



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

NFS at 10Gb/sec

Brent Callaghan
Sun Microsystems, Inc.

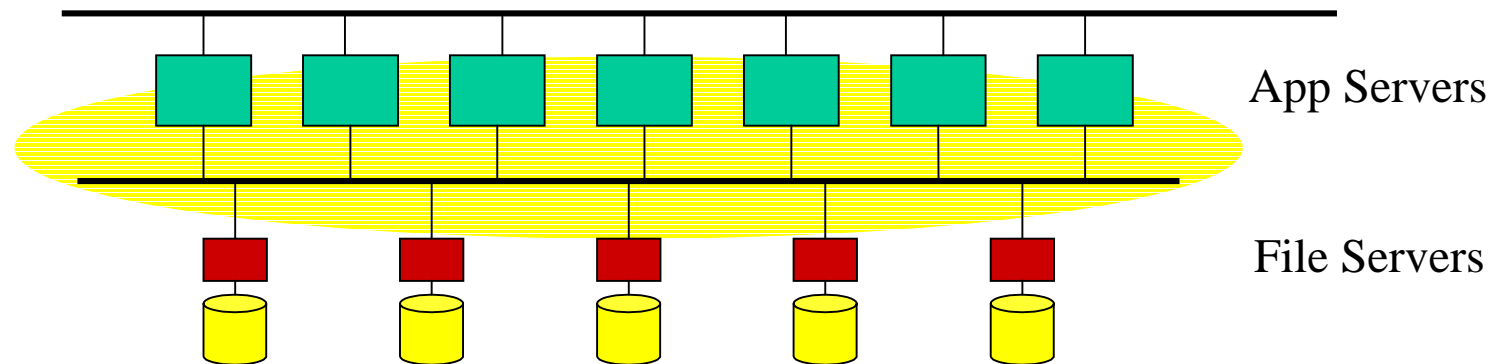
October 22-23, 2002



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

Data Center NFS

- Transaction processing clusters
 - Database, Web, Email, eCommerce
- Dataless application servers
- Apps & Data are meters apart

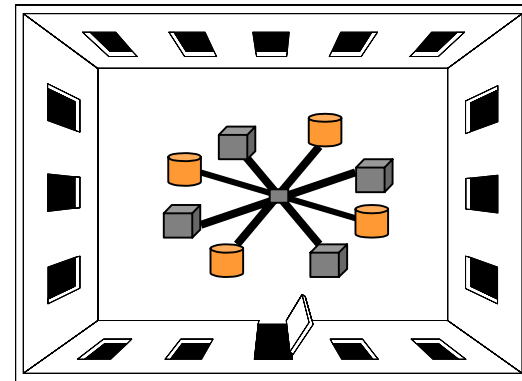




**N I C
F N O
S D N
U S F
T R E
R Y N
C E**

Data Center Network

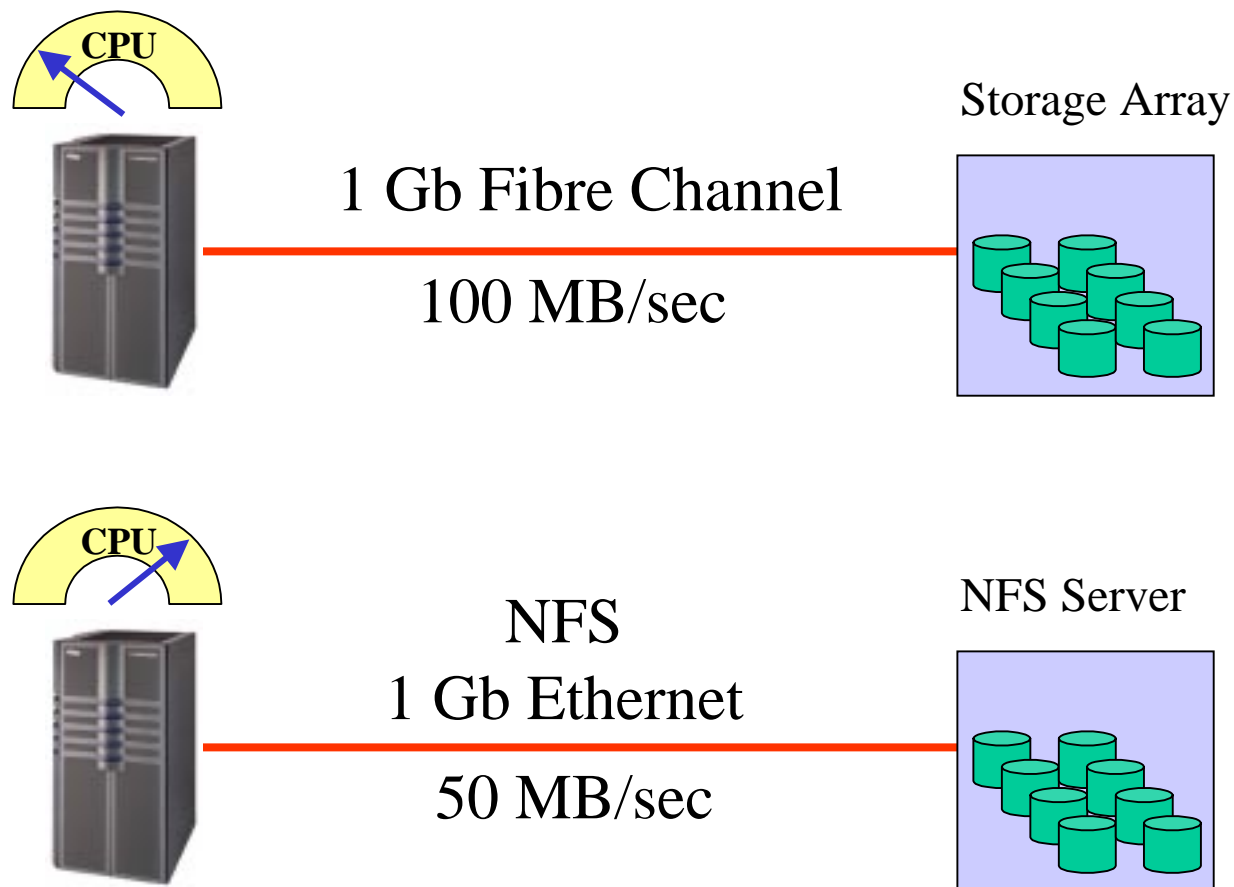
- Low latency - nodes are meters apart
- High bandwidth - runs are short, cheap
- Low error rate
- Simple network
- Physical security
- Tightly configured & controlled





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

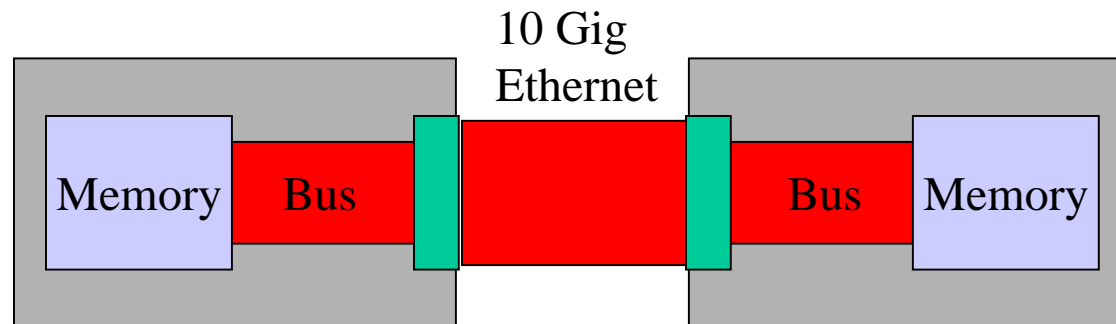
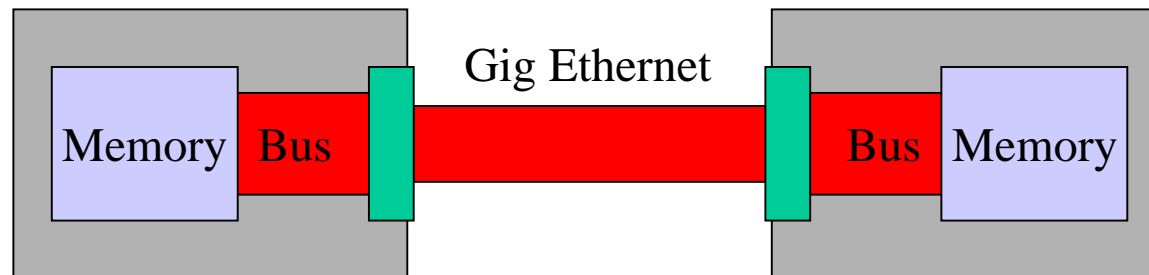
A Problem: Data Center Performance





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

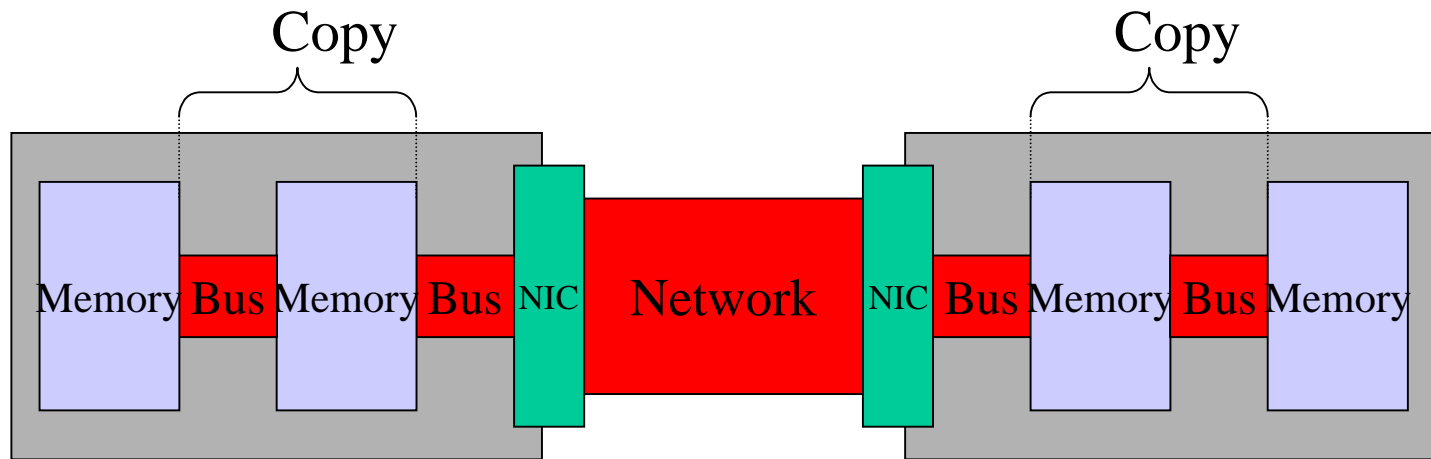
Network vs Bus Performance





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

Network vs Bus Performance

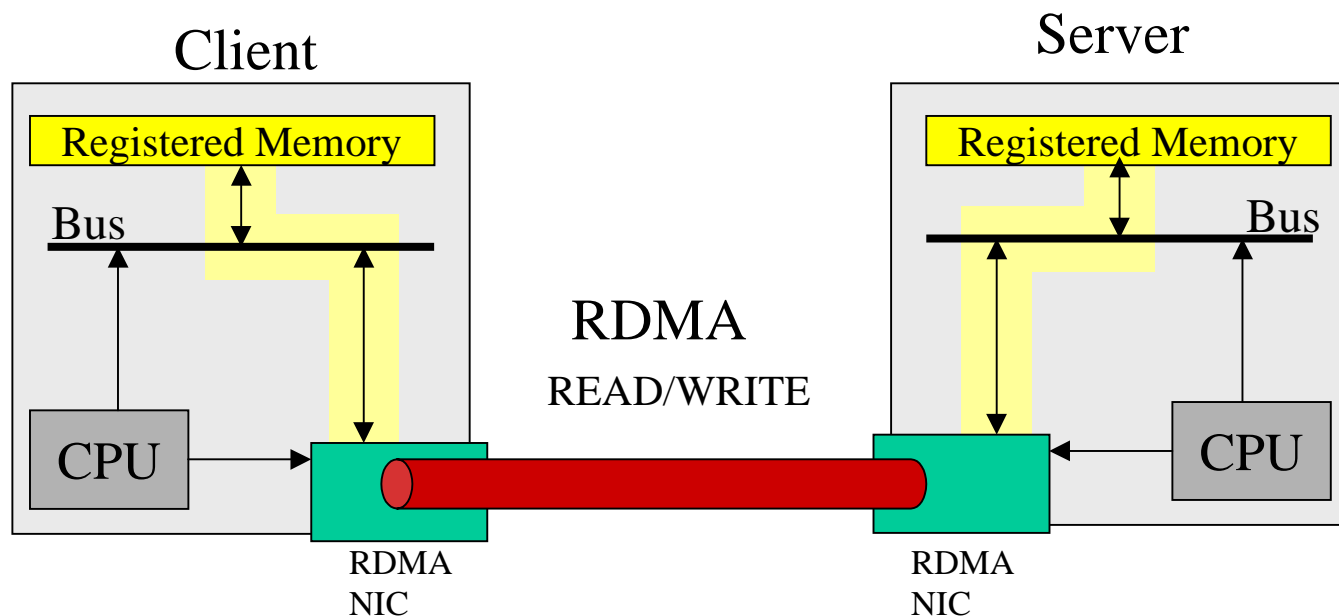


Latency gets worse with each memory copy
which further loads the CPU.



What is RDMA ?

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





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

RDMA Sweet Spot

- High bandwidth
 - > 1Gb links
- Big chunks of data
 - More than 1KB
- Short distance (low latency)
 - 10's of Meters
- Busy CPU
- Protocols
 - NFS, replication, database, backup





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

RDMA Operations

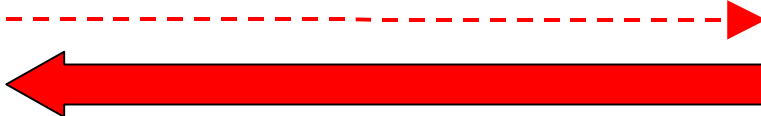
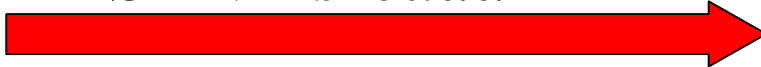
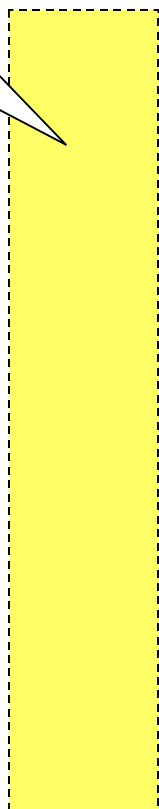
Registered memory
Region referenced
by handle,
Address, length

No
Completion

SEND *srcaddr*

WRITE *srcaddr, dstaddr*

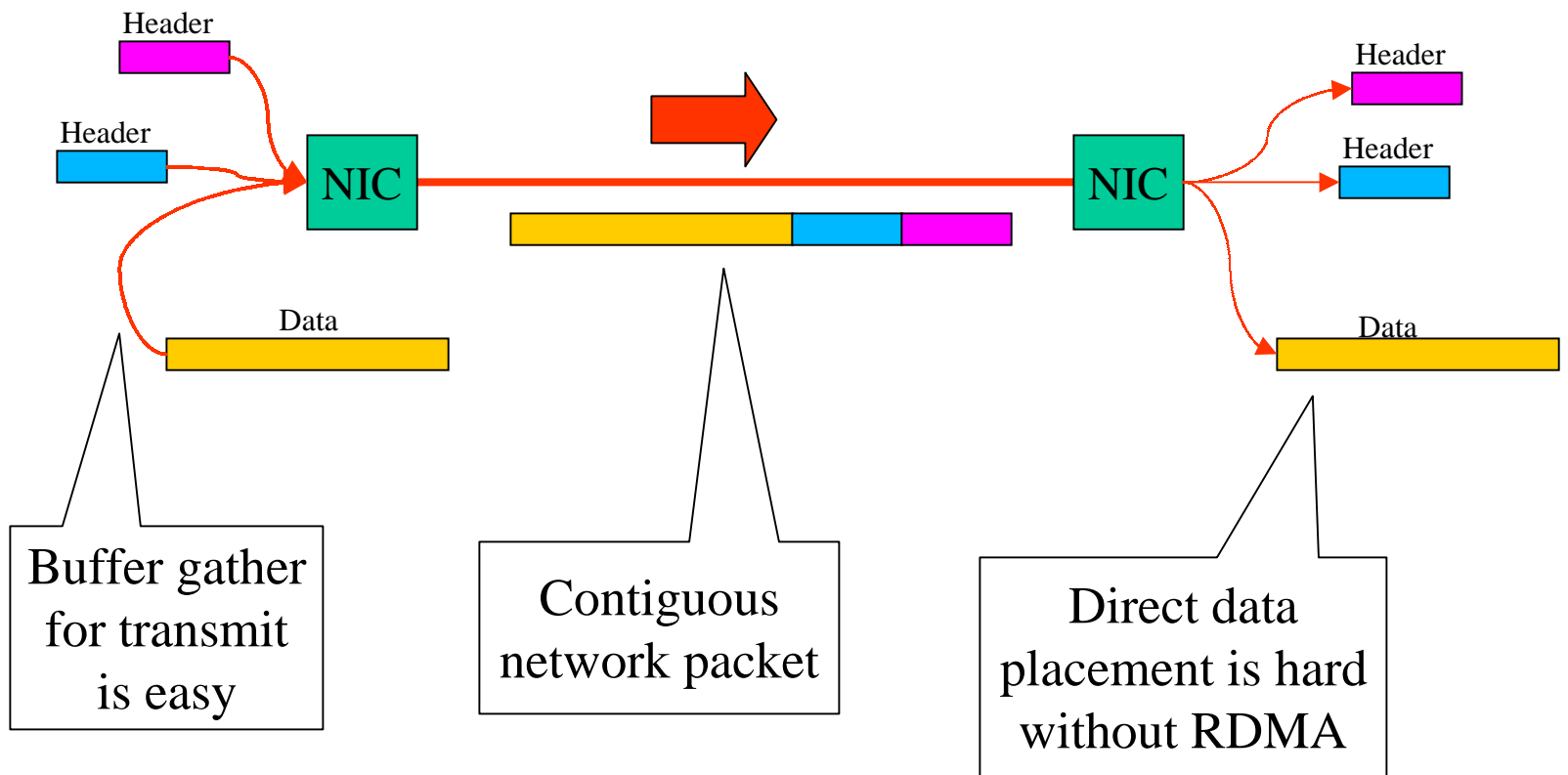
READ *srcaddr, dstaddr*





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

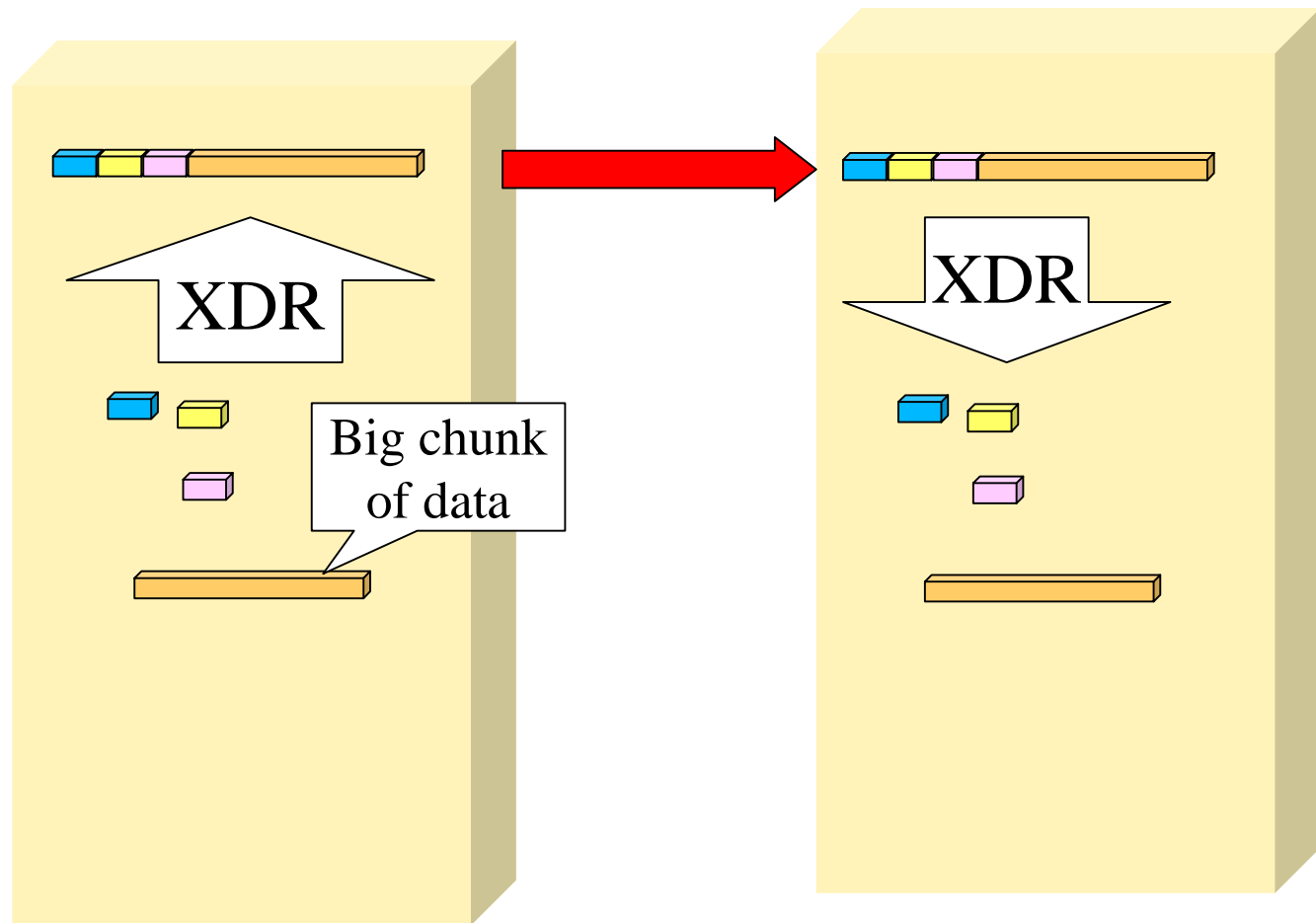
RDMA Provides Direct Data Placement





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

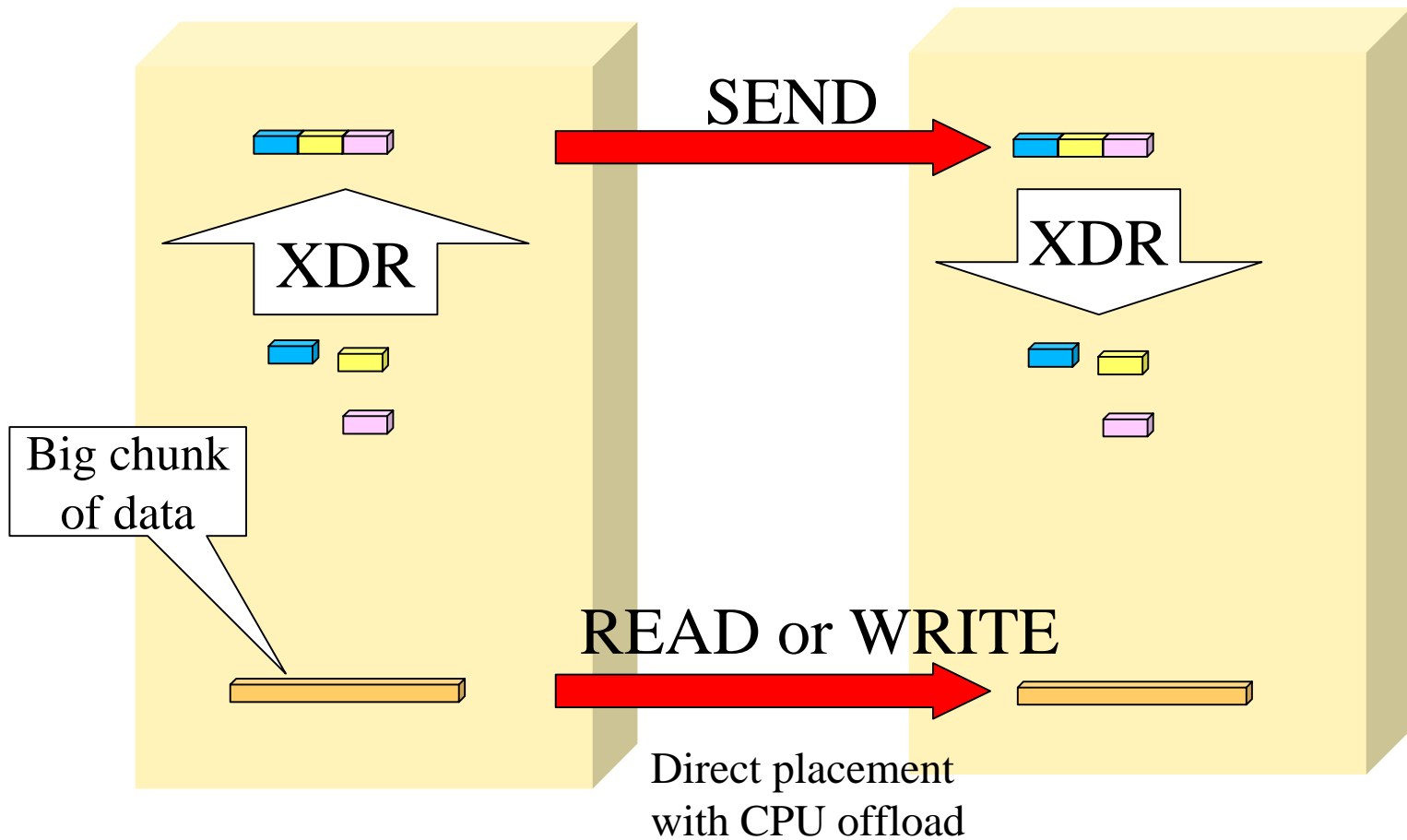
Conventional RPC Data Movement





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

RDMA RPC Data Movement

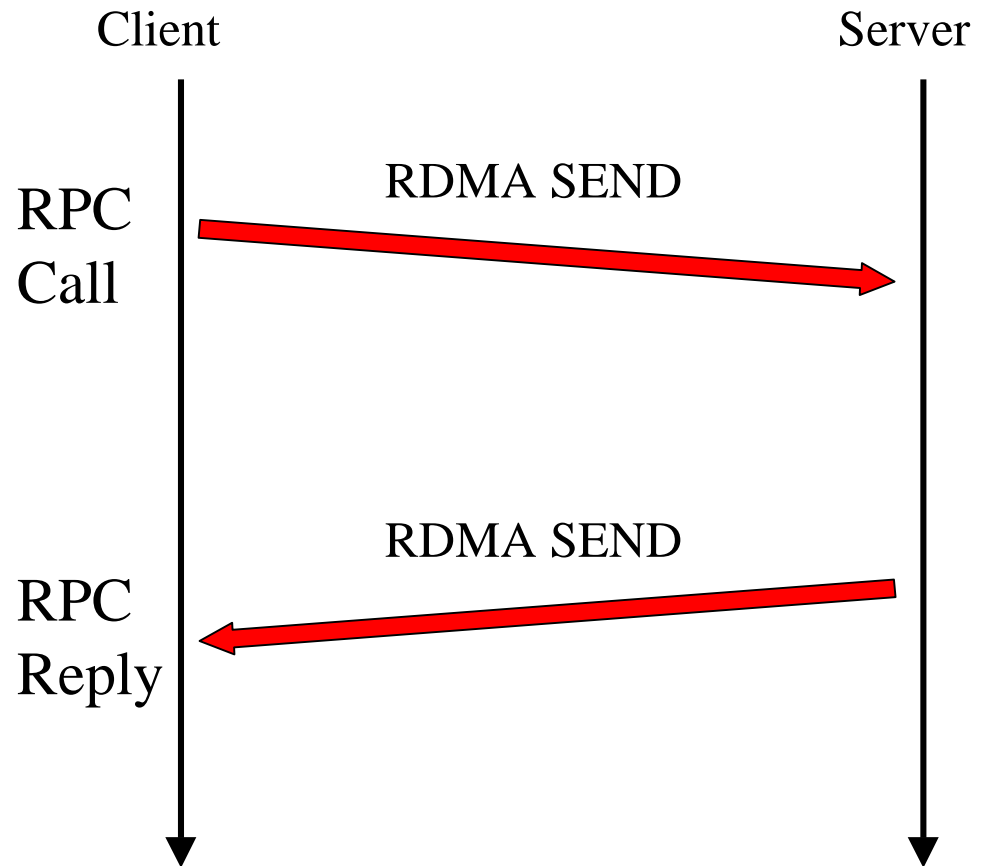




**N I C
F N O
S D N
U S F
T R E
R Y N
C E**

Small RPC Messages

Most RPC Messages are small.
Examples:
LOOKUP
GETATTR



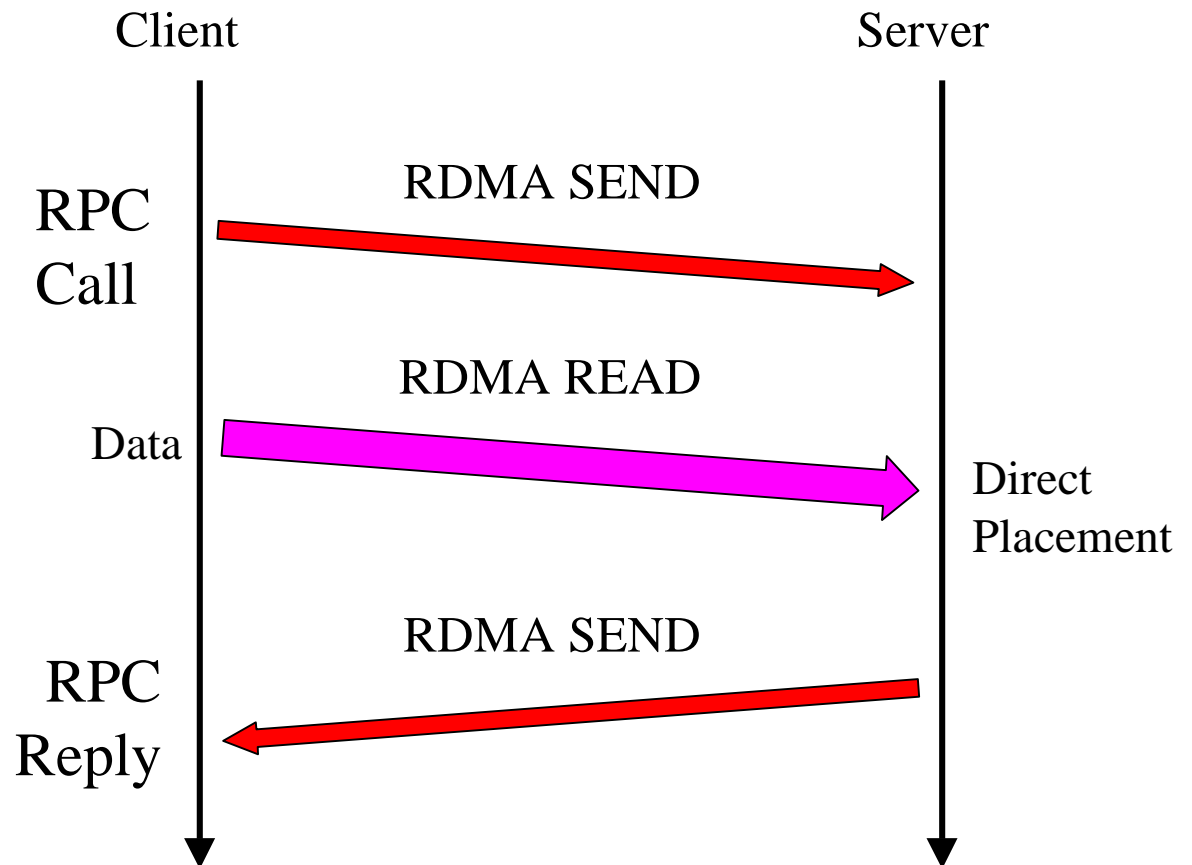


**N I C
F N D
S U S
T R E
R Y N
C E**

Big RPC Call



Example:
WRITE



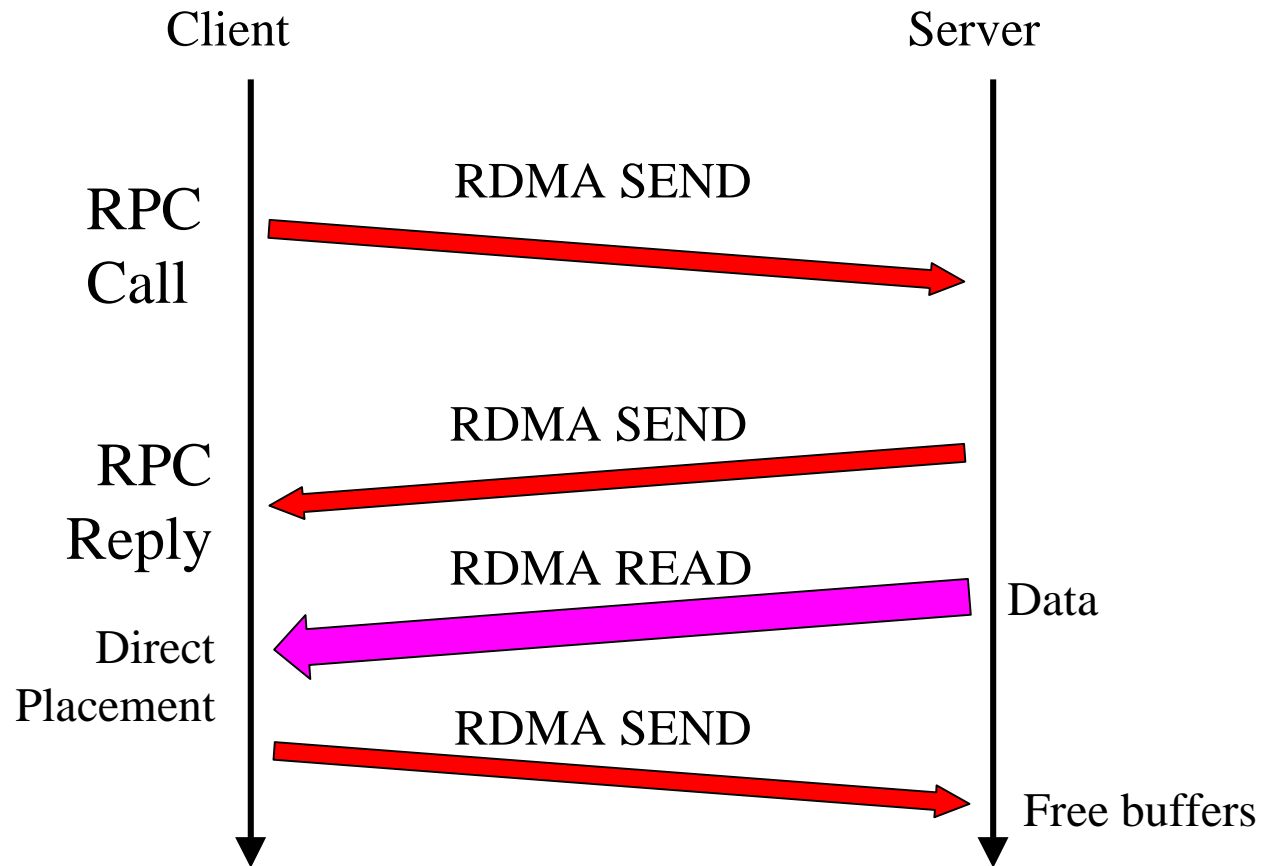


**N I C
F N D
S U S
T R E
R Y N
C E**

Big RPC Reply



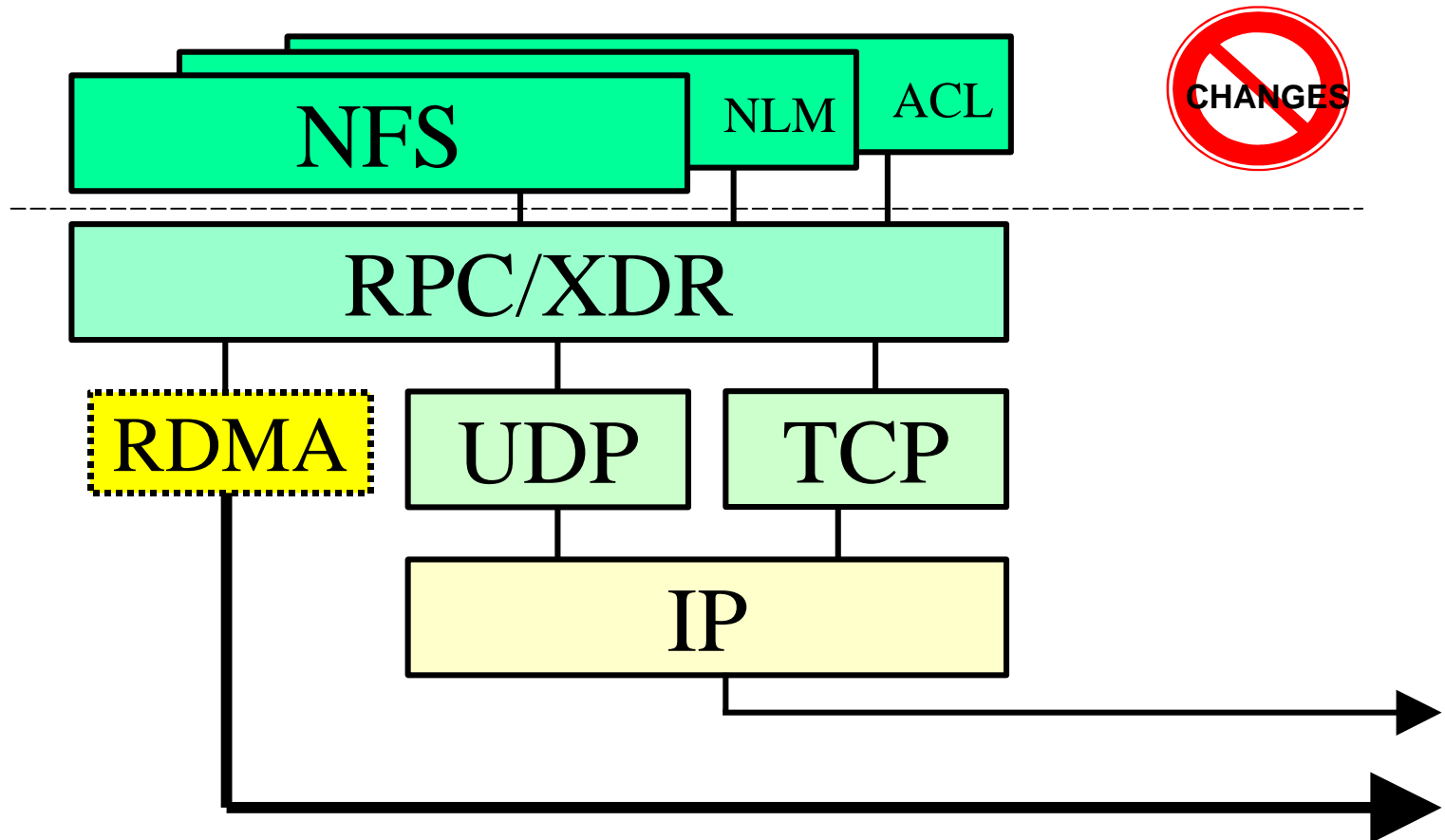
Example:
READ





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

Adapting NFS to RDMA





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

RDMA Flavors



- On Ethernet
 - Emulex GN9000/VI - VI/TCP via 1Gb fibre
 - IETF RDDP Working Group
 - An interoperable, Internet standard
 - RDMA via a Direct Data Placement (DDP) layer
 - Defined for SCTP and TCP
- On Infiniband
 - Supported natively by all IB hardware
- Other
 - Myrinet, Fibre channel, CLan, ...



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

Solaris Prototype



- Extension to Solaris kernel RPC
 - Supports all NFS versions, 2, 3 & 4
 - All kernel RPC: Lock Manager, NFS_ACL
- Behaves like a normal NFS mount
 - only a LOT faster!
- Supporting two RDMA flavors
 - kVIPL with Emulex GN9000/VI over Gigabit Ethernet
 - Infiniband



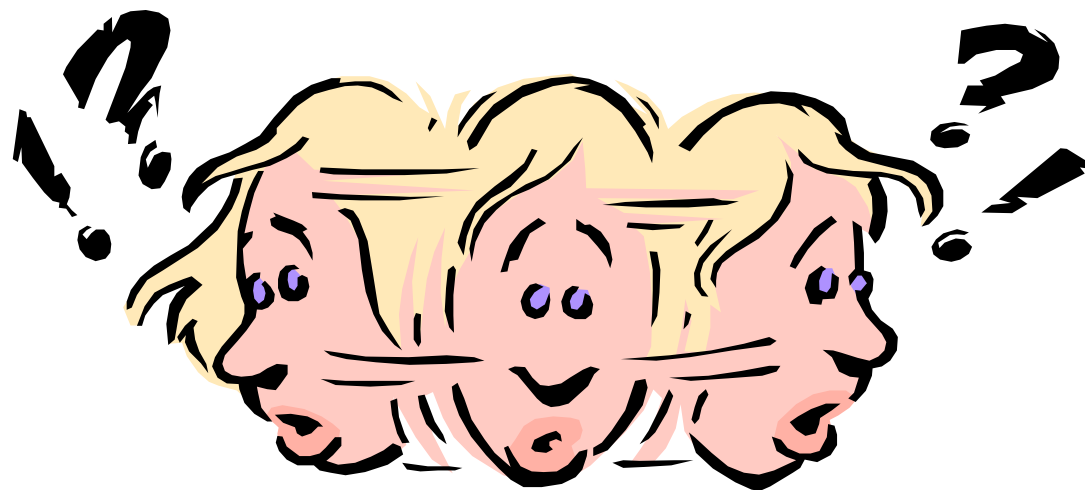
**N I C
F N O
S D N
U F
S E
T R
R E
Y N
C
E**

NFS/RDMA Standard

- NFS will continue to be an open, interoperable protocol!
- Need to publish standards for doing NFS/RDMA via Ethernet & Infiniband
- Storage Network Industry Association (SNIA)
 - NFS/RDMA Technical working group
- Abstract NFS/RDMA protocol
 - based on generic RDMA operations & features
- Transport mappings
 - for NFS/RDDP & NFS/Infiniband



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



Questions & Answers