**
**N**

**1 sample every 125us gives a 64Kbits/s channel**

---

# Back to virtual circuits

❑Virtual circuit refers to a connection oriented network/link layer: e.g. X.25, Frame Relay, ATM

Virtual
Circuit
Switching:
a path is defined
for each connection

R3

R1

A

R4

D

B

E

C

R2

R5

F

**But IP is connectionless!**

# Virtual circuit principles

Connections &
Virtual circuits table

label

R3

| Label IN | Link IN | Label OUT | Link OUT |
|---|---|---|---|
| 23 | 1 | 34 | 3 |
| 45 | 2 | 78 | 4 |

23   78
Link 1   Link 4
Link 2   Link 3
45   34

R3

R1   R4   D
A

Virtual
Circuit
Switching

B   E

C   R2   R5

---

# End-to-end operation (1)

| VCI E. | Lien E. | VCI S. | Lien S. |
|---|---|---|---|
|  |  |  |  |

| 17 | A | 1 |
|---|---|---|
| 34 | B | 2 |
| 23 | C | 3 |

| VCI E. | Lien E. | VCI S. | Lien S. |
|---|---|---|---|
| 13 | 1 | 05 | 0 |

A

| 05 | A | 2 |
|---|---|---|
| 23 | B | 0 |
| 41 | C | 1 |

| 05 | A | 0 |
|---|---|---|
| 67 | B | 2 |
| 05 | C | 2 |

0   CR@B   5

CR@B   0

2

B

| VCI E. | Lien E. | VCI S. | Lien S. |
|---|---|---|---|
| 0 | 0 | 45 | 2 |
|  |  |  |  |

1   CR@B   13

CR@B   45

| 12 | A | 2 |
|---|---|---|
| 15 | B | 1 |
| 62 | C | 0 |

C

| VCI E. | Lien E. | VCI S. | Lien S. |
|---|---|---|---|
| 45 | 2 | 13 | 1 |
|  |  |  |  |

# End-to-end operation (2)

---

# Why virtual circuit?

❑ Initially to speed up router's forwarding tasks: X.25, Frame Relay, ATM.

We're fast enough!

Now: Virtual circuits for traffic engineering!

# Virtual circuits in IP networks

❑ Multi-Protocol Label Switching
  ❑ Fast: use label switching➔ LSR
  ❑ Multi-Protocol: above link layer, below network layer
  ❑ Facilitate traffic engineering

IP

MPLS

LINK

| PPP Header(Packet over SONET/SDH) | PPP Header | MPLS Header | Layer 3 Header |
|---|---|---|---|
| Ethernet | Ethernet Hdr | MPLS Header | Layer 3 Header |
| Frame Relay | FR Hdr | MPLS Header | Layer 3 Header |

MPLS                                                                    43

---

# Label structure

```
0                   1                   2                   3
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
       Label                 | Exp|S|      TTL
```

**Label = 20 bits**
**Exp = Experimental, 3 bits**
**S = Bottom of stack, 1bit**
**TTL = Time to live, 8 bits**

❑ More than one label is allowed -> Label Stack
❑ MPLS LSRs always forward packets based on the value of the label at the top of the stack

MPLS                                                                    44

# From multilayer networks...

FR access

ATM

**SONET/SDH ring**

DCS or ADM

**SONET/SDH ring**

DCS or ADM

MPLS

45

# ...to IP/MPLS networks

Multi-access

Label Switch Router

**IP/MPLS**

link 1

Multi-access

UNIVERSITY

MPLS

46

23

# MPLS operation

**1a. Routing protocols (e.g. OSPF-TE, IS-IS-TE) exchange reachability to destination networks**

**4. LSR at egress removes label and delivers packet**

**1b. Label Distribution Protocol (LDP) establishes label mappings to destination network**

Label Switch Router

IP

IP 10

IP 20

IP 40

IP

134.15.8.9

| src | dest | out |
|-----|------|-----|
| * | 134.15/16 | 1/10 |
| * | 140.134/16 | 1/26 |

link 1

**2. Ingress LSR receives packet and "label"s packets**

Source Yi Lin, modified C. Pham

**3. LSR forwards packets using label switching**

MPLS

47

---

# Label Distribution

| In Link | In label | dest | out link | out label |
|---------|----------|------|----------|-----------|
| 0 | - | 134.15/16 | 1 | **10** |

| In Link | In label | dest | out link | out label |
|---------|----------|------|----------|-----------|
| 0 | **20** | 134.15/16 | 1 | 40 |

Use label 20 for 134.15/16

Use label 40 for 134.15/16

link 1

Use label 10 for 134.15/16

134.15.8.9

| In Link | In label | dest | out link | out label |
|---------|----------|------|----------|-----------|
| 0 | **10** | 134.15/16 | 1 | 20 |

| In link | In label | dest | out link | out label |
|---------|----------|------|----------|-----------|
| 0 | **40** | 134.15/16 | 1 | - |

Unsolicited downstream label distribution

MPLS

48

---

24

# Label Distribution (con't)

| In Link | In label | dest | out link | out label |
|---------|----------|------|----------|-----------|
| 0 | – | 134.15/16 | 1 | 10 |

| In Link | In label | dest | out link | out label |
|---------|----------|------|----------|-----------|
| 0 | 20 | 134.15/16 | 1 | 40 |

Use label 10 for 134.15/16

Use label 20 for 134.15/16

Use label 40 for 134.15/16

request label for 134.15/16

request label for 134.15/16

request label for 134.15/16

| In Link | In label | dest | out link | out label |
|---------|----------|------|----------|-----------|
| 0 | 10 | 134.15/16 | 1 | 20 |

134.15.8.9

| In link | In label | dest | out link | out label |
|---------|----------|------|----------|-----------|
| 0 | 40 | 134.15/16 | 1 | – |

On-demand downstream label distribution

MPLS

49

---

# Dynamic circuits for grids

E

B

I need 2.5 Gbps between:
A & B
B & C
D & C
E & A

VISA

MPLS

50

# Forwarding Equivalent Class: high-level forwarding criteria

**Table B**
```
L4: (FEC E) C, L6
    (FEC F) D, L7
L3: (FEC X) A, L8
    (FEC Y) D, L9
L5: (FEC Z) C, L10
```

**Table C**
```
L24:(FEC X) B, L3
L25:(FEC Y) F, pop
L10:(FEC Z) E, pop
L14:(FEC Z) E, pop
L19:(FEC Z) E, pop
```

**Table A**
```
L6: (FEC F) D, L11
L8: (FEC X) A, pop
    (FEC Y) D, L12
    (FEC Z) B, L5
```

**Table E**
```
(FEC D) C, L22
(FEC F) C, L23
(FEC X) C, L24
(FEC Y) C, L25
```

**Table D**
```
L7: (FEC F) F, pop
L11:(FEC F) F, pop
L18:(FEC X) A, pop
L9: (FEC Y) F, pop
L12:(FEC Y) F, pop
    (FEC Z) C, L14
```

**Table F**
```
(FEC D) D, pop
(FEC E) C, L17
(FEC X) D, L18
(FEC Z) C, L19
```

LSR B   LSR C   LSR E   Z
L5   L10
X   LSR A
L14   L19
LSR D   LSR F   Y

MPLS                                                                51

---

# Forwarding Equivalent Class

A FEC aggregates a number of individual flows with the same characteristics: IP prefix, router ID, delay or bandwidth constraints…

**Table A**
```
L6: (FEC F)
L8: (FEC X)
    (FEC Y)
    (FEC Z)
```

```
) B, L3
) F, pop
```

**One possible utilization of FEC**

FTP

Application Traffic

E-mail

Web Browsing

Voice

**FEC Classification**

Ingress LSR

FEC A L34

FEC B L45

FEC C L07

X                                                                   Z

```
C, L22
C, L23
C, L24
C, L25
```

```
D, pop
C, L17
D, L18
C, L19
```

MPLS                                                                52

---

26

# MPLS FEC for the grid



FTP

Egress
Egress

Egress

Egress

Egress
Egress

Egress

| | FEC A: time constraint applications |
| --- | --- |
| | FEC B: best effort traffic |

53

---

# Label & FEC

❑ Independent LSP control
  ❑ An LSR binds a label to a FEC, whether or not the LSR has received a label from the next-hop for the FEC
  ❑ The LSR then advertises the label to its neighbor

❑ Ordered LSP control
  ❑ An LSR only binds and advertises a label for a particular FEC if:
    • it is the egress LSR for that FEC or
    • it has already received a label binding from its next-hop

54

# Label Distribution Protocols

❏ LDP
  - Maps unicast IP destinations into labels
❏ RSVP-TE, CR-LDP
  - Used in traffic engineering
❏ BGP
  - External labels (VPN)
❏ PIM
  - For multicast states label mapping

# MPLS for resiliency
## MPLS FastReroute

❏ Intended to provide SONET/SDH-like healing capabilities
❏ Selects an alternate route in tenth of ms, provides path protection
❏ Traditional routing protocols need minutes to converge!
❏ FastReroute is performed by maintaining backup LSPs

# MPLS for resiliency, con't
## Backup LSPs

❑ One-to-one
❑ Many-to-one: more efficient but needs more configurations

C      G

A

D      F      H

B      One LSP from A to H
       3 backup LSPs
E

# MPLS for resiliency, con't
## Recovery on failures

❑ Suppose E or link B-E is down...
❑ B uses detour around E with backup LSP

C      G

A

D      F      H

B
E

# MPLS for VPN
## (Virtual Private Networks)

❑ Virtual Private Networks: build a secure, confidential communication on a public network infrastructure using routing, encryption technologies and controlled accesses

TOP SECRET

Shared network

MPLS

59

# MPLS for VPN, con't
## The traditional way of VPN

❑Uses leases lines, Frame Relay/ATM infrastructures...

MPLS

60

# MPLS for VPN, con't
## VPN over IP/MPLS

❑ IP/MPLS replace dedicated  networks

❑ MPLS reduces VPN complexity by reducing routing information needed at provider's routers



```
134.13/16    S0
134.15/16    s0
140.11.45/8  s1
```

Do not know VPNs at all

134.15/16

134.13/16

VPN B

VPN B

VPN A

VPN A

**Ingress LSR**

**Egress LSR**

VPN B

**backbone**

140.11.45/8

IP          MPLS          IP

MPLS

61

---

# MPLS for optical networks
## Before MPLS



| Application | | | Application |
| Transport | | | Transport |
| Network | Network | Network | Network |
| Link | WDM | WDM | Link |

Terminals        IP router        IP router        Terminals

**Source J. Wang, B. Mukherjee, B. Yoo**

MPLS

62

# MPLS for ON, con't
## MPλS=MPLS+λ lightpath



Application

Transport

Network

Link

Terminals

Optical Label Switch

λ is viewed as a label

λ
Routing Control

Fabric

$\lambda_1 \lambda_2 \ldots \lambda_n$

$\lambda_1 \lambda_2 \ldots \lambda_n$

$\lambda_1 \rightarrow \lambda_2$

$\lambda_1 \lambda_2 \ldots \lambda_n$

$\lambda_1 \lambda_2 \ldots \lambda_n$

Application

Transport

Network

Link

Terminals

MPLS

63

---

# MPLS for ON, con't
## GMPLS

❑ GMPLS stands for "Generalized Multi-Protocol Label Switching"

❑ Extends the concept of MPLS beyond data networks to address legacy transport networks

❑ Reduce OPEX cost for operators

❑ A suite of protocols that provides a common set of control functions for disparate transport technologies (IP, ATM, SONET/SDH, DWDM)

❑ Hot issue at IETF!

MPLS

64

# MPLS for ON, con't
## GMPLS control plane

| | |
|---|---|
| **LINK MANAGEMENT:** Link Management Protocol (LMP) | -Neighbor discovery<br>-Maintain control channel connectivity<br>-Verify data link connectivity<br>-Correlate link property information<br>-Suppress downstream alarms<br>-Localize link failures |
| **ROUTING:** Open Shortest Path First-Traffic Engineering (OSPF-TE) | -Distribute TE link information<br>-Advertise nodes in the network and create topology<br>-Calculate constrained shorted path (CSPF)<br>-Routing information for control and data plane |
| **SIGNALING:** Resource ReserVation Protocol-Traffic Engineering (RSVP-TE) | -Signals setup/teardown/refresh of paths with QoS requirements (e.g., circuit size)<br>-Uses control channel to setup an optical LSP<br>-Supports refresh reduction<br>-Supports Explicit Route Object (ERO) and Record Route Object (RRO) |

**Source S. Kinoshita, R. Rabbat, APNOMS 2005**

MPLS

65

---

# Ex: Service Provisioning

Typical service provisioning          With GMPLS

| | Typical service provisioning | With GMPLS |
|---|---|---|
| CLIENT | Contract | Signaling |
| SALES | Contract handling | UNI |
| ADMINISTRATION | | |
| PROJECT MNGT | Create work packages | Call control → OK? |
| NETWORK OPERATION LOCAL DOMAIN | Plan | Call control |
| NETWORK OPERATION LOCAL DOMAIN | Plan | RSVP signaling |
| EXTERNAL SUPPLIER | Plan   Install   Configure | RSVP signaling |

PATH          RESERV

Knowledge of each technology domain!

Provisioing of each layer!

Separate mngt of per-domain operations functions!

From Pascalini et al., IEEE Comm. Mag. July 2005

MPLS

66

# DiffServ over (G)MPLS
## map DiffServ class on MPLS FEC



MPLS

67

# Some words on inter-domain



MPLS

68

Summary
Towards IP/(G)MPLS/DWDM



Summary
Technology scope

# Want to know more?

- GMPLS: IEEE Comm. Mag., Vol. 43(7), July 2005
- Optical Control Plane for the Grid Community: IEEE Comm. Mag., Vol. 44(3), March 2006.
- "Optical Transport Systems/Networks" by S. Kinoshita & R. Rabbat, APNOMS 2005. http://www.apnoms.org/2005/tutorial/Tutorial%202.pdf
- « Inter-domain Traffic Engineering for QoS-guaranteed DiffServ Provisioning », Young-Tak Kim, APNOMS 2005. http://www.apnoms.org/2005/tutorial/Tutorial%203.pdf
- See Tutorial IV of HOTI 2006: Dynamic Optimal Networks for Grid Computing

71

---

End of part 1, go to part 2

72