



**N I C
F N D
S U S
T R E
R Y N
C E**

NFS over RDMA

Sweet Spot

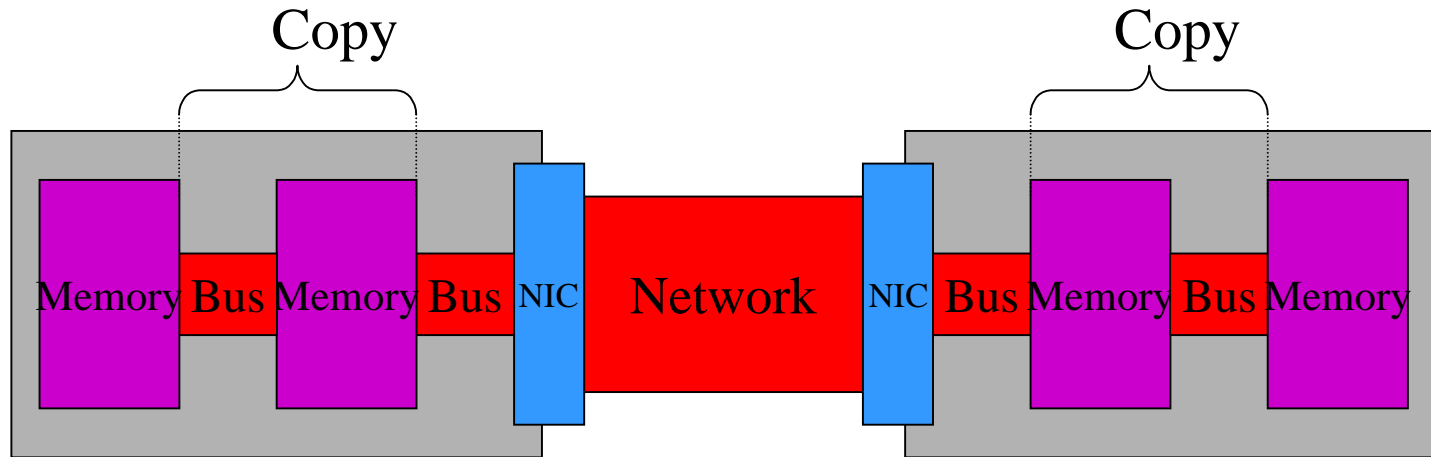
Brent Callaghan
Sun Microsystems, Inc.

September 22-24



**N I C
F N O
S D N
U S F
T R E
R Y N
C E**

Network vs Bus Performance

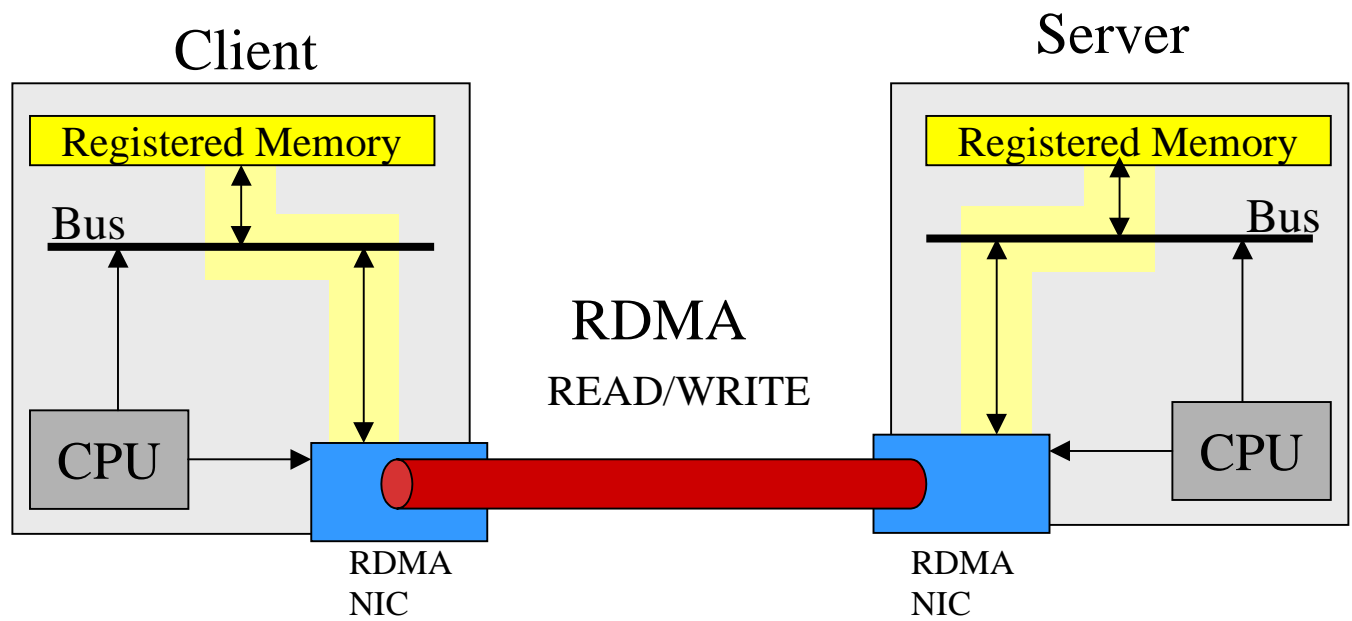


Latency gets worse with each memory copy
which further loads the CPU.



What is RDMA ?

- DMA: Direct Memory Access
- RDMA: *Remote* Direct Memory Access
- Supports Direct Placement
- Networking offload for CPU

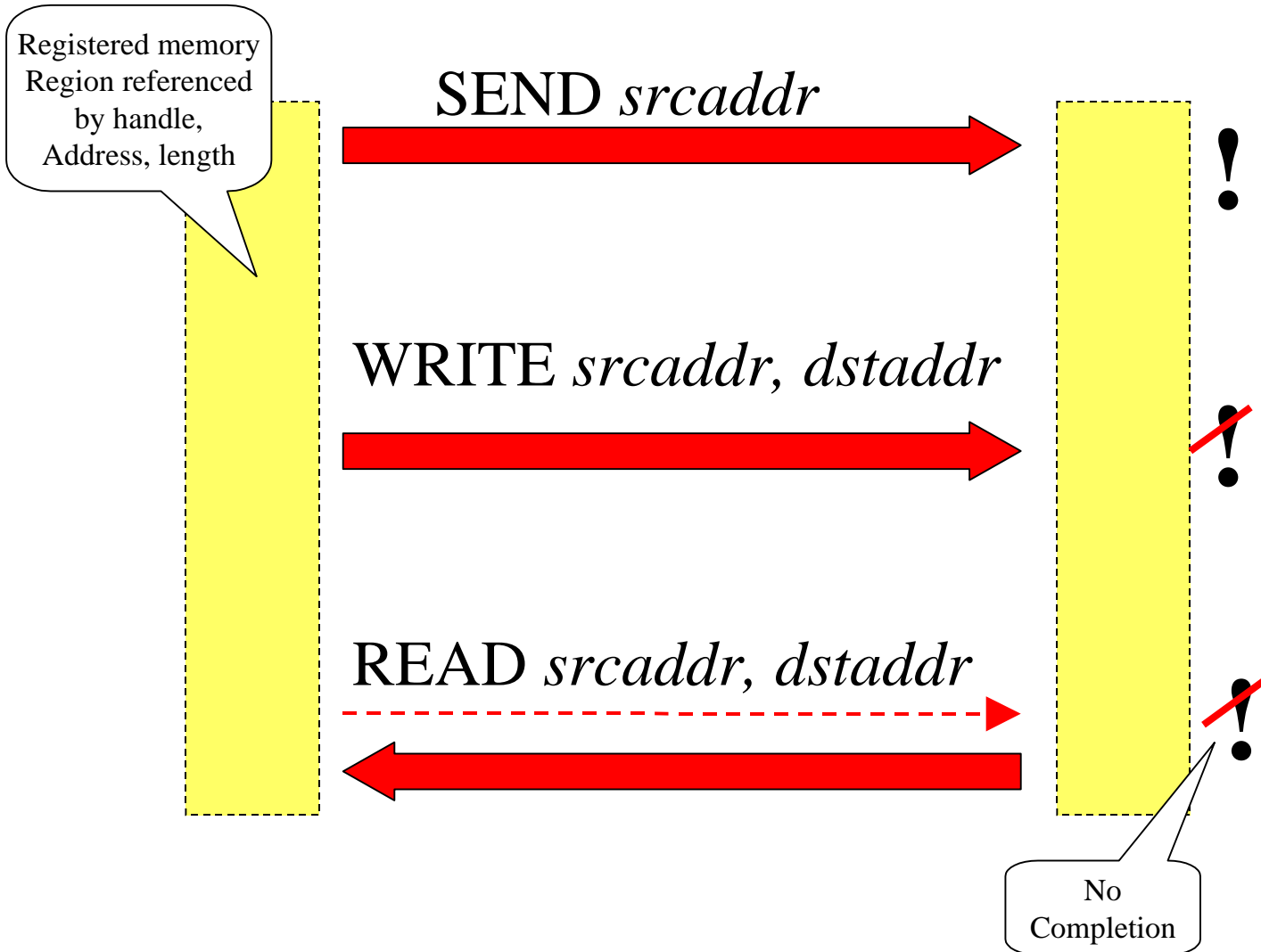




**N I C
F N D
S U S
T R E
R E N
C E**

September 22-24

RDMA Operations



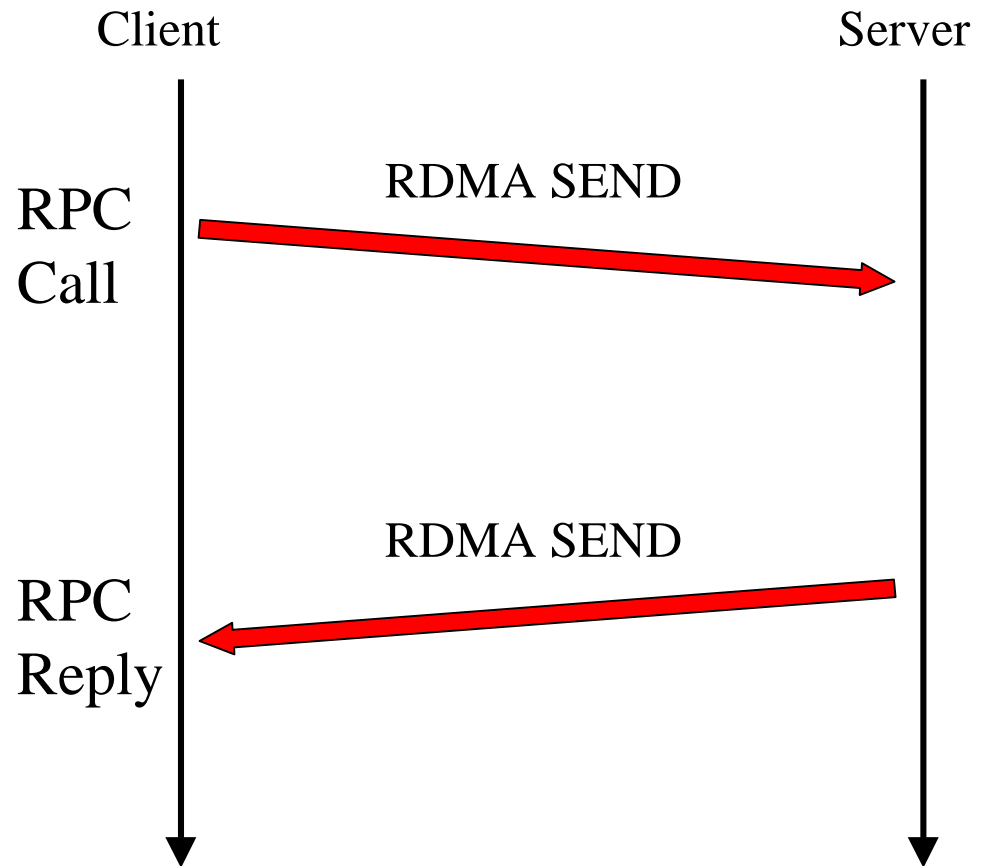


NFS INDUSTRY CONFERENCE

September 22-24

Small RPC Messages

Most RPC Messages are small.
Examples:
LOOKUP
GETATTR

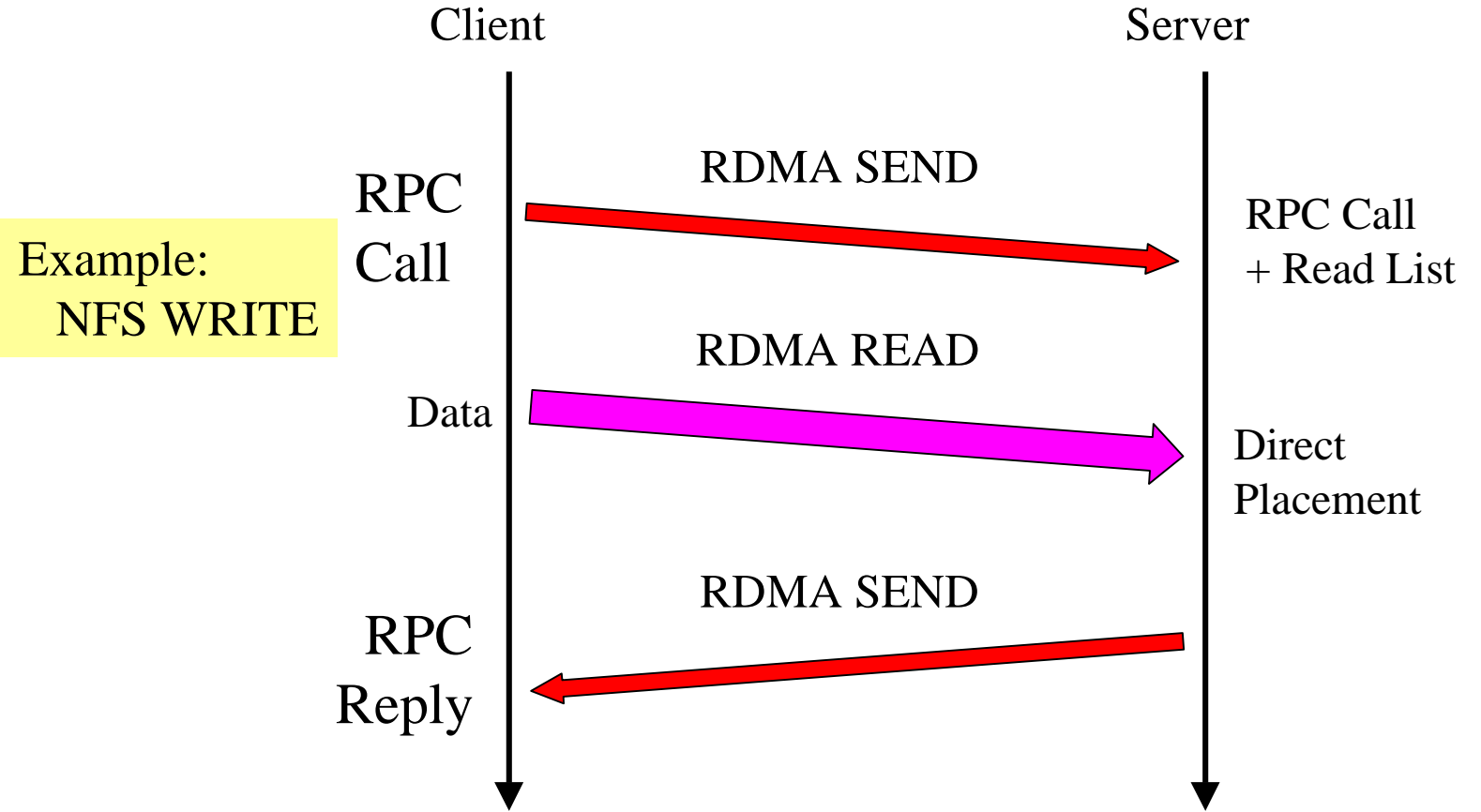




**N I C
F N D
S U S
T R E
R Y N
C E**

September 22-24

Big RPC Call





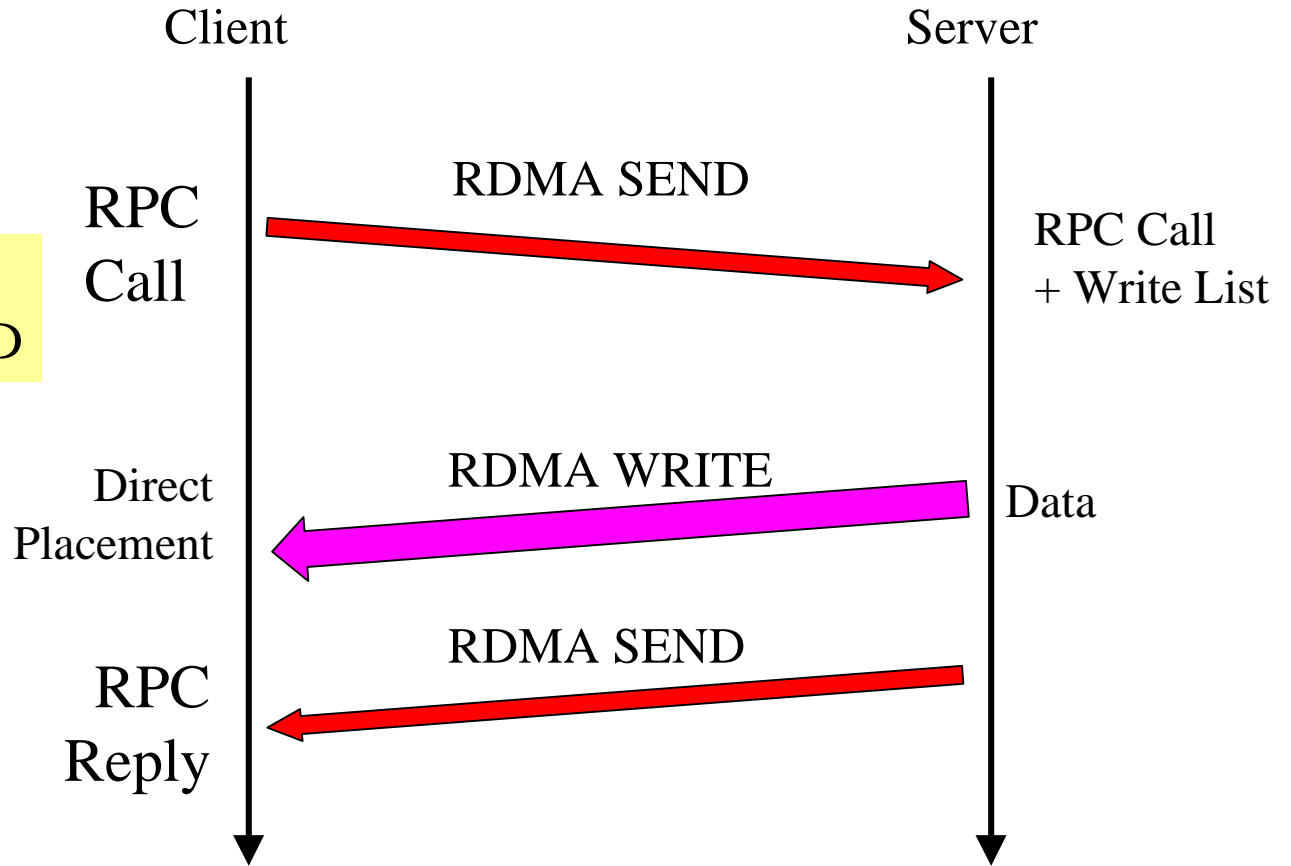
NFS INDUSTRY CONFERENCE

September 22-24

Big RPC Reply



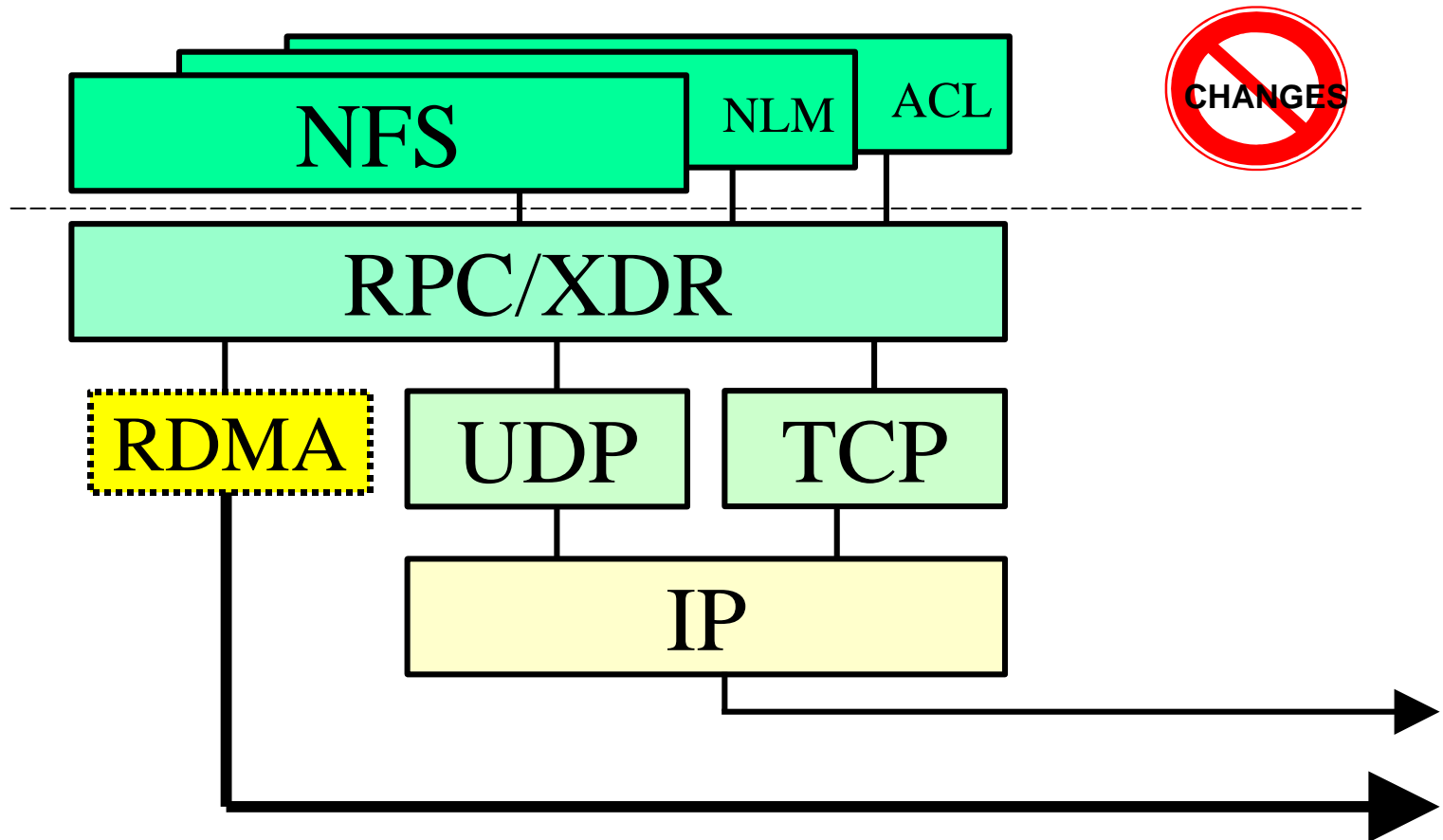
Example:
NFS READ





**N I C
F N D
S U S
I N D
T R E
R Y N
C E**

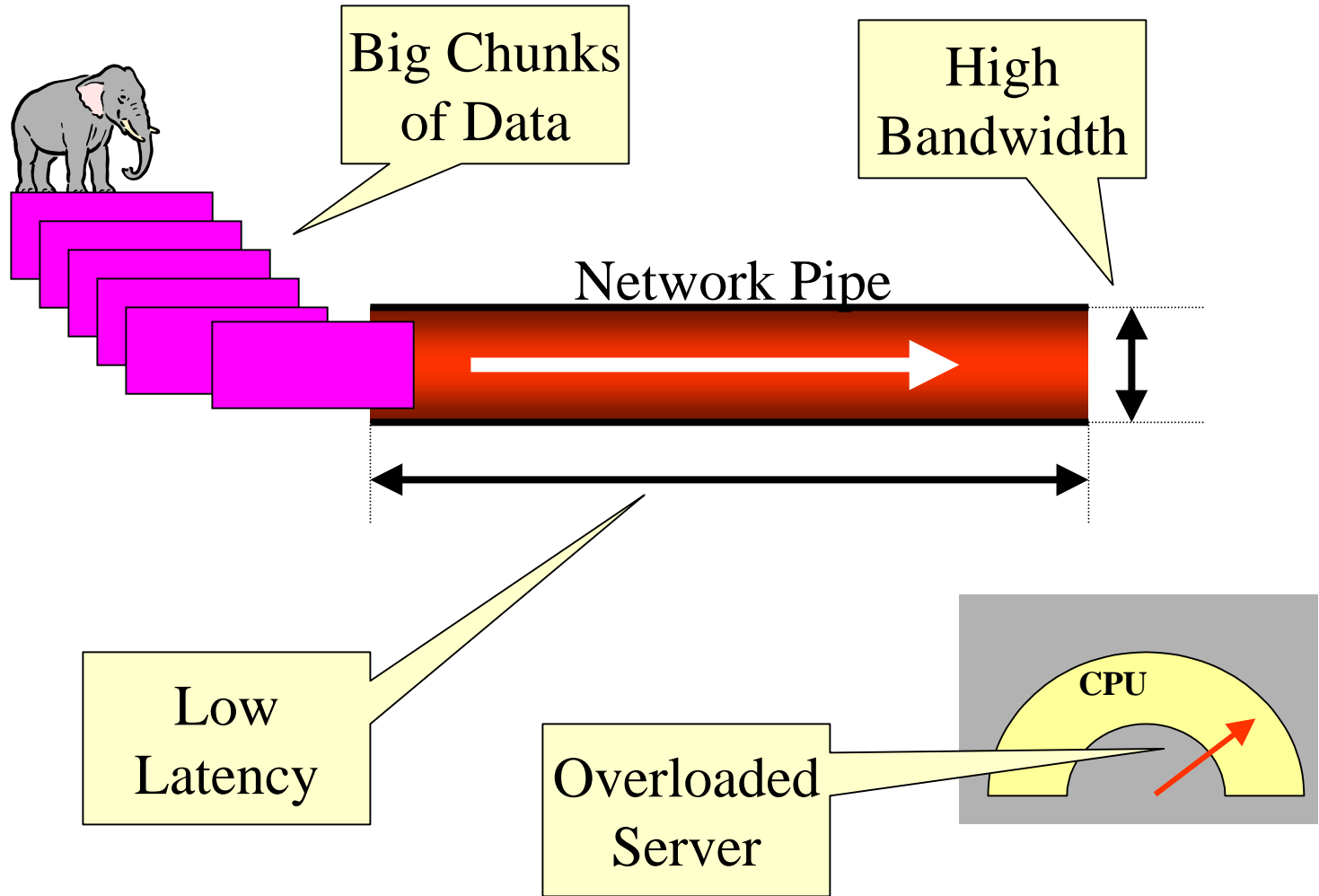
Adapting NFS to RDMA





**N I C
F N D O
S U S F
T R T R
R E
E
N C
E**

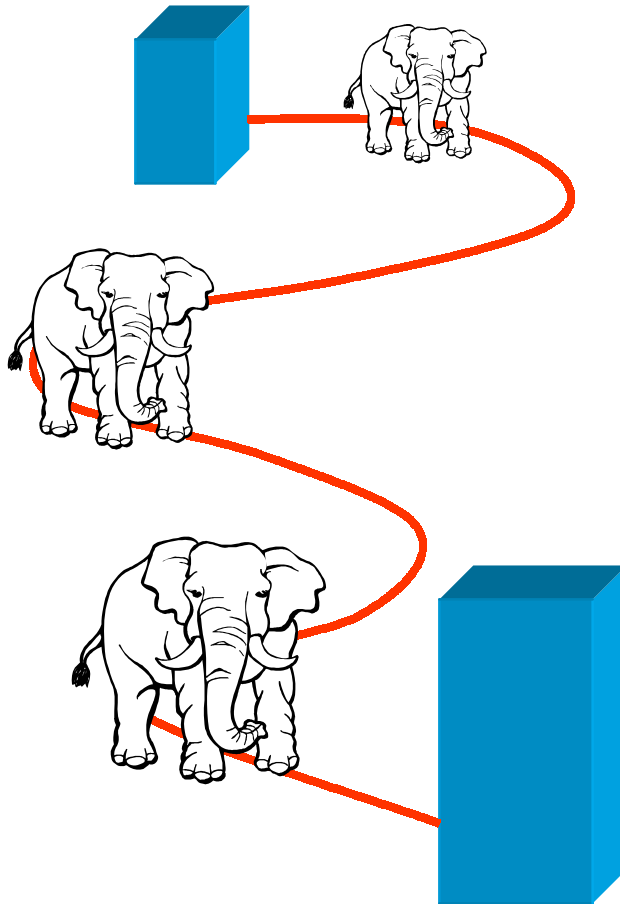
RDMA *Sweet Spot*



Big Chunks of Data



**N I C
F N D O
S U S F
T R E
R Y E
N C
E**

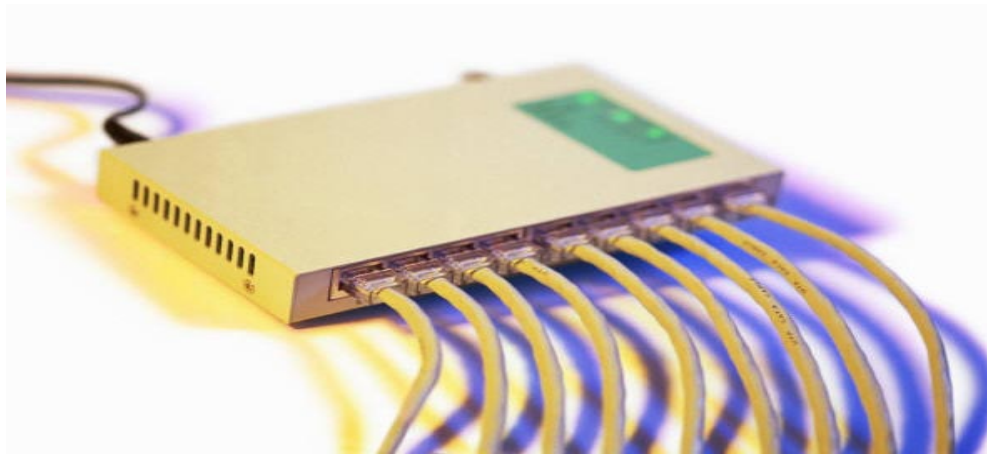


- RDMA Excels at Moving big data
 - NFS
 - Database
 - Backup
 - Replication
- Small data ?
 - May still be some improvement from protocol offload.

High Bandwidth



**N I C
F N O
S D N
U S F
T R E
Y R E
C E**



- 10 Mb/sec
- 100 Mb/sec
- 1,000 Mb/sec
- 10,000 Mb/sec

- 10 & 100Mb Ethernet - easily handled by CPU & bus
- 1 Gb Ethernet drains CPU & bus cycles
- At 10Gb, CPU & bus overwhelmed.
Must have RDMA.

Low Latency



**N I C
F N O
S D N
I U F
N S E
D T R
I R E
N C
E**



- RDMA Protocols more “chatty”
 - More network messages
 - Per message latency is cumulative
- RDMA performance advantages may fall off with distance

Overloaded Servers



N I C
F N O
S D N
I U F
N S E
D T R
E
R
E
N
C
E

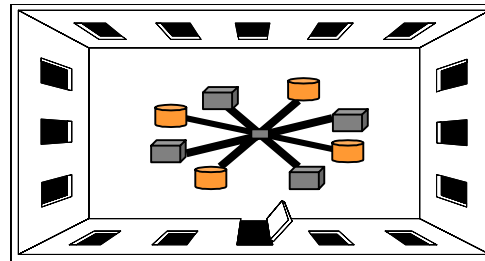


- RDMA Hardware offloads protocol processing from busy CPUs.
- **RDMA Direct Data Placement** relieves congested memory and I/O buses.

RDMA in the Data Center



**N I C
F N D O
S U S F
T R E
R Y R E
N C
E**



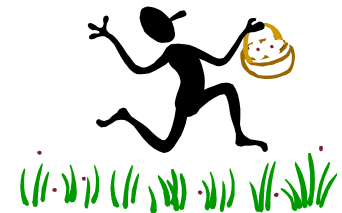
- Low latency - nodes are meters apart
- High bandwidth - runs are short, cheap
- Lots of data movement
- Infiniband clusters

NFS Storage Upgrade



**N I C
F N O
S D N
U S F
T R E
R Y N
C E**

- Many customers already using NFS for data center storage.
- RDMA over Ethernet
 - Just a NIC upgrade.
 - Use existing fabric, switches, routers
 - No change to administration

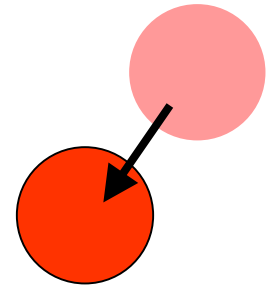




**N I C
F N O
S D N
I U F
N S E
D T R
I R E
N
C
E**

Sweet Spot will shift

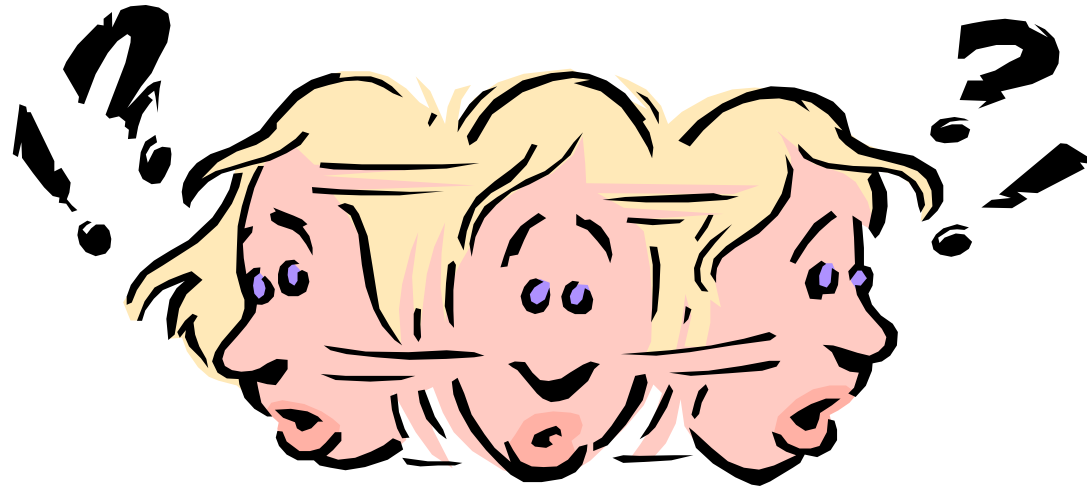
- As NFS/RDMA becomes available on Ethernet
- As the price drops
 - RDMA on the motherboard
- As available bandwidth grows
- As big data increases
- As CPU performance improves





**N I C
F N D
S U S
T R E
R E N
C E**

September 22-24



Questions & Answers

2003 NFS Industry Conference

Page 17 of 17