# New Internet and Networking Technologies and
# Their Application on Computational Sciences

COSCI 2004

Dai Hoc Bach Khoa, Ho Chi Minh City,
Vietnam, March 3, 2004

C. Pham

University Lyon, France

LIP (CNRS-INRIA-ENS-UCBL)

# Computational Sciences

❑ Use of computers to solve complex problems
  - ❑ Modelling techniques
  - ❑ Simulation tehniques
  - ❑ Analytic & Mathematic methods
  - ❑ …

❑ Large problems require huge amount of processing power: supercomputers, high-performance clusters, etc.

# Earth simulator: #1 TOP500

©JAMSTEC

- Intensive numerical simulations
- Ex: Super High Resolution Global Atmospheric Simulation



**Super High Resolution Global Atmospheric Simulation**
AFES T1279L96
Precipitation [mm/hour] 03 SEP/15 15Z

# A large variety of applications

**Astrophysics:**
Black holes,
neutron stars,
Supernovae…

**Mechanics:**
Fluid dynamic,
CAD, simulation.

**High-Energy Physics:**
Fundamental particles of matter,
Mass studies…

**Chemistry&biology:**
Molecular simulations,
Genomic simulations…

# This talk is about…

❑ How the Internet revolution could be beneficial to computational sciences

From isolated resources to Internet-based resources



Pre-PC

PC

WorkStation

Super Computer

Clusters

Computational Grids

It's not a talk about grids, see www.ggf.org for pointers on grid computing

# The big-bang of the Internet

# # Internet host



Hobbes' Internet Timeline Copyright ©2003 Robert H Zakon
http://www.zakon.org/robert/internet/timeline/

| DATE | HOSTS | | DATE | HOSTS |
|------|-------|---|------|-------|
| 12/69 | 4 | | 05/82 | 235 |
| 06/70 | 9 | | 08/83 | 562 |
| 10/70 | 11 | | 10/84 | 1,024 |
| 12/70 | 13 | | 10/85 | 1,961 |
| 04/71 | 23 | | 02/86 | 2,308 |
| 10/72 | 31 | | 11/86 | 5,089 |
| 01/73 | 35 | | 12/87 | 28,174 |
| 06/74 | 62 | | 07/88 | 33,000 |
| 03/77 | 111 | | 10/88 | 56,000 |
| 12/79 | 188 | | 07/89 | 130,000 |
| 08/81 | 213 | | 10/89 | 159,000 |

■ New Survey
◆ Old Survey

# www.web-the-big-bang.org

Hobbes' Internet Timeline Copyright ©2004 Robert H Zakon
http://www.zakon.org/robert/internet/timeline/

| DATE | SITES | | DATE | SITES |
|------|-------|---|------|-------|
| 12/90 | 1 | | 12/93 | 623 |
| 12/91 | 10 | | 06/94 | 2,738 |
| 12/92 | 50 | | 12/94 | 10,022 |
| 06/93 | 130 | | 06/95 | 23,500 |
| 09/93 | 204 | | 01/96 | 100,000 |
| 10/93 | 228 | | 06/96 | 252,000 |

# The Internet in Vietnam

| Year | 2000 | 2001 | 2002 | oct2003 |
|---|---|---|---|---|
| Subscribers | 103,751 | 166,616 | 350,000 | 650,654 |
| Users | 430,000 | 700,000 | 1.4 mil | 2.6 mil |
| Penetration rate | 0.5% | 0.9% | 1.7% | 3.2% |

**Source: VNNIC**



**Jan 2004**

868,059
3.5 mil
4.31%

- 2600000 (2003)
- 1400000 (2002)
- 700000 (2001)
- 430000 (2000)

China ← 155 Mbps

Korea ← 2.5 Mbps

Japan → 2 Mbps

Hongkong ← 231.5 Mbps

Hongkong ← 155 Mbps

2.5Gbps

Singapore → 159 Mbps

**Total International traffics: 705 Mbps**
**Total backbone traffics: 2.5 Gbps**

1Mbps=1 million bits/s

LAN technology=10-100Mbps

# Internet usage: e-mail…

- ❑ Convenient way to communicate in an informal manner

- ❑ Attachments as a easy way to exchange data files, images…

# ...and surfing the web

- ❑ A true revolution for rapid access to information
- ❑ Increasing number of apps:
  - ❑ e-science,
  - ❑ e-commerce, B2B, B2C,
  - ❑ e-training, e-learning,
  - ❑ e-tourism
  - ❑ ...

# Towards all IP

# A whole new world for IP

# The optical revolution



CPU Processing Power
2x / 18 months

From McKeown

Link Speed
2x / 7 months

TDM    DWDM

Demand: about 111 million km of cabled optical fiber / year

# DWDM, bandwidth for free?

DWDM: Dense Wavelength Division Multiplexing



10Gbps

2Gbps

2.5, 10, 40 Gbps are available!

From Computer Desktop Encyclopedia
Reproduced with permission.
© 2001 Metromedia Fiber Network

# The information highways

Truck of tapes
5PByte

DWDM

1600 Gbyte/s

320λ
40Gbps

## Revisiting the truck of tapes

**Consider one fiber**

- Current technology allows for 320 λ in one of the frequency bands

- Each λ has a bandwidth of 40 Gbit/s

- Transport: $320 * 40*10^9 / 8 = 1600$ GByte/sec

- Take a 10 metric ton truck

- One tape contains 50 Gbyte, weights 100 gr

- Truck contains $( 10000 / 0.1 ) * 50$ Gbyte = 5 PByte

- Truck / fiber = 5 PByte / 1600 GByte/sec = **3125 s ≈ one hour**

- For distances further away than a truck drives in one hour (50 km)
  minus loading and handling 100000 tapes **the fiber wins!!!**

Example from A. Tanenbaum, slide from Cees De Laat

# Fibers everywhere?



offices

residentials

FTTH
FTTC

10Gbps

metro ring

Internet
Data
Center

Network Provider
2.5Gbps

2.5Gbps

10Gbps

Network Provider

campus

1Gbps
GigaEth

Core
40Gbps

# High Performance Routers



PRO/8812

©cisco

©Procket Networks

©Lucent

©Alcatel

©Nortel Networks

**and more…**

# Operator's infrastructure

❑ Backbones are optical: OC48 (2.5Gbps), OC192 (10Gbps), OC768 (40Gbps) soon

❑ New technologies deployed by operators, POPs available worldwide

# In a near future?

2.5 Gbps

2.5 Gbps

2.5 Gbps

2.5 Gbps

2.5 Gbps

2.5 Gbps

10 Gbps

2.5 Gbps

# New applications on the information highways

## Think about...

- video-conferencing
- video-on-demand
- interactive TV programs
- remote archival systems
- tele-medecine
- virtual reality, immersion systems
- high-performance computing, grids
- distributed interactive simulations

# Computational grids

## user application



1PFlops

Virtually unlimited resources

# High Energy Physics at CERN



**LEP**

**LHC**

**CMS**
Compact
Muon
Solenoid

**ATLAS**

Images from EDG (DataGrid) project

**3.5 Petabytes/year $\approx 10^9$ events/year**

# Distributed Databases



**1 TIPS = 25,000 SpecInt95**

**PC (1999) = ~15 SpecInt95**

Online System

~100 MBytes/sec

~TBytes/sec

*Bunch crossing per 25 nsecs.*
*100 triggers per second*
*Event is ~1 MByte in size*

Offline Farm
~20 TIPS

~100 MBytes/sec

**Tier 0**

CERN Computer
Center > ~20 TIPS

HPSS

~622
Mbits/sec
or Air Freight

**Tier 1**

**~ 4 TIPS**

France Regional
Center

HPSS

UK Regional
Center

HPSS

Italy Regional
Center

HPSS

Fermilab

HPSS

• • •

~2.4 Gbits/sec

**Tier 2**

Tier2 Center

Tier2 Center

Center

2 Center

2 Center

~622 Mbits/sec

**Tier 3**

Institute
~0.25 TIPS

tute

stitute

Institute

Physics data cache

100 - 1000
Mbits/sec

*Physicists work on analysis "channels".*

*Each institute has ~10 physicists*
*working on one or more channels*

*Data for these channels should be*
*cached by the institute server*

## Large data transfers
## require high bandwidth

# Wide-area interactive simulations



display

computer-based
plane simulator

$(x,y,z,t)$

INTERNET

airport simulator

**Interactive applications
require low latencies**

human in the loop
flight simulator

# Limitations of the current Internet

- Bandwidth
    - Raw bandwidth is not a problem: DWDM
    - Provisioning bandwidth on demand is more problematic
- Latency
    - Mean latencies on Internet is about 80-160ms
    - Bounding latencies or ensuring lower latencies is a problem
- Loss rate
    - Loss rate in backbone is very low
    - End-to-End loss rates, at the edge of access networks are much higher
- Communication models
    - Only unicast communications are well-defined: UDP, TCP
    - Multi-parties communication models are slow to be deployed

# New technologies addressed in this talk

❑ More Quality of Service: Differentiated Services, who pays more gets more!

❑ Bandwidth provisioning: MPLS for virtual circuit in the core networks

❑ Multicast: enhancing the communication model

# Revisiting the *same service for all* paradigm

IP packet

No delivery guarantee

INTERNET

Regular mail

## Enhancing the best-effort service

IP packet

URGENT

Introduce
Service Differentiation

**Prioritaire**
First Class
Letter **A**

FRAGILE

# 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

# 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

# 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, then evolve
  - ❑ Build simple system at first, expand if needed in future

# IP implementation: DiffServ

RFC 2475

No per flow state in the core

IP packet

Flow 1
Flow 2
Flow 3
Flow 4
…

10Gbps=2.4Mpps
with 512-byte packets

**Stateful approaches**
**scalable**
**at gigabit rates**

1981    1993    1997

IP TOS

IntServ/
RSVP

DiffServ

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

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

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

| DSCP | | CU | | Padding |

# Traffic Conditioning

# 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

# 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

# Premium Service Example



Drop always

10Mbps

Fixed Bandwidth

source Gordon Schaffee

# Assured Service Example



Drop if congested

10Mbps

Assured Service

Uncongested

Congested

source Gordon Schaffee

# 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

**Assured Service**

Token Bucket

| 0<br>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 |

Packet Input → **Test if token**

No token

Token → **Set A-bit** → Data Queue → Packet Output

source Gordon Schaffee, modified by C. Pham

# 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--

Packets Out

RED In/Out Queue Management

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

source Gordon Schaffee, modified by C. Pham

# Practical realization



Drop probalility

WRED Queue 0

WRED Queue 1

1/4    1/2    3/4    Queue filling

30 %
30 %
30 %
10 %

Queue 0
Queue 1
Queue 2
Queue 3

Classifier

Source VTHD

| Prec. 0 | BE + AF UDP out profile |
| Prec. 1 | AF UDP in profile |
| Prec. 2 | AF TCP out profile |
| Prec. 3 | AF TCP in profile |
| Prec. 4 | |
| Prec. 5 | EF |
| Prec. 6 | Control |
| Prec. 7 | Control |

# DiffServ for grids



Wide-area interactive simulations

FTP

scheduling

Ingress/Ingress

Egress

marking

Egress

Egress

Egress

Egress

Assured Forwarding

Premium

# 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

# Bandwidth provisioning

- 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 |

# Provider's view



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

# 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

**IP is connectionless!**

# Virtual circuit explained

### Connections & Virtual circuits table

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

label

R3

23

78

45

34

Link 1

Link 2

Link 3

Link 4

R3

R1

R4

A

B

C

D

E

Virtual Circuit Switching

R2

R5

# Why virtual circuit?

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

We're fast enough!

Now: Virtual circuits for bandwidth provisioning!

# MPLS

❑ 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 operation

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

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

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

Label Switch Router

IP

link a

IP 10

IP 20

IP 40

IP

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

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

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

Source Yi Lin, modified C. Pham

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

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

LSR A
LSR B
LSR C
LSR D
LSR E
LSR F

L10
L5
L14
L19

X
Y
Z

# Forwarding Equivalent Class

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

**One possible utilization of FEC**

Table A

L6: (FEC F)
L8: (FEC X)
    (FEC Y)
    (FEC Z)

X

B, L3
F, pop

Z

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

FTP

Application
Traffic

E-mail

Web
Browsing

Voice

**FEC
Classification**

Ingress
LSR

FEC A
L34

FEC B
L45

FEC C
L07

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

# MPLS & VPN

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

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

TOP SECRET

# MPλS: MPLS+optical



| Application |
| ----------- |
| Transport |
| Network |
| Link |

Terminals

MP...

Application

...nsport

...work

...ink

...inals

λ is viewed as a label

**Optical Label Switch**

λ
Routing Control

$\lambda_1 \, \lambda_2 \, \ldots \, \lambda_n$

$\lambda_1 \, \lambda_2 \, \ldots \, \lambda_n$

Fabric

$\lambda_1 \rightarrow \lambda_2$

$\lambda_1 \, \lambda_2 \, \ldots \, \lambda_n$

$\lambda_1 \, \lambda_2 \, \ldots \, \lambda_n$

# Towards IP/MPLS/DWDM



From cisco

# Ex: MPLS circuits on grids



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

# Ex: MPLS FEC for the grid



Egress
Egress

Egress

Egress

Egress
Egress

FEC A: time constraint applications

FEC B: best effort traffic

# Unicast, the current (Internet) communication model

FTP

TCP

TCP

TCP

❑ There are applications that naturally need multi-destination communication model

❑ Collaborative works
❑ Visio-conferencing
❑ Software distribution

❑ Video-on-Demand
❑ Virtual Reality
❑ Distributed Simulation

# From unicast to multicast

# Multicast in example

## The user's perspective



224.2.0.1

Multicast IP address range
224.0.0.0 … 239.255.255.255

from UREC, http://www.urec.fr

━━ **Multicast address group 224.2.0.1**

# What's behind the scene?

domain

peering point

access router

Internet router

224.2.0.1

# IP multicast TODO list

- ✓ **Receivers must be able to subscribe to groups, need group management facilities**
- ✓ **A communication tree must be built from the source to the receivers**
- ✓ **Branching points in the tree must keep multicast state information**
- ✓ **Inter-domain routing must be reconsidered for multicast traffic**
- ✓ **Need to consider non-multicast clouds**

# Ex: Reliable multicast on grids

Data replications

Code & data transfers, interactive job submissions

Data communications for distributed applications (collective & gather operations, sync. barrier)

Databases, directories services

**224.2.0.1**

SDSC IBM SP
1024 procs
5x12x17 =1020

NCSA Origin Array
256+128+128
5x12x(4+2+2) =480

ENS cluster
48 nodes

**Multicast address group 224.2.0.1**

# Conclusions

❑ There's a lot more technologies going on that have impact on computational science

    ❑ Pure optical networks, broadband wireless

    ❑ Peer-to-Peer, Overlays

    ❑ Web services…

❑ The future will be all connected, all IP, anytime, anywhere, for more…