# Performance limiting issues over high latency & high-bandwidth networks

## Kei Hiraki

## University of Tokyo

# Computing System for real Scientists

- Fast CPU, huge memory and disks, good graphics
  - Cluster technology, DSM technology, Graphics processors
  - Grid technology
- Very fast remote file accesses
  - Global file system, data parallel file systems, Replication facilities

- Transparency to local computation
  - No complex middleware, or no small modification to existing software

- **Real Scientists are not computer scientists**

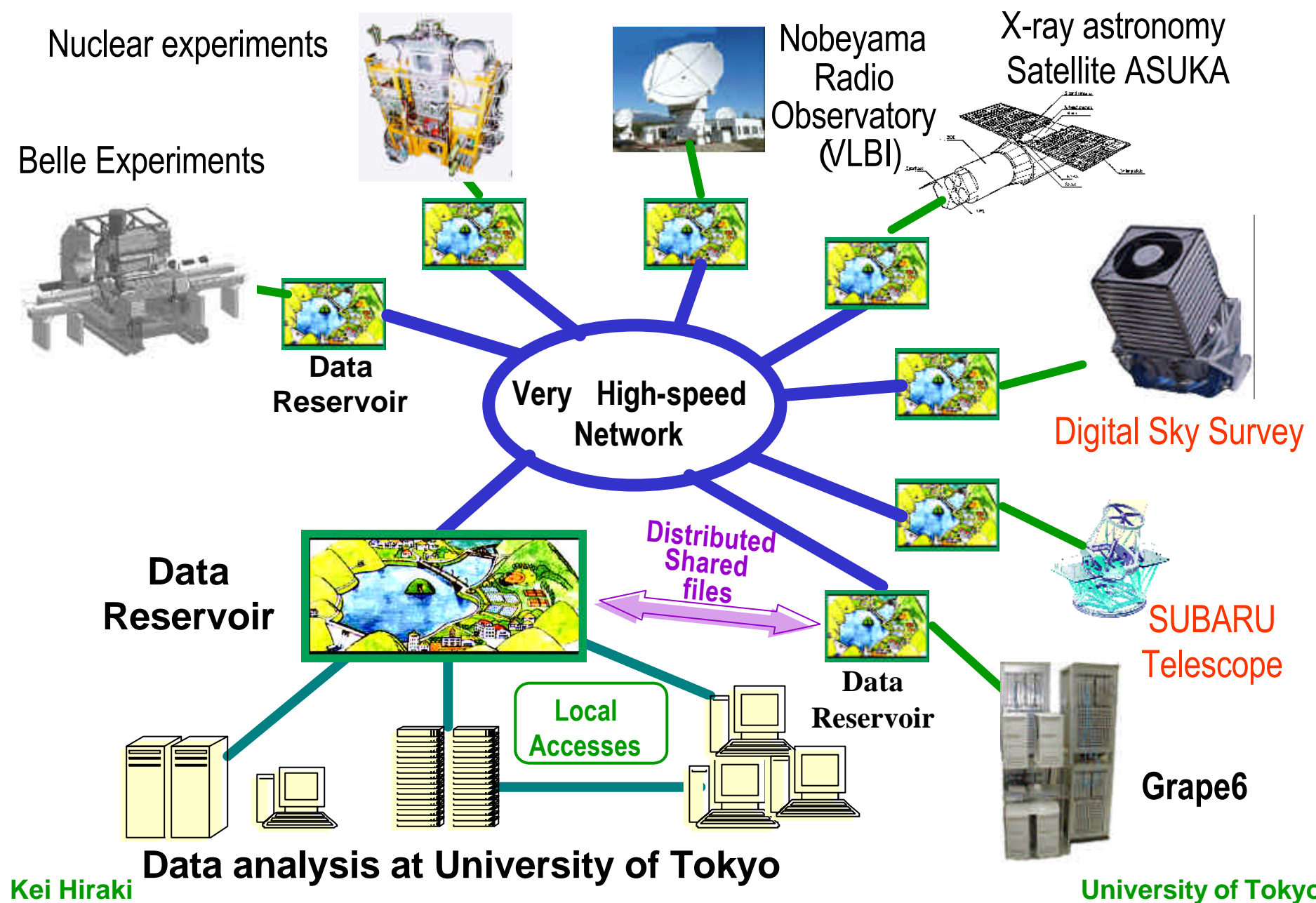- **Computer scientists are not work forces for real scientists**

# Data Reservoir

- **Sharing Scientific Data between distant research institutes**
  - Physics, astronomy, earth science, simulation data

- **Very High-speed single file transfer on Long Fat pipe Network**
  - > 10 Gbps, > 20,000 Km, > 400ms RTT

- **High utilization of available bandwidth**
  - Transferred file data rate > 90% of available bandwidth
    - Including header overheads, initial negotiation overheads

- **OS and File system transparency**
  - Storage level data sharing (high speed iSCSI protocol on stock TCP)
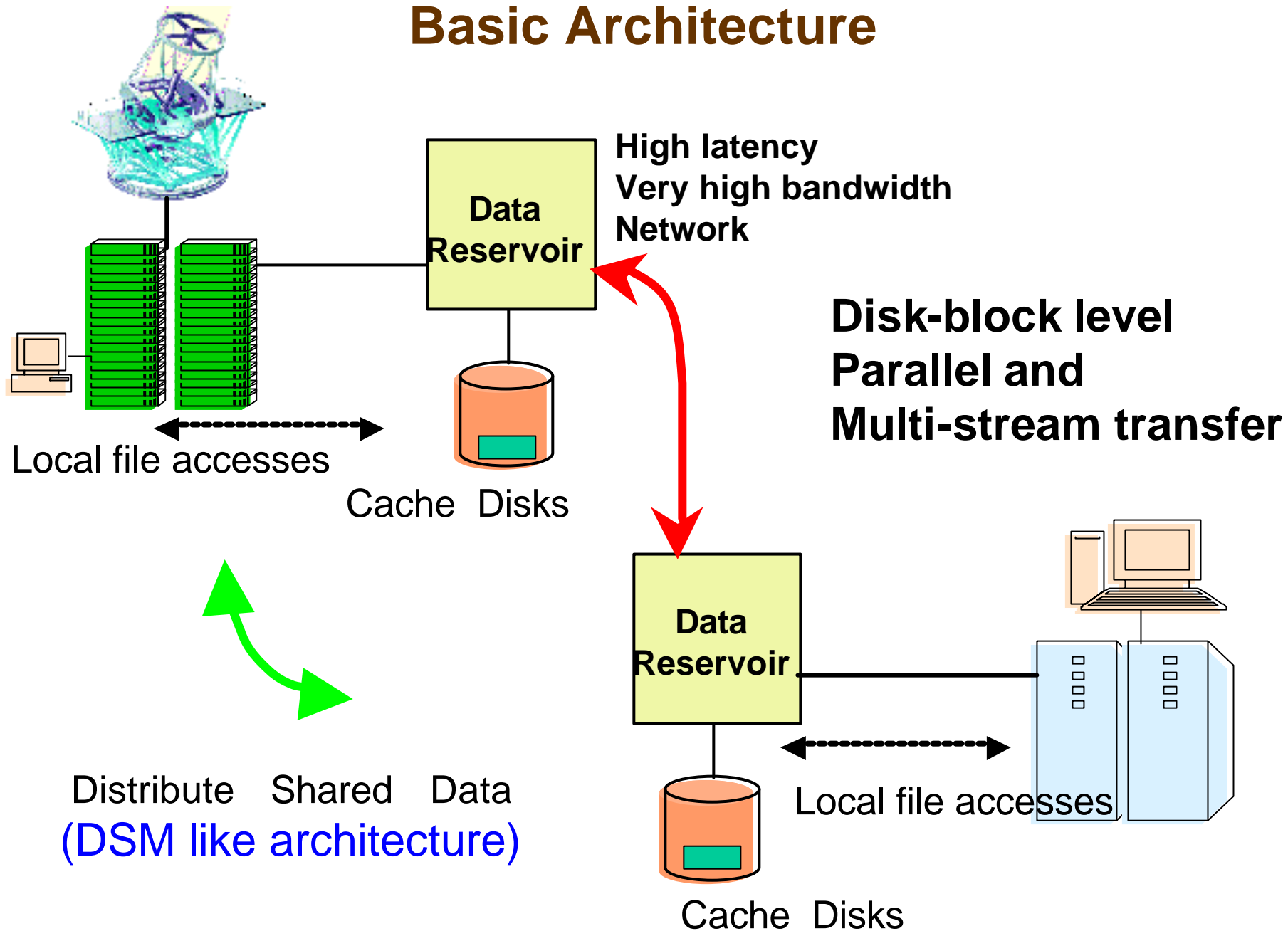  - Fast single file transfer

# GRAPE-DR

- GRAPE-DR:Very high-speed attached processor to a server
  - 2004 – 2008
  - Successor of Grape-6 astronomical simulator

- 2PFLOPS on 512 node cluster system
  - 1G FLOPS / processor
  - 512 processor / chip
  - 4 chips / PCI card
  - 2 PCI card / serer

  - 2 M processor / system

- Semi-general-purpose processing
  - N-body simulation, gridless fluid dyinamics
  - Linear solver, molecular dynamics
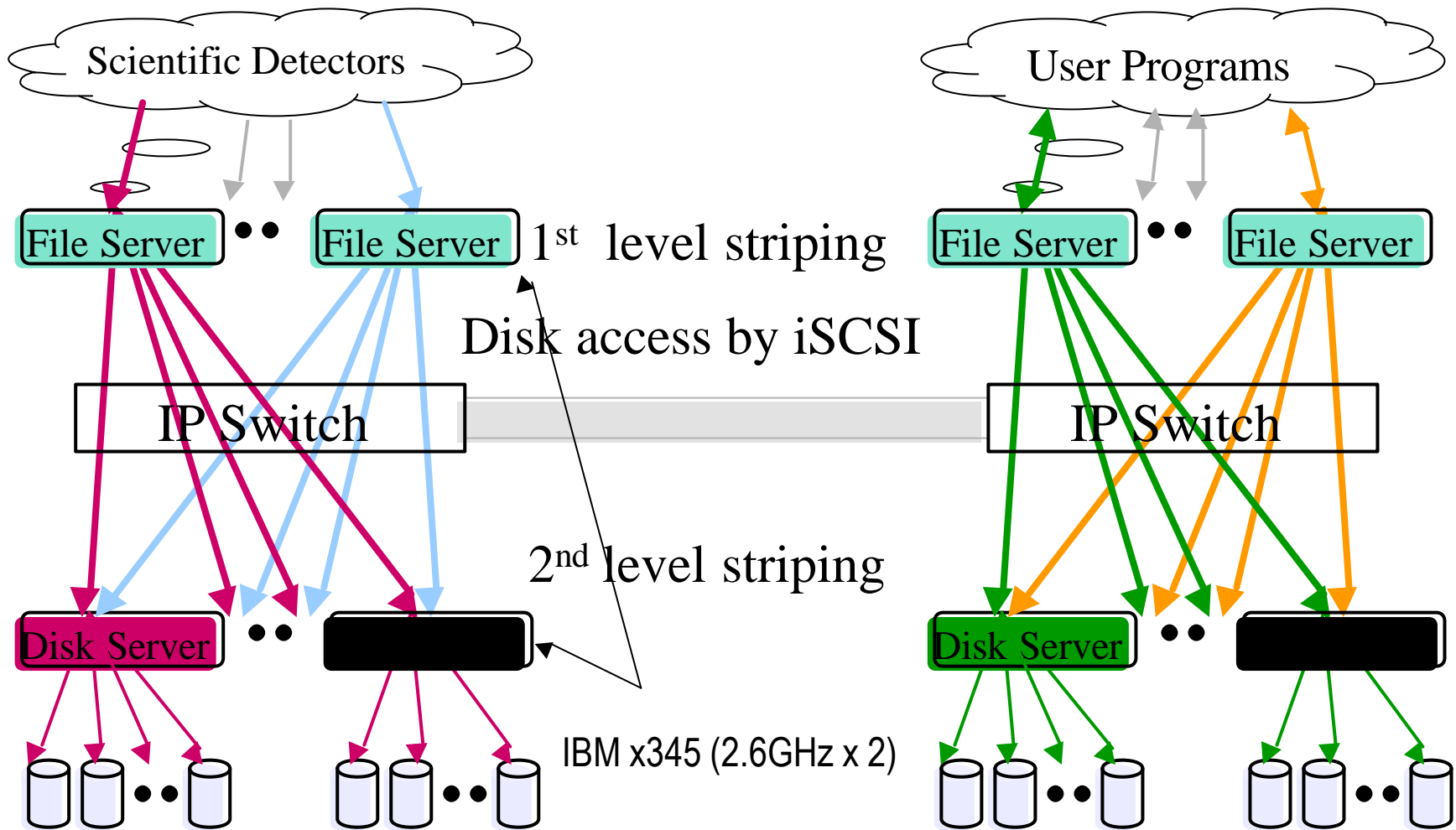  - High-order database searching (genome, protein data etc.)

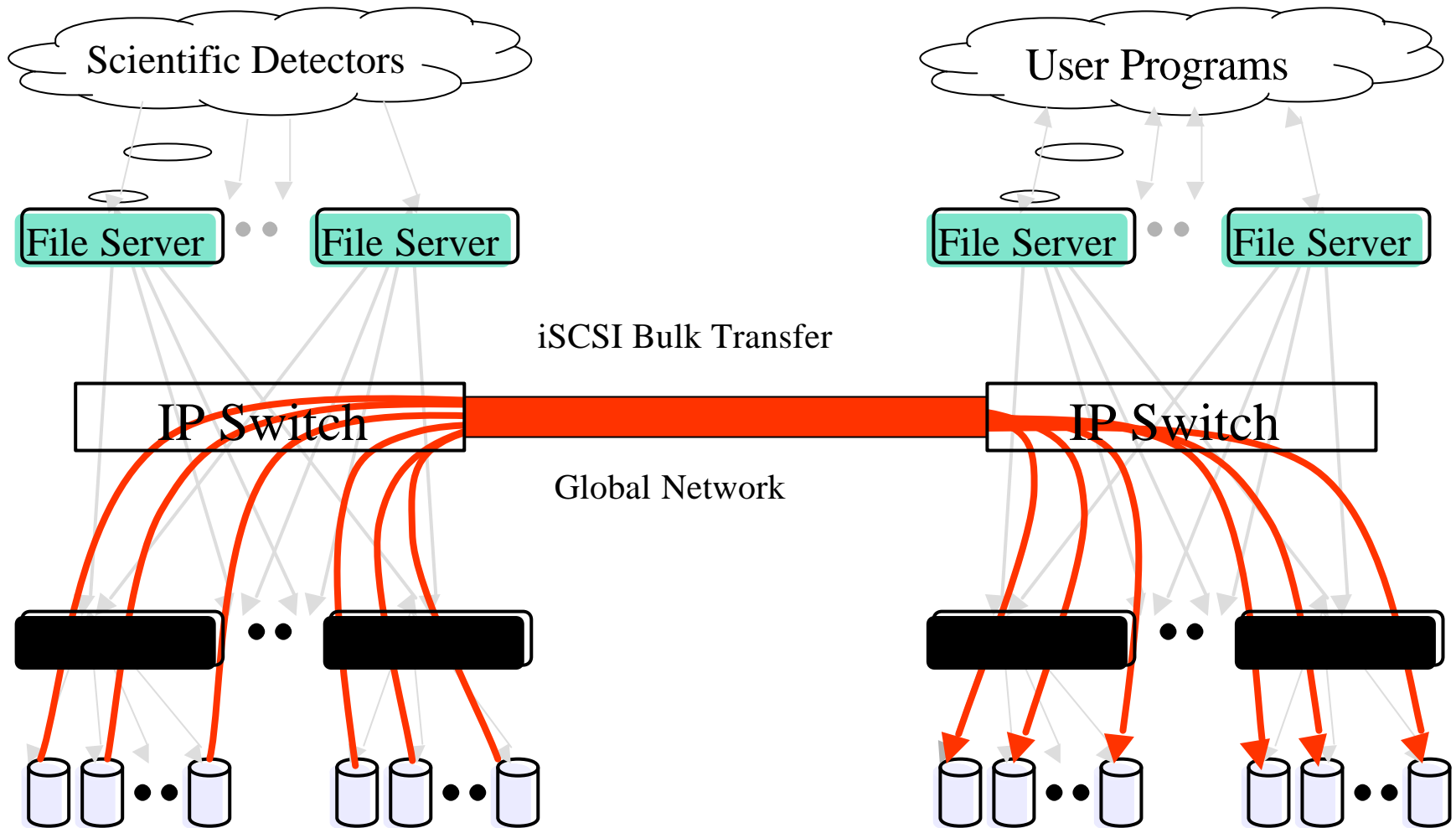# Data intensive scientific computation through global networks

Nuclear experiments

Belle Experiments

Nobeyama Radio Observatory VLBI)

X-ray astronomy Satellite ASUKA



**Data Reservoir**

**Very High-speed Network**

Digital Sky Survey

**Data Reservoir**

**Distributed Shared files**

**Local Accesses**

**Data Reservoir**

SUBARU Telescope

**Grape6**

**Data analysis at University of Tokyo**

Kei Hiraki

# Basic Architecture

**Data Reservoir**

**High latency Very high bandwidth Network**

**Disk-block level Parallel and Multi-stream transfer**

Local file accesses

Cache Disks

**Data Reservoir**

Local file accesses

Distribute Shared Data
(DSM like architecture)

Cache Disks

# File accesses on Data Reservoir



Scientific Detectors

User Programs

File Server ● ● File Server   1st level striping   File Server ● ● File Server

Disk access by iSCSI

IP Switch

IP Switch

2nd level striping

Disk Server ● ●

Disk Server ● ●

IBM x345 (2.6GHz x 2)

# Global Data Transfer

Scientific Detectors

User Programs

File Server • • File Server

File Server • • File Server

iSCSI Bulk Transfer

IP Switch

IP Switch

Global Network

# 1st Generation Data Reservoir

- **1st generation Data Reservoir (2001-2002)**
  - Efficient use of network bandwidth (SC2002 technical paper)
  - Single filesystem image (by striping of iSCSI disks)
  - Separation of local file accesses by users to remote disk to disk replication
  - Low level data-sharing (disk block level, iSCSI protocol)

- **26 servers for 570 Mbps data transfer**
  - High sustained efficiency ( 93
  - Low TCP performance
    - 21Mbps/server --- very slow
  - Unstable TCP performance

# Problems

- Low TCP bandwidth due to packet losses
  - TCP congestion window size control
  - Very slow recovery from fast recovery phase (>20min)

- Unbalance among parallel iSCSI streams
  - Packet scheduling by switches and routers
  - User and other network users have interests only to total behavior of parallel TCP streams

- Unstable network behavior from application soft

# Our starting point (1)

- **Fast Ethernet vs. GbE**
    - ? Iperf in 30 seconds
    - ? Min/Avg: Fast Ethernet > GbE

# Our starting point (2)

- **Delay emulator v.s. actual network**
  - Performance on a delay emulator is much better than actual network



Delay emulator

Clock level accuracy
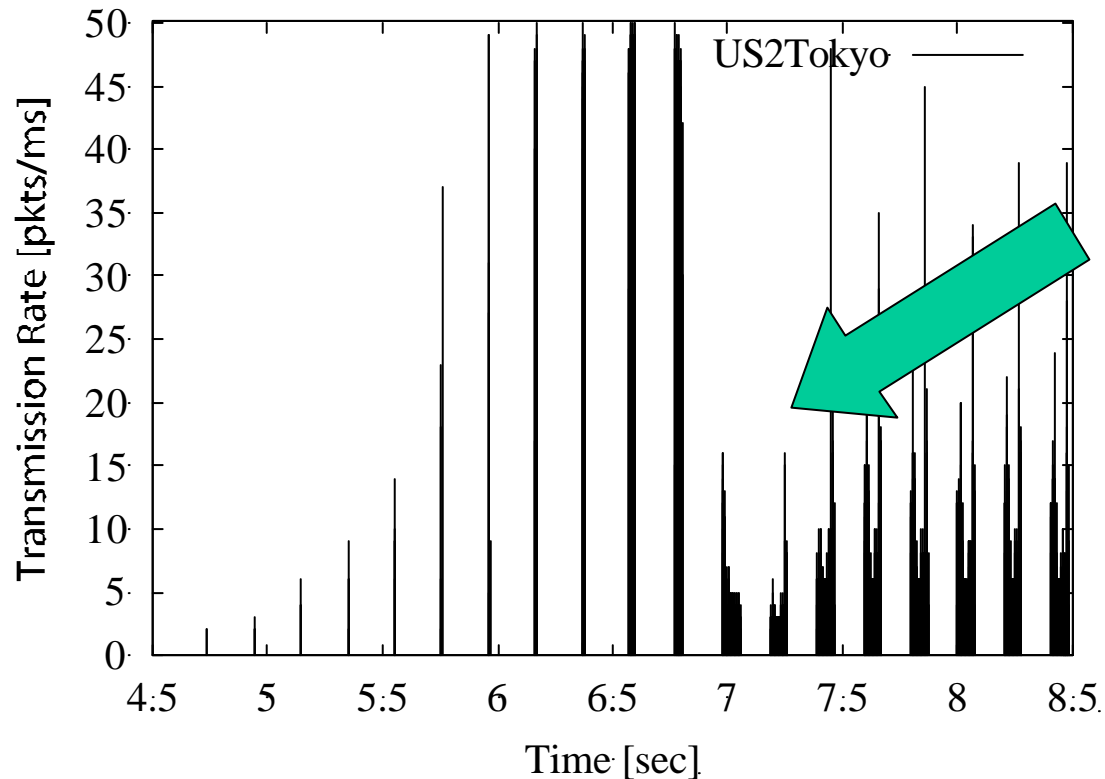
Real network
500ms RTT



SC|05 network path

# Packet Transmission Rate

? Bursty behavior

? Transmission in 20ms against RTT 200ms
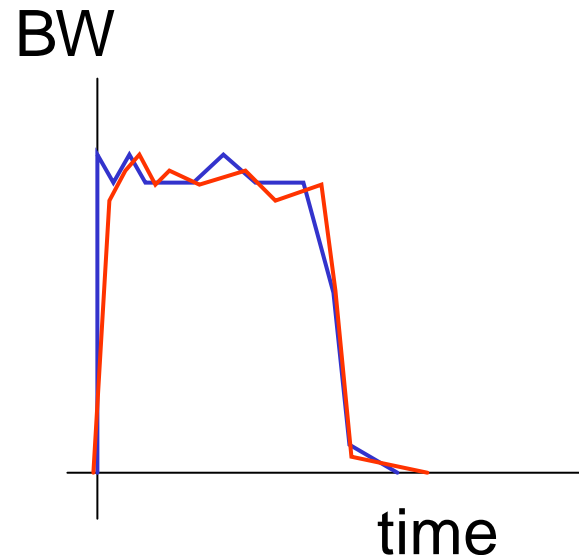
? Idle in rest 180ms

Packet loss occurred

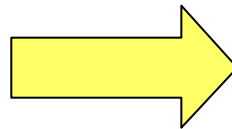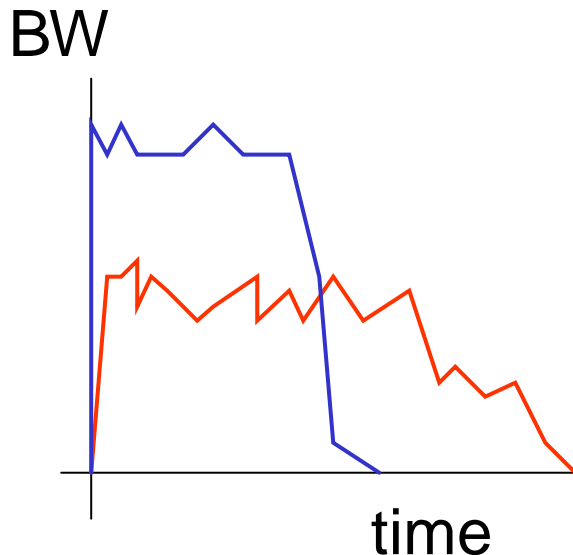# Unbalance within parallel TCP streams

? Unbalance among parallel iSCSI streams
- ? Packet scheduling by switches and routers
- ? Meaningless unfairness among parallel streams
- ? User and other network users have interests only to total behavior of parallel TCP streams
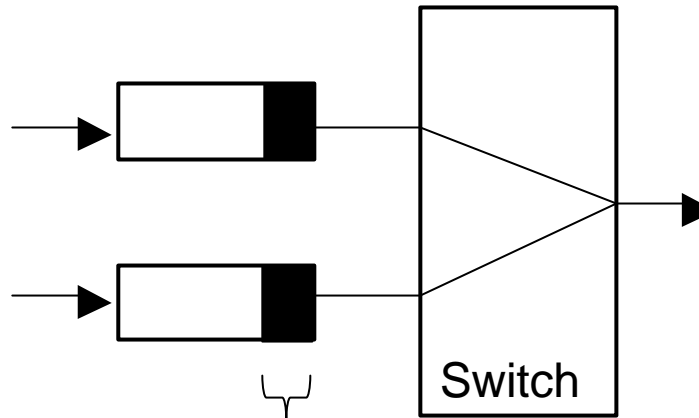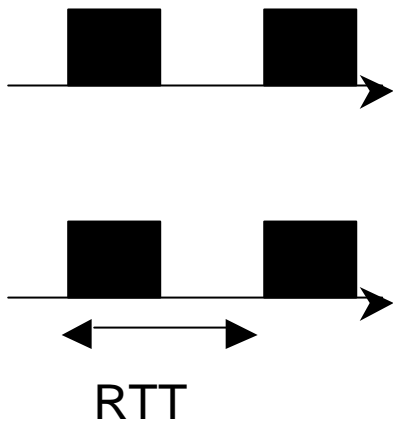
? Our approach
- ? Constant $cwnd_i$ for fair TCP network usage to other users
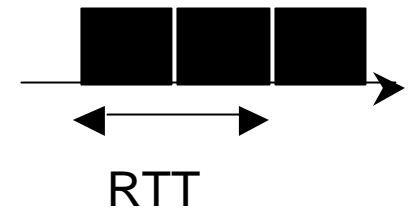- ? Balance each $cwnd_i$ communicating between parallel TCP streams

# Merging TCP streams

- Worst case is 2 TCP streams -> 1 TCP stream
- Bursty traffic just after slow start is the worst case

Traffic = ½ BW

Traffic = BW

Switch

RTT

RTT

Buffer size = ¼ BW   RTT   Network Utilization
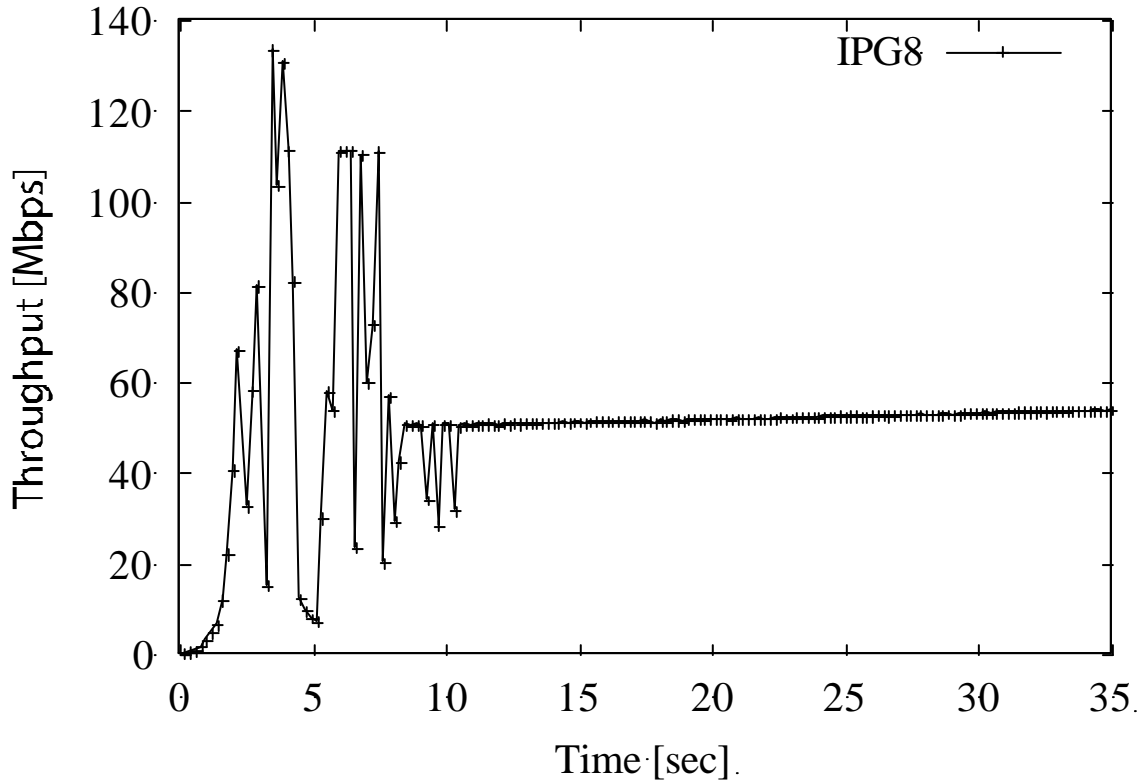
10 Gbps, 200msRTT      62.5 MB * Utilization factor

Many switches do not have large buffer
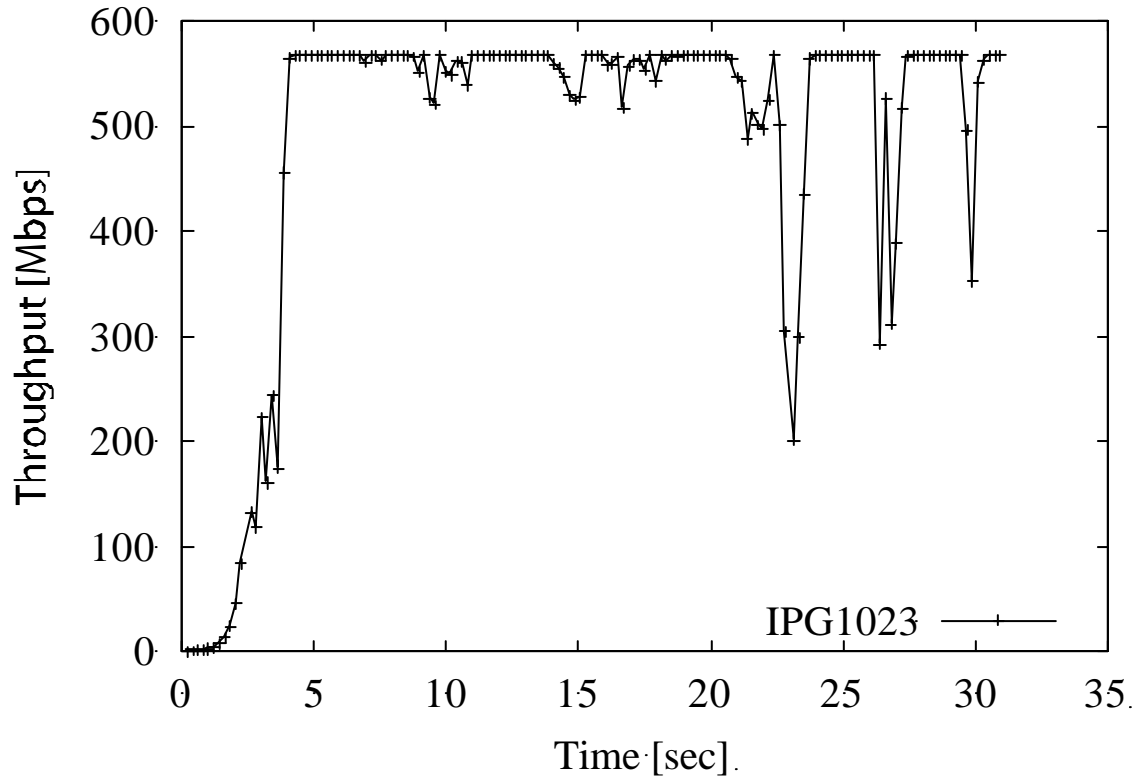Flow control cannot apply at intermediate switches

# Example Case of 8 IPG

? Success on Fast Retransmit

? Smooth Transition to Congestion Avoidance

? CA takes 28 minutes to recover to 550Mbps

# Best Case of 1023B IPG

? Like Fast Ethernet case

? Proper transmission rate

? Spurious Retransmit due to Reordering

# Buffer size problem



Kei Hiraki

University of Tokyo

# Sender side pacing

- BW bottleneck in the middle of network



- Too small buffer size of edge switch and trans-ocean gateway switch



Edge switch                Trans-ocean gateway

- Coordination between MAC and TCP layer (variable packet pace)
  - Avoid unnecessary packet loss at slow start phase

# Difficulty (1)

- Artificial packet losses
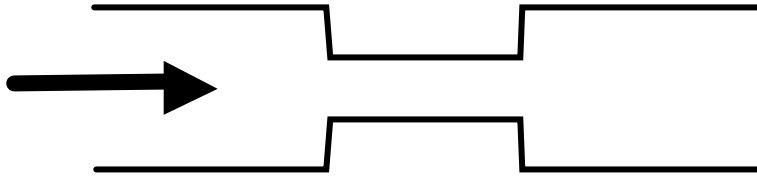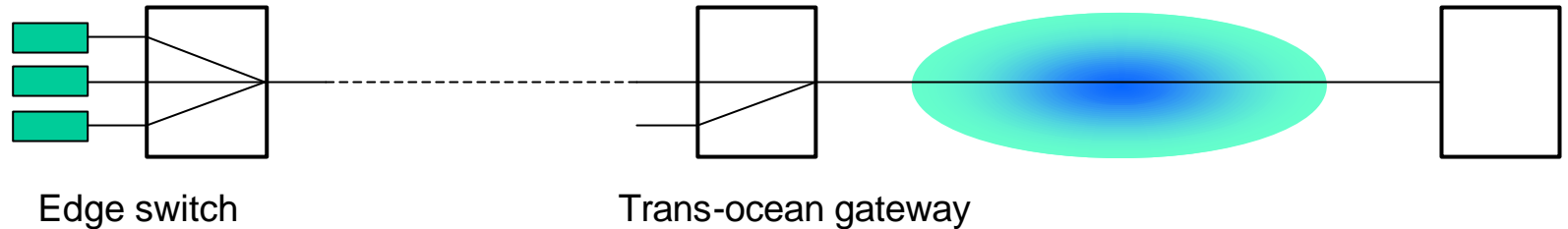  - Average
- Bursty behavior of TCP traffic
  - Improvement in TCP or TCP-like protocol is not effective
  - Pacing works in some situation

- Unbalance between parallel TCP streams
  - Application performance may decided by the worst stream
  - No good balancing protocol below network layer

- Too small buffer size at switches
  - Merging TCP streams cause artificial packet losses
  - Minimum buffer size > ¼ RTT * BW is necessary

# 3rd Generation Data Reservoir

- **<u>Hardware and software basis for 100Gbps Distributed Data-sharing systems</u>**


- 10Gbps disk data transfer by a single Data Reservoir server
- Transparent support for multiple filesystems (detection of modified disk blocks)
- Hardware(FPGA) implementation of Inter-layer coordination mechanisms
- 10 Gbps Long Fat pipe Network emulator and 10 Gbps data logger

# Utilization of 10Gbps network

- ## A single box 10 Gbps Data Reservoir server

  - Quad Opteron server with multiple PCI-X buses (prototype, SUN V40z server)

  - Two Chelsio T110 TCP off-loading NIC

  - Disk arrays for necessary disk bandwidth

  - Data Reservoir software (iSCSI deamon, disk driver, data transfer maneger)

PCI-X bus

Chelsio
T110
TCP NIC

10GBASE-SR

10G
Ethernet
Switch

PCI-X bus

Chelsio
T110
TCP NIC

Quad Opteron
Server
(SUN V40z)
Linux 2.6.6

PCI-X bus

SCSI
adaptor

Ultra320SCSI

PCI-X bus

SCSI
adaptor

Data
Reservoir
Software

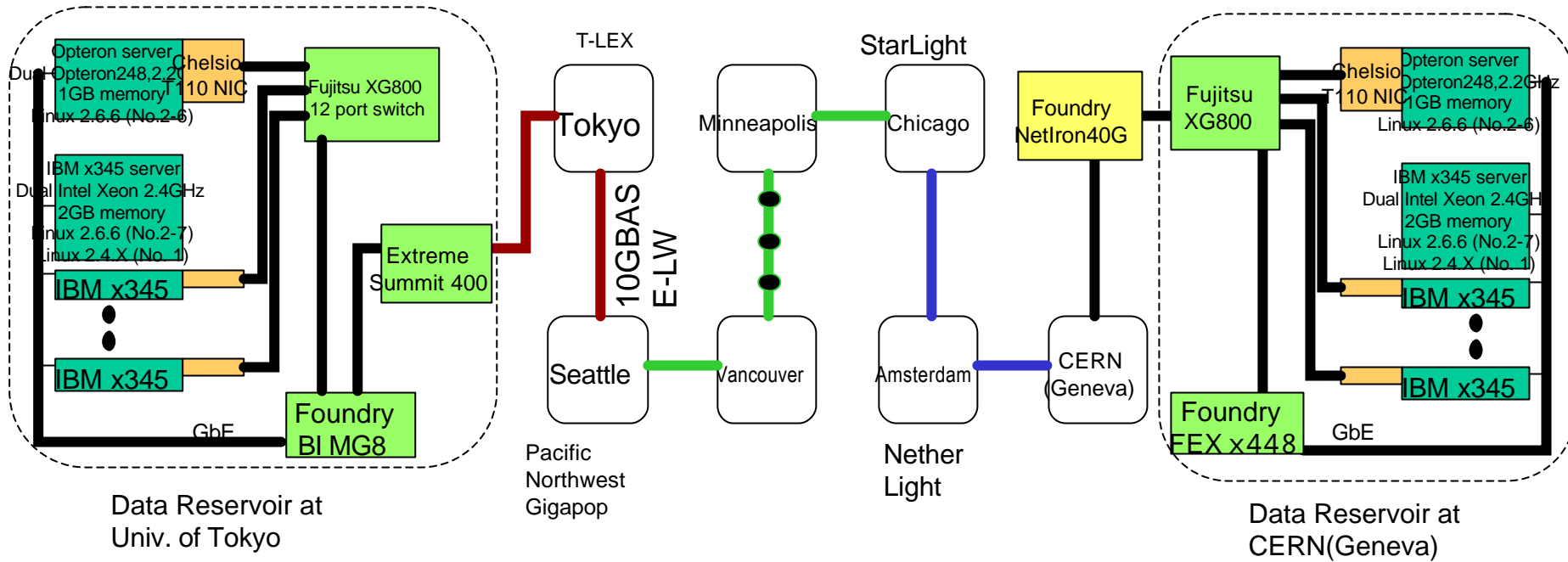Network used in the experiment

■ End Systems

□ A L1 or L2 switch

Tokyo-CERN Network connection

# Network topology of CERN-Tokyo experiment



**Data Reservoir at Univ. of Tokyo**

- Opteron server
- Dual Opteron248,2.2GHz
- 1GB memory
- Linux 2.6.6 (No.2-6)
- Chelsio T 110 NIC
- Fujitsu XG800 12 port switch
- IBM x345 server
- Dual Intel Xeon 2.4GHz
- 2GB memory
- Linux 2.6.6 (No.2-7)
- Linux 2.4.X (No. 1)
- IBM x345
- IBM x345
- Extreme Summit 400
- Foundry BI MG8
- GbE

**T-LEX**

- Tokyo
- Seattle
- 10GBASE-LW
- Pacific Northwest Gigapop

- Minneapolis
- Vancouver

**StarLight**

- Chicago
- Amsterdam
- CERN (Geneva)
- Nether Light

**Data Reservoir at CERN(Geneva)**

- Foundry NetIron40G
- Fujitsu XG800
- Chelsio T 10 NIC
- Opteron server
- Opteron248,2.2GHz
- 1GB memory
- Linux 2.6.6 (No.2-6)
- IBM x345 server
- Dual Intel Xeon 2.4GHz
- 2GB memory
- Linux 2.6.6 (No.2-7)
- Linux 2.4.X (No. 1)
- IBM x345
- IBM x345
- Foundry FEX
- GbE

Legend:
- ▬▬ WIDE / IEEAF
- ▬▬ CA*net 4
- ▬▬ SURFnet
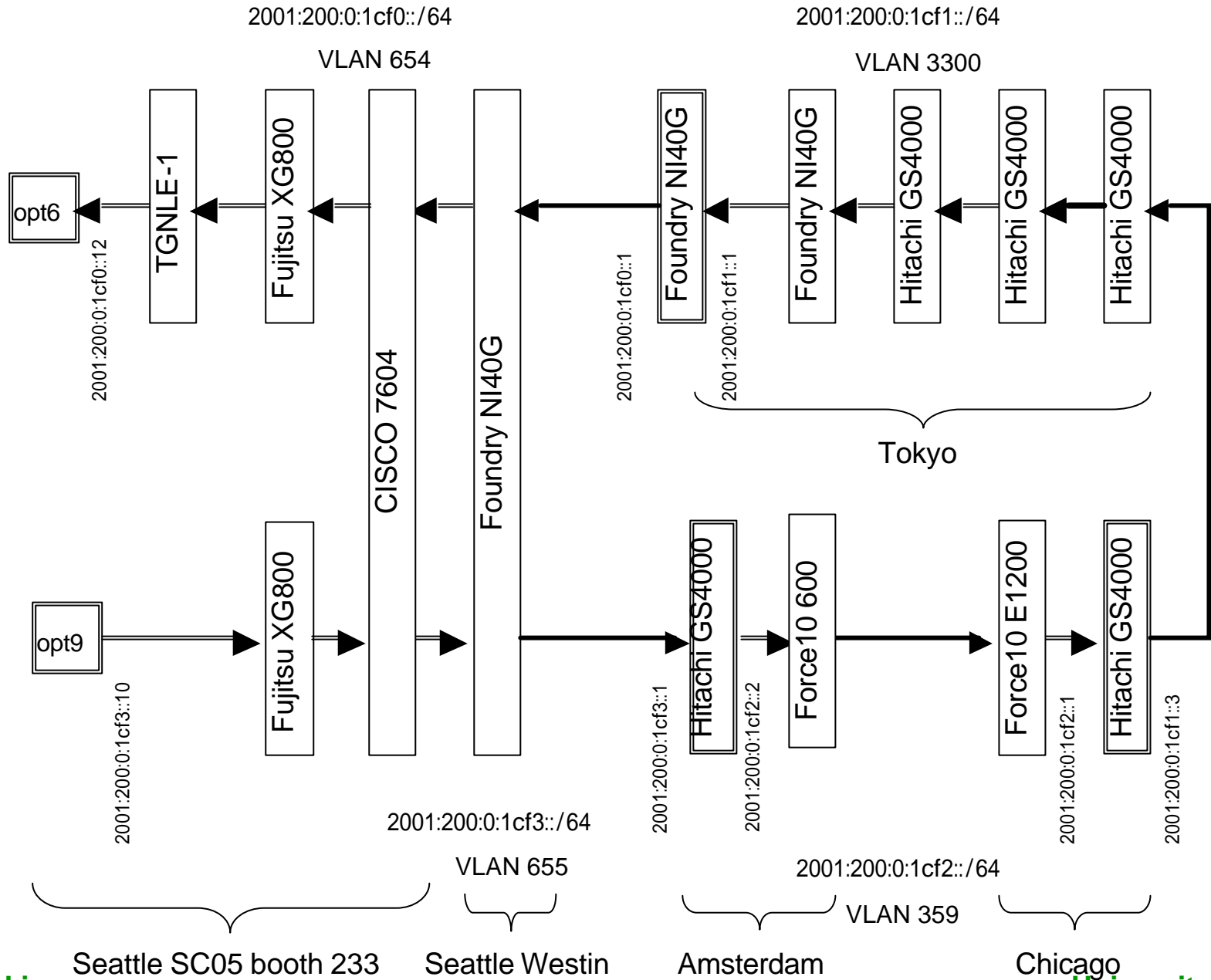
# Difficulty (2)

- ## Long-distance L2 network
  - Difficulty in debugging
  - Unstableness due to Spanning Tree algorithm
  - Long latency MAC address detection

- ## Our opinion
  - Avoid long distance L2 network
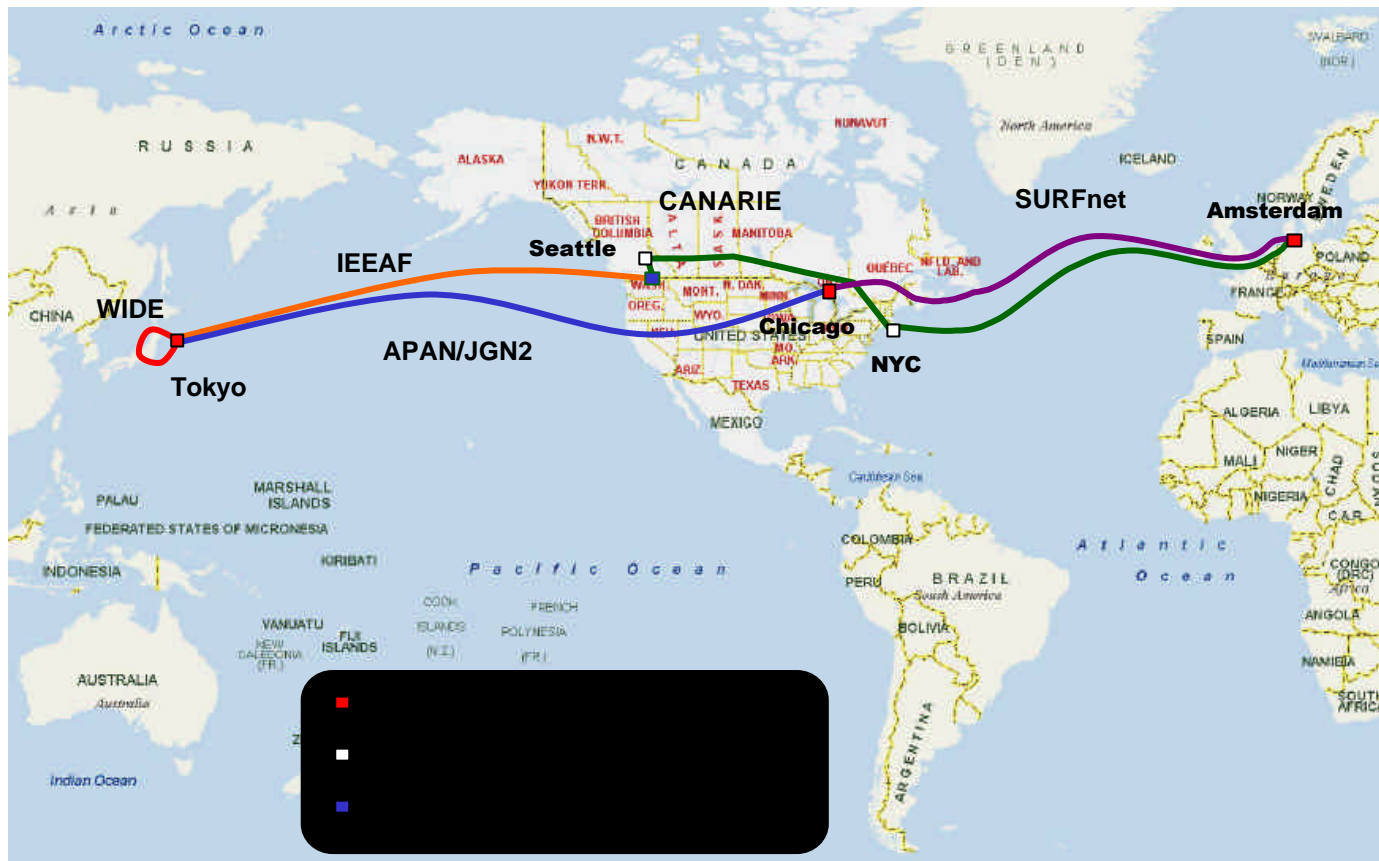  - Especially on trans-ocean connection

# SC|05 network path



**Kei Hiraki**

**University of Tokyo**

# SC|05 address arrangement

2001:200:0:1cf0::/64
VLAN 654

2001:200:0:1cf1::/64
VLAN 3300

opt6 — TGNLE-1 — Fujitsu XG800 — CISCO 7604 — Foundry NI40G — Foundry NI40G — Foundry NI40G — Hitachi GS4000 — Hitachi GS4000 — Hitachi GS4000

2001:200:0:1cf0::12

2001:200:0:1cf0::1

2001:200:0:1cf1::1

Tokyo

opt9 — Fujitsu XG800 — Hitachi GS4000 — Force10 600 — Force10 E1200 — Hitachi GS4000

2001:200:0:1cf3::10

2001:200:0:1cf3::1

2001:200:0:1cf2::2

2001:200:0:1cf2::1

2001:200:0:1cf1::3

2001:200:0:1cf3::/64
VLAN 655

2001:200:0:1cf2::/64
VLAN 359

Seattle SC05 booth 233     Seattle Westin     Amsterdam     Chicago

# SC|05 network path map

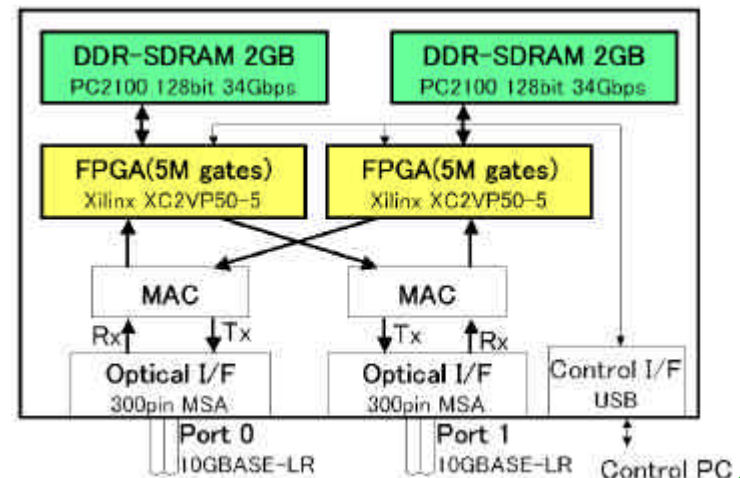# Difficulty (3)

- Observation
  - 500ms precise delay shows much better performace than real network
    - Anue 10Gbps delay emulator       same level to local performance
    - Seattle -> Tokyo -> Chicago -> Amsterdam -> Seattle network
        Unstable and less performance to local performance

  - Unnecessary packet loss at receiving server
    - Bursty packet stream at receiver
    - Bottleneck at I/O bus, CPU utilization, and amount of buffer

  - Performance difference by intermediate switch / router
    - Difference is about 10%
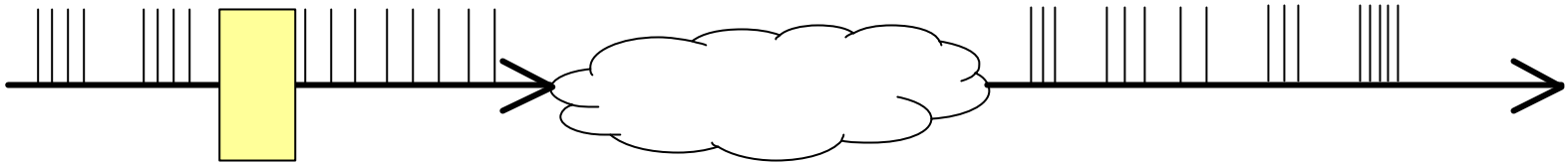    - Sender side pacing is not effective

# Receiver side pacing

- Objective
  - Reduce packet losses due to receiver bottleneck
    - Act as large receiving buffer at receiving Network Interface Card
    - Set to the maximum receiving speed of the system

- Implementation
  - SANSEI-system TGNLE-1 FPGA based buffering system
  - FPGA is used to implement fine-grain pacing mechanism
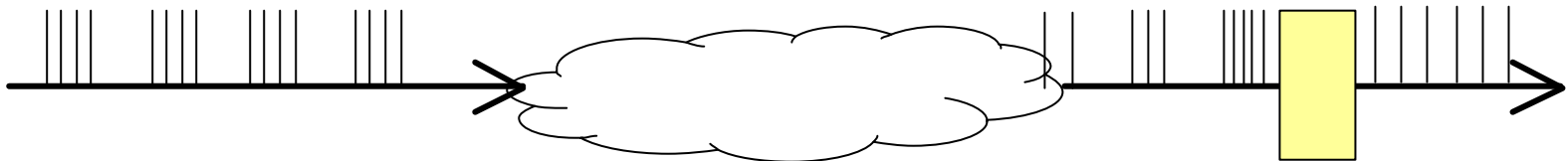  - 2GB buffer memory (DDR memory)

# Pacing

- Sender side pacing
  - Effective for the bottleneck middle of the network
    - Complex hardware to cooperate with TCP protocol
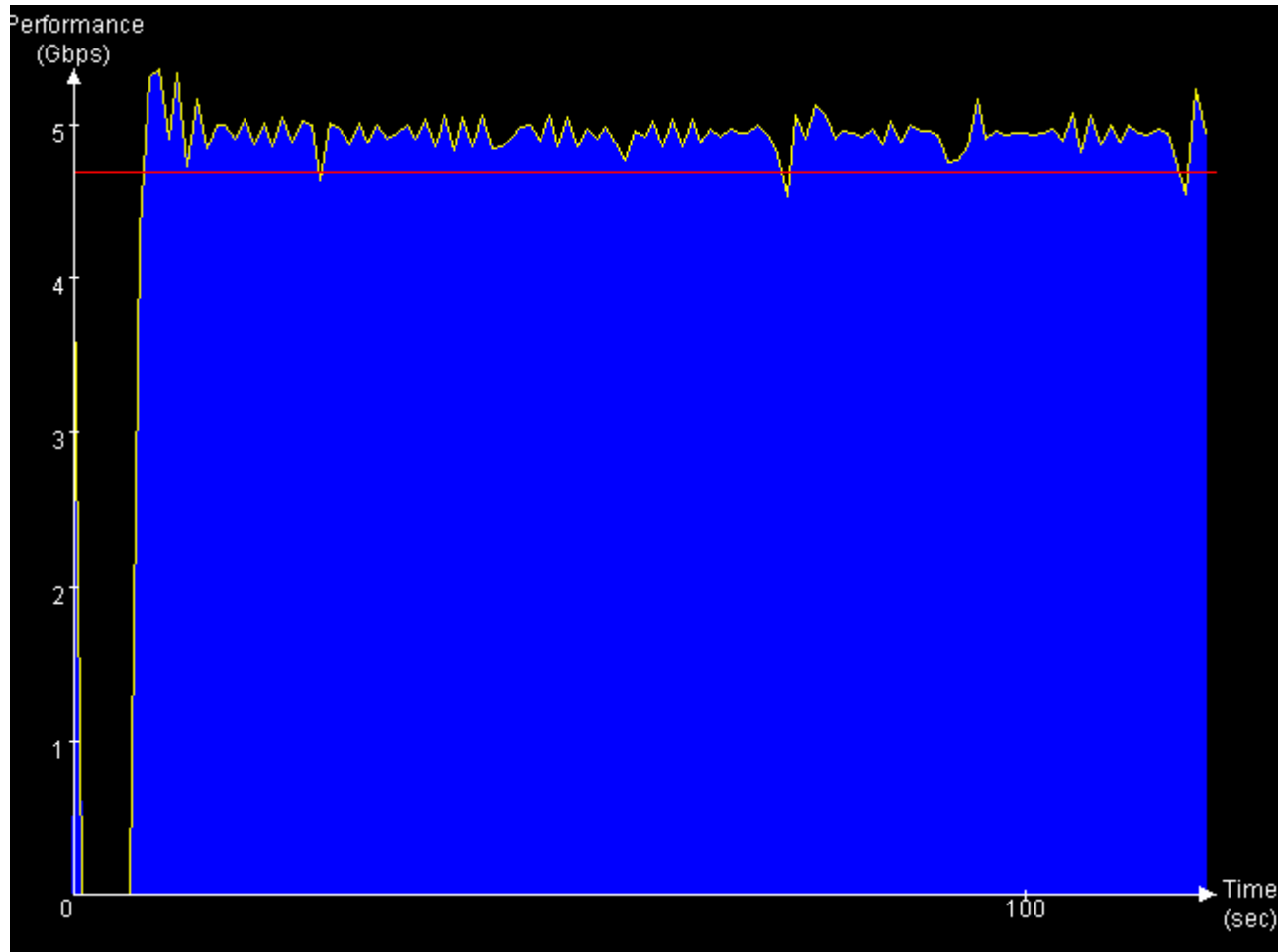    - Not effective to some 10G switches

- Receiver side pacing
  - Effective to NIC whose maximum bandwidth is less than 10G
  - Simple hardware. It can be implemented in NIC
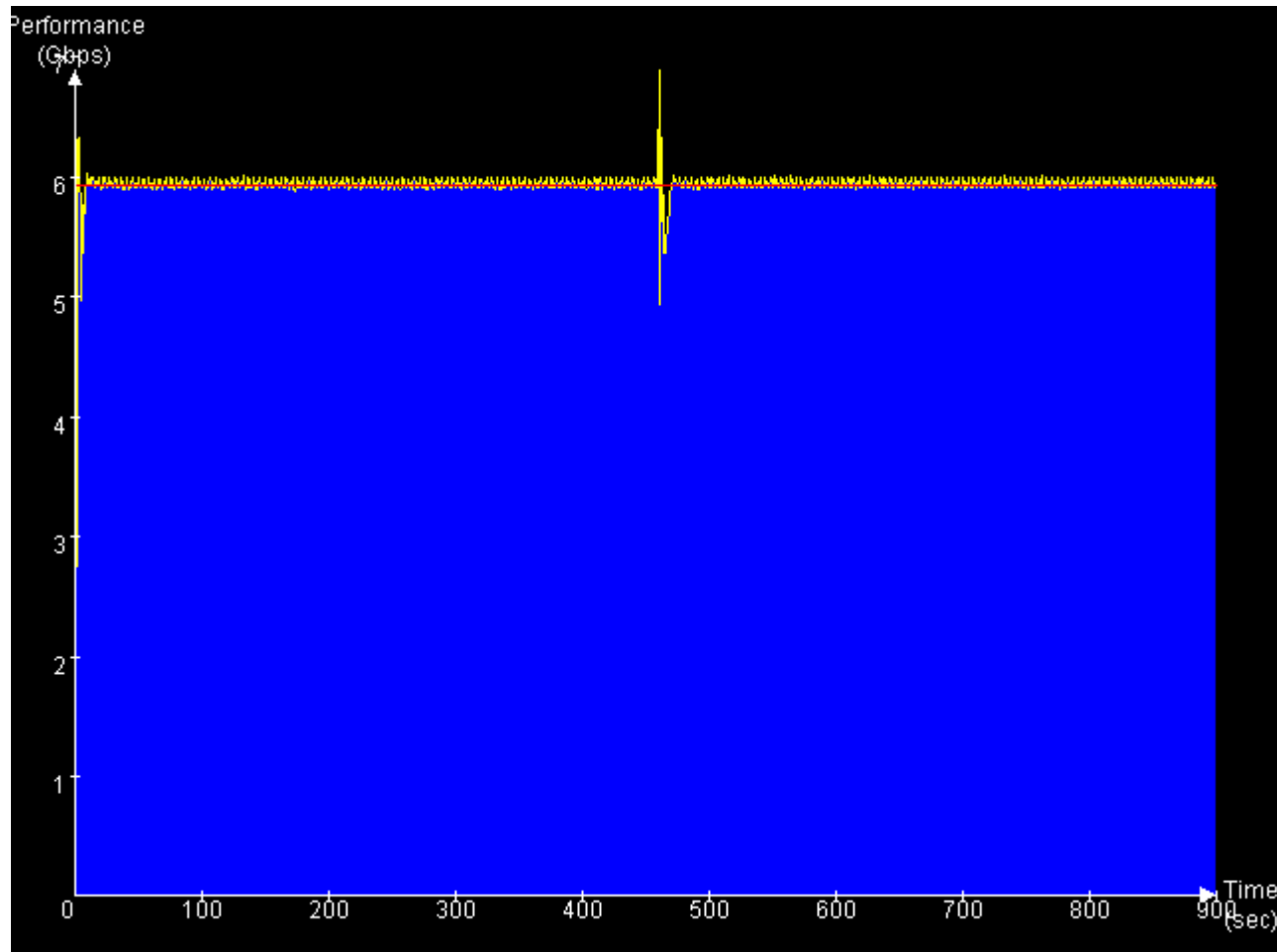  - Increase latency when bandwidth change

# Preliminary Results
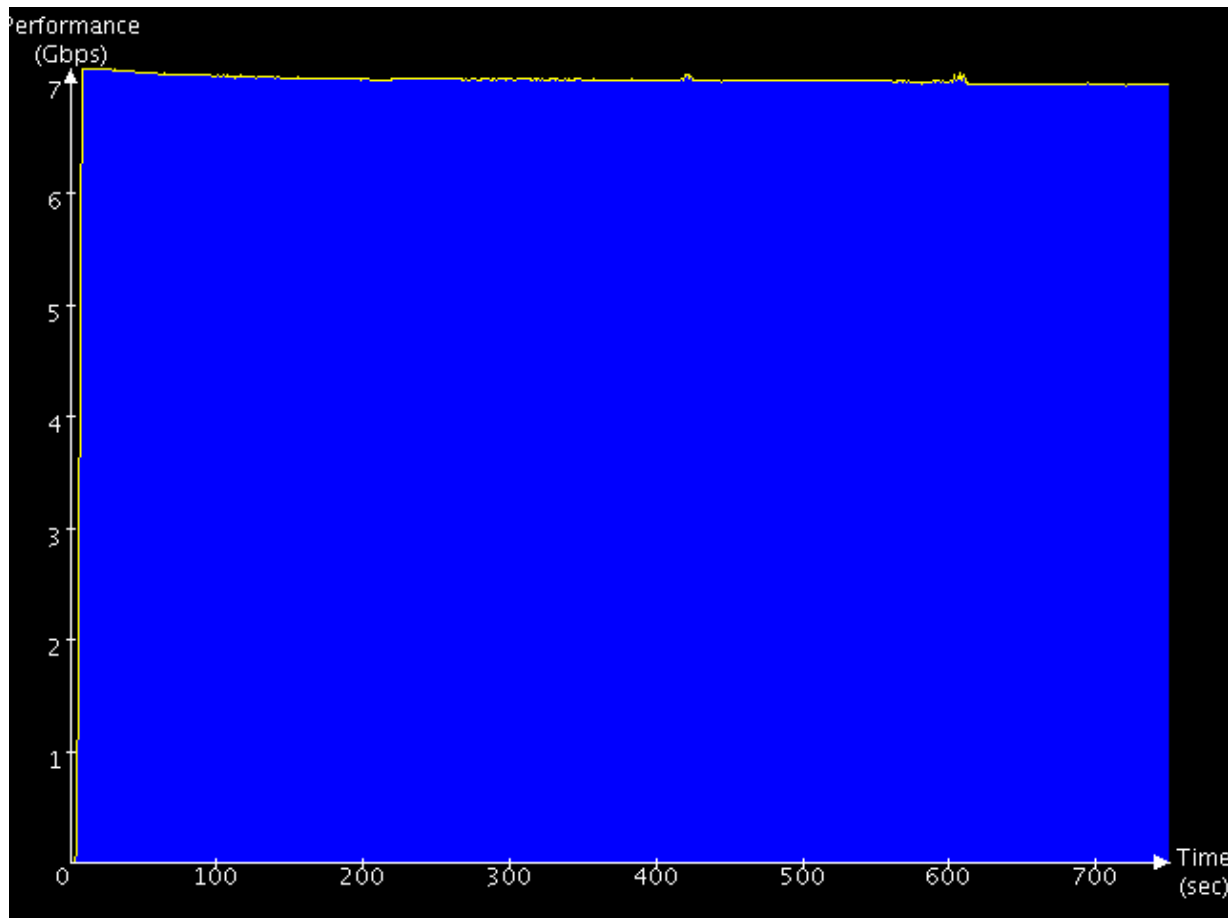
- Without receiver side pacing

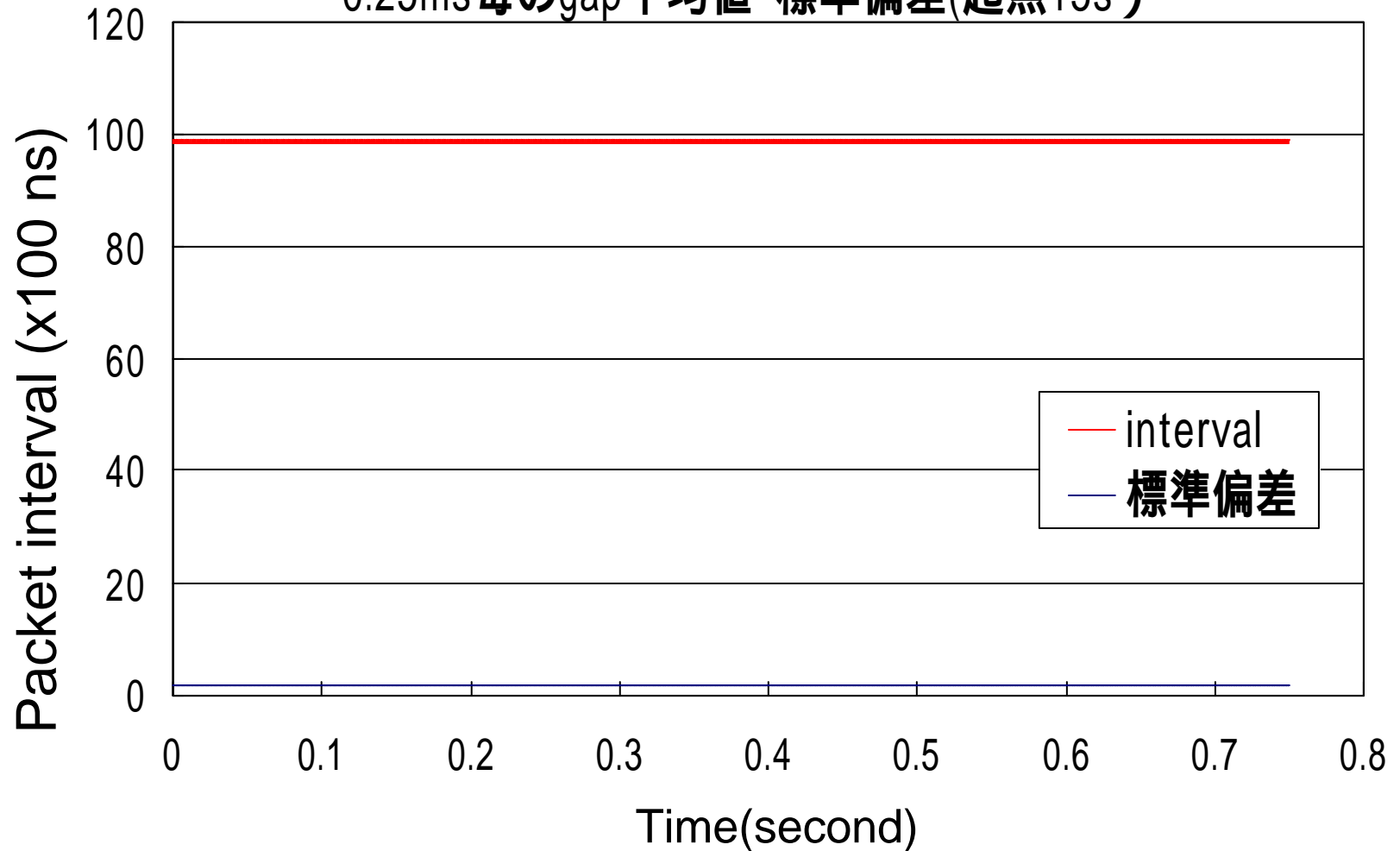# Preliminary Results

- Without receiver side pacing

# Preliminary Results

- IPv6 TCP single stream
  - 6.96Gbps  (use of "receiver side pacing")
  - 5.58Gbps  (without "receiver side pacing")

# Local  Data packets@Sender-TX



0.25ms        gap                    (        15s

Packet interval (x100 ns)

Time(second)

interval

# Local ACKpackets@Sender-RX



**Local ACKpackets@Sender-RX**

0.25ms gap ( 15s

Packet interval (x100 ns) vs Time(second)

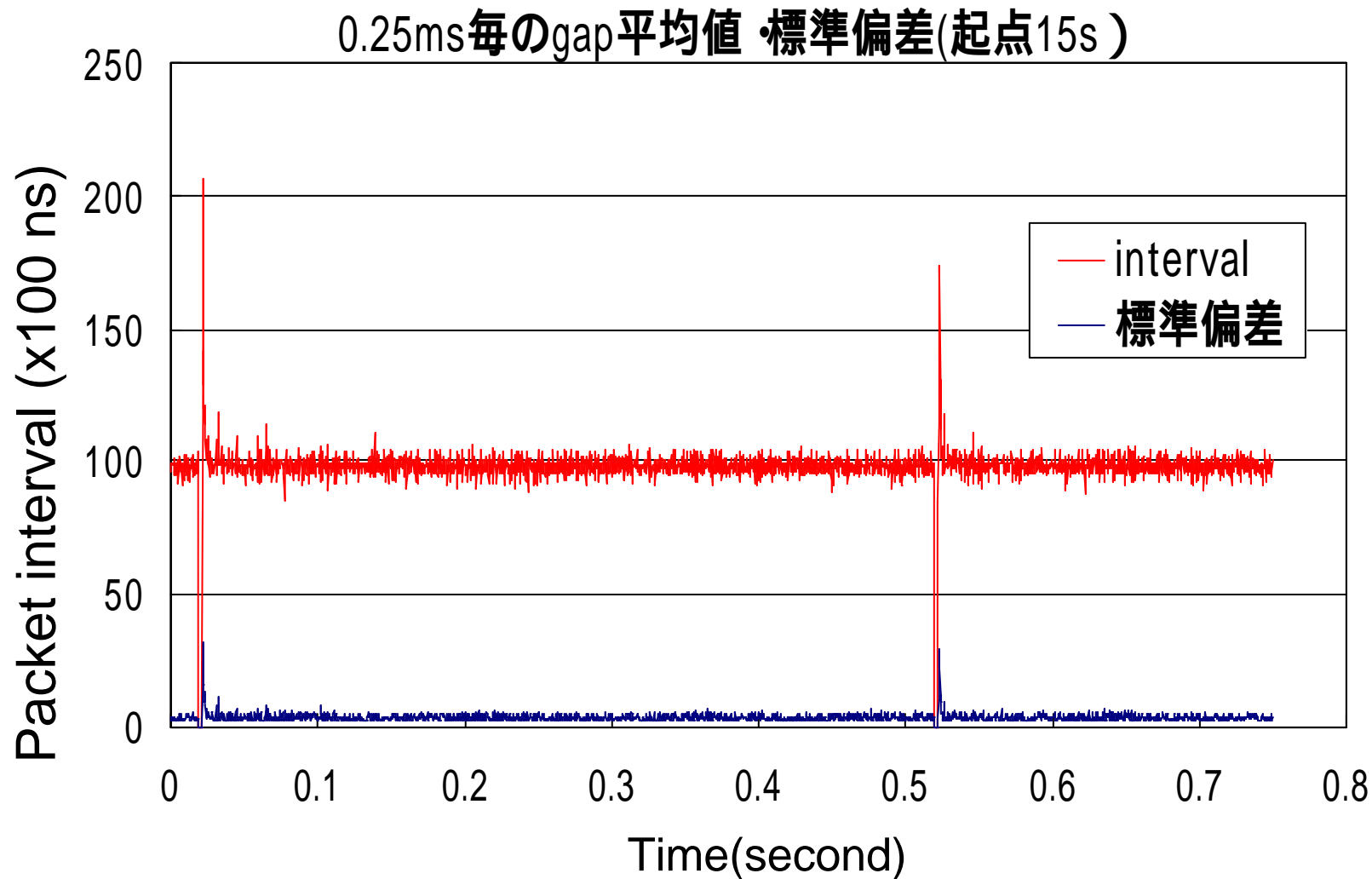- interval

# One way, 500ms network  Data Packets@Sender-TX

# Data packets @Receiver-RX

# One way, 500ms ACKPackets@Receiver-TX

# One way, 500ms RTT  ACK packets@Sender-RX

# Chicago loopback, 280ms  Data packets  @Sender-TX

# Chicago loopback ACK packets@Sender-RX

# MRTG along the network path

Measuing point


Seattle NI40G


T-LEX output


NEZU NI40G
to JGN2


JGN2 Chicago GS4000


Amsterdam UvA Force10

Kei Hiraki

# Performance dip by MAC learning



Frame Rate * Frame Size
Flooder -> T-LEX ——> SEA (RTT=177ms) (64B frame)

Legend: WAN-PHY limit for 64B frame — Rx — Tx

# Short(14+46+4=64Byte)フレームでの速度

Tx server -> NI40G -> F10 ——> Chicago F10



凡例: 送信側TxA ── 受信側TxB(RxBと思って良い)

# System for IPv6

- Server configuration
  - CPU:  Dual AMD Opteron 248, 2.2GHz
  - Mother Board: RioWorks HDAMA rev. E
  - Memory: 1G bytes, Corsair Twinx CMX512RE-3200LLPT x 2,  PC3200, CL=2
  - OS: Linux kernel 2.6.12
  - Disk Seagate IDE 80GB, 7200r.p.m. (disk speed not essential for performance)
- Network interface card
  -  Chelsio T110 (10GBASE-SR), TCP offload    OFF
  -  PCI-X/133 MHz I/O bus connection

Technical points
- Traditional tuning and optimization methods
  - Tuning and optimization by Inter-layer coordination  Univ. of   Tokyo)
  - Fine-grained pacing by offload engine
  - optimization at slow start phase
  - Utilization of flow control
  - Receiving side pacing

- Performance of single stream TCP on 30000 km circuit

  - 6.96 Gbps (TCP payload), standard Ethernet frame
    - 208,800 terabit meter / second

# System for IPv4

- Server configuration
  - CPU:  Quad  Intel Xeon,3.6GHz
  - Server  IBM x366
  - Memory: 32GB DDR2
  - OS: Windows Server 2003
  - Network interface card
  -  Netrion Xframe II
  -  PCI-X V2.0 266MHz I/O bus connection
  -  No TOE function

Technical points
- Jumbo frame (9014B)

-  Tuning and optimization by Inter-layer coordination  Univ. of Tokyo)
  - optimization at slow start phase
  - Utilization of flow control

- Performance of single stream TCP on 30000 km circuit

  - 7.96 Gbps (TCP payload), standard Ethernet frame
                 239,820 terabit meter / second
  10%  more than previous Land Speed Record

# Internet2 Land Speed Record

The Internet2 Land Speed Record (I2-LSR) competition for the highest-bandwidth, end-to-end networks is an open and ongoing contest.

# Current Records

## IPv6 Category

**Single Stream Class**: 167,400 terabit-meters per second by a team consisting of members from the University of Tokyo, the WIDE Project, Fujitsu Computer Technologies LTD, and others accomplished by transferring 585 gigabytes of data across 30,000 kilometers of network in 15 minutes at **an average rate of 5.58 gigabits per second**.

**Multiple Stream Class**: 167,400 terabit-meters per second by a team consisting of members from the University of Tokyo, the WIDE Project, Fujitsu Computer Technologies LTD, and others accomplished by transferring 585 gigabytes of data across 30,000 kilometers of network in 15 minutes at **an average rate of 5.58 gigabits per second**.

## IPv4 Category

**Single Stream Class**: 216,300 terabit-meters per second by a team consisting of members from the University of Tokyo, the WIDE project, and Chelsio Communication and other organizations by sending 1485 gigabytes of data across 30,000 kilometers of network over 30 minutes at **an average rate of 7.21 gigabits per second**.

**Multiple Stream Class**: 216,300 terabit-meters per second by a team consisting of members from the University of Tokyo, the WIDE project, and Chelsio Communication and other organizations by sending 1485 gigabytes of

# Lessons learnt from LSR experiments

- Difficult combination -- WAN PHY, IPv6, Jumbo frame, L3 switching
  - There is no trouble-free switch
  - Unexpected packet losses by many reason
  - Further investigation on "Flow Control" is essential

- Unnecessary packet losses by packet clustering
  - Max receiving speed is normally less than wire speed
  - Behavior is different from switch to switch
  - Receiving side pacing is quite useful.

- Unnecessary packet losses by bursty TCP traffic
  - Sending side pacing is effective if the number of intermediate switch is small
  - Intermediate switches and routers erase effect of pacing
  - Difference between WAN PHY and LAN PHY may make trouble

# Lessons learnt from LSR experiments

- Use of wide-area L2 network
  - Spanning Tree algorithm may make unstableness
  - MAC address learning may cause packet losses
  - Difficulty in debugging
  - Switches for trans-ocean (trans-pacific, etc.) should have very large buffers and pacing capability

- 10G network interface
  - Pacing capability is essential
  - Large input buffer (2*RTT*BW) or receiving side pacing is useful
  - Proper setting of window size, buffer size and queue length is essential

# Toward end-to-end 10Gbps internet

- Purchase of 10Gbps NIC is always disappinting
  - Performance of 10 Gbps NIC may be worse than GbE
  - Packet losses!!

- Users want disk to disk, client to server performance
  - Disk performance
  - CPU bottleneck

- Importance of wire-rate switch
- Switches should be stable under any user packet sequences
- There may be a good switch and a bad switch. But the reason is unknown.

- Our next target is about 9Gbps.

# Observation on Protocol

- We used traditional? TCP (loss-based)
  - If traffic is controlled, standard TCP works well
  - Selection of coefficients in AIMD algorithm may improve behavior

- Delay based TCP may not work properly
  - Major cause of variation of RTT is not congestion but meaningless jitter
  - Artificial packet clustering by intermediate switches/routers
  - Insufficient amount of intermediate buffers

- First several second may have problem
  - Experience from "Receiving side pacing"

# Observation on Protocol

- Pacing is essential
  - Current switches/routers are not compatible to LFN
  - We need good cooperation algorithm between network layer and MAC layer

- Control of bursty behavior is also essential
  - Currently, users cannot control behavior of intermediate switches/routers

- Large buffer size at receiving NIC is very important
  - About RTT * BW size?

# Conclusion

We thanks all the people who support our experiments

Next target is 9Gbps through WAN PHY network

Current 10Gbps devices and software technology is still far from satisfaction

Buffer size, burstness control, large-scale L2 network

The next step is disk to disk transfer on single TCP stream

Optimization of cache memory and memory buses

After that, 10Gbps http service

We are developing Hardware tools
– Wire-rate packet capture
– Traffic multiplier (for pushing networks)
– Packet filtering

# Photo of SC|05 University of Tokyo booth

# Devices at SC|05

# Devices at Pacific Northwest Gigapop

# History of single-stream IPv4 Internet Land Speed Record

Distance bandwidth product
Pbit m / s

**10 Gbps * 30,000km**

2005/11/10
Data Reservoir project
WIDE project
240 Pbit m / s

2004/11/9
Data Reservoir project
WIDE project
149 Pbit m / s

Year

# History of single-stream IPv6 Internet Land Speed Record

Distance bandwidth product
Pbit m / s

**10 Gbps * 30,000km**

2005/11/13
Data Reservoir project
WIDE project
208 Pbit m / s

2004/10/29
Data Reservoir project
WIDE project
167 Pbit m / s

1,000

100

10

1

2000  2001  2002  2003  2004  2005  2006  2007

**Year**

# Thanks

# Land Speed Record

| | | | |
|---|---|---|---|
| 2000   3   21 | 4278 | Microsoft<br>Qwest Communications<br>University of Washington<br>USC Information Sciences Institute | IPv4   Single Stream |
| 2000   3   29 | 5384 | Microsoft<br>Qwest Communications<br>University of Washington<br>USC Information Sciences Institute | IPv4   Multiple Stream |
| 2002   4   9 | 4933 | University of Alaska at Fairbanks<br>Faculty of Science of the University of Amsterdam<br>SURFnet | IPv4   Single Stream |
| 2002   8   22 | 40 | Oregon Gigapop<br>NYSERNet<br>University of Oregon | IPv6   Multiple Stream |
| 2002   9   27 | 2517 | ARNES<br>DANTE<br>RedIRIS | IPv6   Single Stream |
| 2002   10   9 | 5154 | ARNES<br>DANTE<br>RedIRIS | IPv6   Single Stream |
| 2002   11   19 | 10136 | Nationaal Instituut voor Kernfysica en Hoge-Energiefysica (NIKHEF)<br>Universiteit van Amsterdam (UvA)<br>Stanford Linear Accelerator Center (SLAC)<br>California Institute of Technology (Caltech) | IPv4   Single and Multiple stream |
| 2003   2   23 | 23888 | California Institute of Technology (Caltech)<br>CERN<br>Los Alamos National Laboratory (LANL)<br>Stanford Linear Accelerator Center (SLAC) | IPv4   Single and Multiple stream |

| | | | |
|---|---|---|---|
| 2003   5   6 | 6947 | California Institute of Technology (Caltech) CERN | IPv6   Single and Multiple stream |
| 2003   10   10 | 38420 | California Institute of Technology (Caltech) CERN | IPv4   Single and Multiple stream |
| 2003   11   11 | 61752 | California Institute of Technology (Caltech) CERN | IPv4   Single and Multiple stream |
| 2004   2   22 | 68431 | California Institute of Technology (Caltech) CERN | IPv4   Multiple stream |
| 2004   4   14 | 69073 | SUNET Sprint | IPv4   Single stream |
| 2004   5   6 | 77699 | California Institute of Technology (Caltech) CERN | IPv4   Multiple stream |
| 2004   6   25 | 104529 | California Institute of Technology (Caltech) CERN | IPv4   Multiple stream |
| 2004   6   28 | 103583 | California Institute of Technology (Caltech) CERN | IPv4   Single stream |

**Kei Hiraki**　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　**University of Tokyo**

| | | | |
|---|---|---|---|
| 2004　9　12 | 124935 | SUNET<br>Sprint | IPv4　Single and Multiple stream |
| 2004　11　8 | 184877 | California Institute of Technology (Caltech)<br>CERN<br>CENEC | IPv4　Multiple stream |
| 2004　11 | 148850 | University of Tokyo<br>Fujitsu Computer Technologies<br>WIDE<br>Chelsio Communications | IPv4　Single stream |
| 2004　12　25 | 216300 | University of Tokyo<br>Fujitsu Computer Technologies<br>WIDE<br>Chelsio Communications<br>APAN<br>JGN2<br>CANARIE<br>SURFnet<br>Universiteit van Amsterdam | IPv4　Single and Multiple stream |

| | | | |
|---|---|---|---|
| | | | |
| 2005　1　19 | 72225 | CERN<br>CALTECH | IPv6　Single and<br>Multiple stream |
| 2005　10　29 | 167400 | University of Tokyo<br>Fujitsu Computer Technologies<br>WIDE<br>Chelsio Communications<br>JGN2 | IPv6　Single and<br>Multiple stream |
| 2005　11　11 | 239820 | University of Tokyo<br>Fujitsu Computer Technologies<br>WIDE<br>Microsoft<br>JGN2<br>CANARIE<br>SURFnet<br>Universiteit van Amsterdam | IPv4　Single and<br>Multiple stream |
| 2005　11　13 | **208800** | University of Tokyo<br>Fujitsu Computer Technologies<br>WIDE<br>Chelsio Communications<br>JGN2<br>CANARIE<br>SURFnet<br>Universiteit van Amsterdam | IPv6　Single and<br>Multiple stream |

# Future of Internet2 Land Speed Record

- 10Gbps era.
  - I/OBus bottleneck
    - PCI-express or  PCI-X  266MH
    - WAN PHY limit is about 9.1Gbps (TCP payload bandwidth)

- 40Gbps era.
  - 40Gbps WAN will be established in several years
    - At first, domestic network in Japan or US
    - 7500 km is necessary to make new record
  - 40GbpsLAN and NIC is still unpredictable
    - 40Gbps  Ethernet  or  100Gbps  Ethernet
    - We need another new I/O bus

- 100Gbps era.
  - 100GbpsEthernet standard.  When?
  - Yet another I/O bus and CPU will be definitely necessary
  - Year 2010? 2015?