# Message Passing Multiprocessors

# Programming Models Realized by Protocols

CAD          Database          Scientific modeling          Parallel applications

Multiprogramming          Shared address          Message passing          Data parallel          Programming models

Compilation or library

Operating systems support

Communication abstraction
User/system boundary

Communication hardware

Hardware/software boundary

Physical communication medium

**Network Transactions**

# Network Transaction Primitive

Communication Network

serialized msg

output buffer
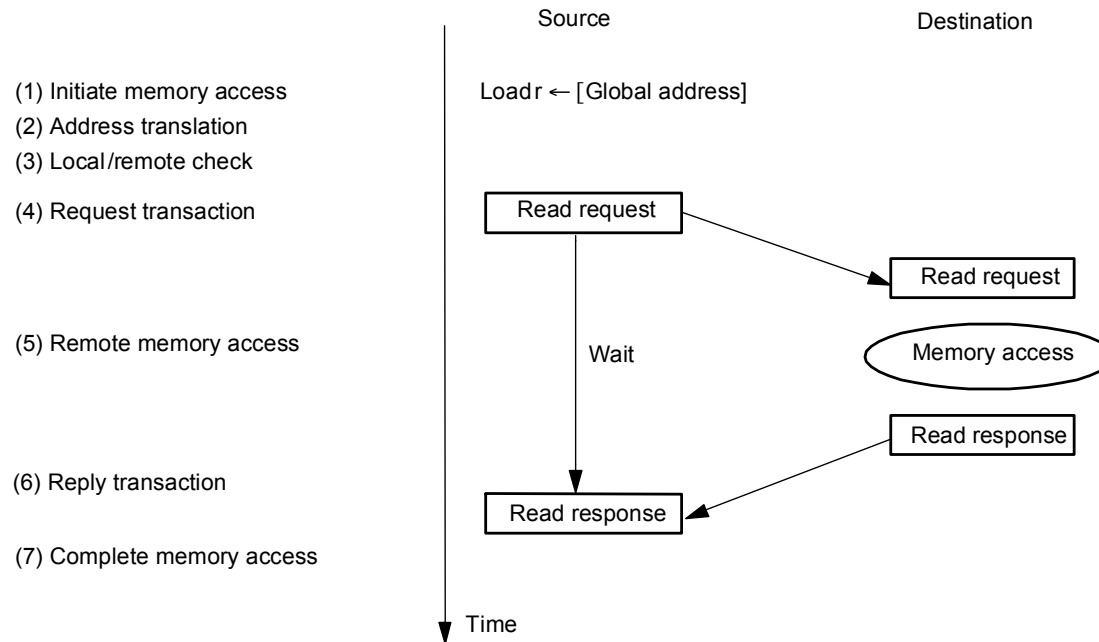
Source Node

input buffer

Destination Node

one-way transfer of information from a source output buffer to a dest. input buffer
  · causes some action at the destination
  · occurrence is not directly visible at source
deposit data, state change, reply

# Shared Address Space Abstraction

Source                                    Destination

(1) Initiate memory access          Load r ← [Global address]

(2) Address translation

(3) Local/remote check

(4) Request transaction              | Read request |

                                                          | Read request |

(5) Remote memory access                      Wait          ( Memory access )

                                                                          | Read response |

(6) Reply transaction                | Read response |

(7) Complete memory access

Time

Fundamentally a two-way request/response protocol
  · writes have an acknowledgement
Issues
  · fixed or variable length (bulk) transfers
  · remote virtual or physical address, where is action performed?
  · deadlock avoidance and input buffer full
coherent?  consistent?

# Message passing

Bulk transfers

Complex synchronization semantics

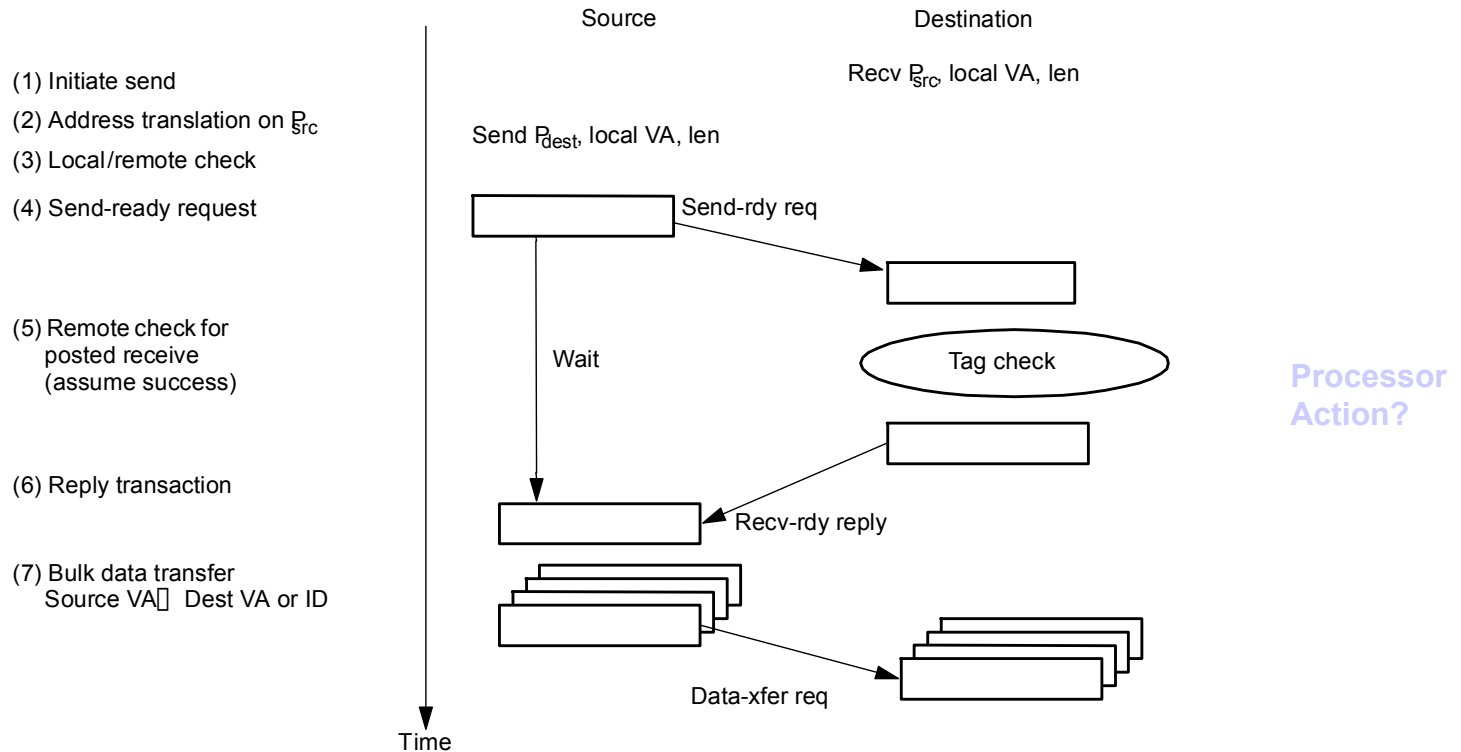- more complex protocols
- More complex action

Synchronous

- Send completes after matching recv and source data sent
- Receive completes after data transfer complete from matching send

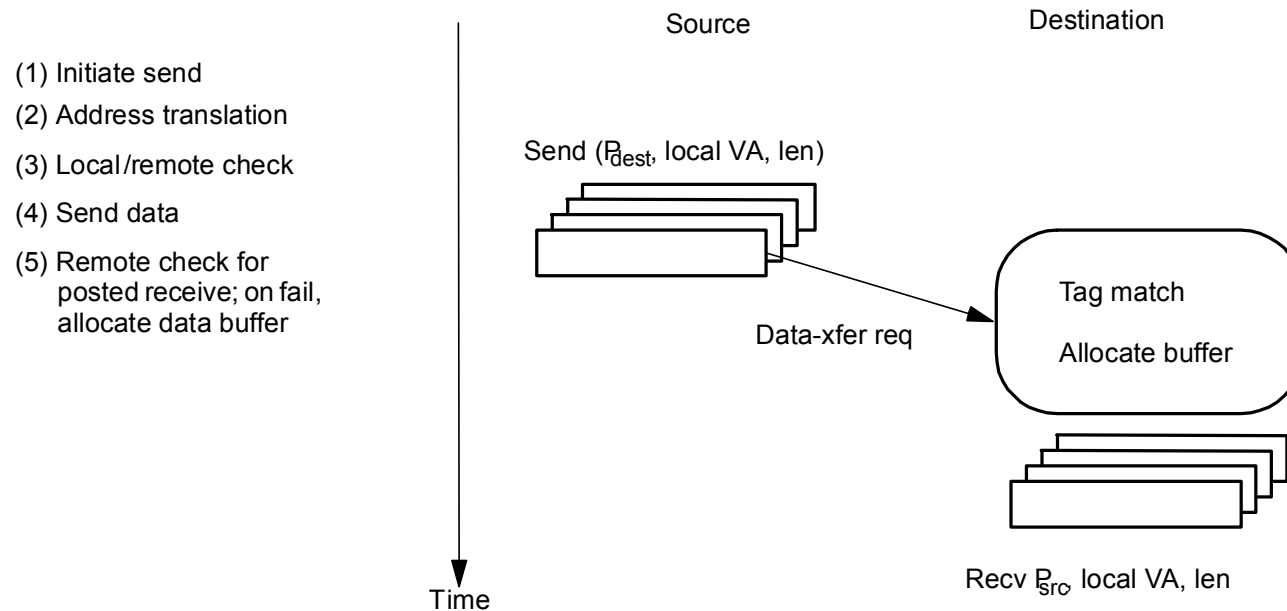Asynchronous

- Send completes after send buffer may be reused

# Synchronous Message Passing

Source             Destination

Recv $P_{src}$, local VA, len

(1) Initiate send

(2) Address translation on $P_{src}$

Send $P_{dest}$, local VA, len

(3) Local/remote check

(4) Send-ready request          Send-rdy req

(5) Remote check for
posted receive         Wait            Tag check
(assume success)                                          **Processor
Action?**

(6) Reply transaction                            Recv-rdy reply

(7) Bulk data transfer
Source VA ☐ Dest VA or ID

Data-xfer req

Time

Constrained programming model.

Deterministic!    What happens when threads added?

# Asynch. Message Passing: Optimistic

Source          Destination

(1) Initiate send

(2) Address translation

(3) Local/remote check

(4) Send data

(5) Remote check for
    posted receive; on fail,
    allocate data buffer

Send ($P_{dest}$, local VA, len)

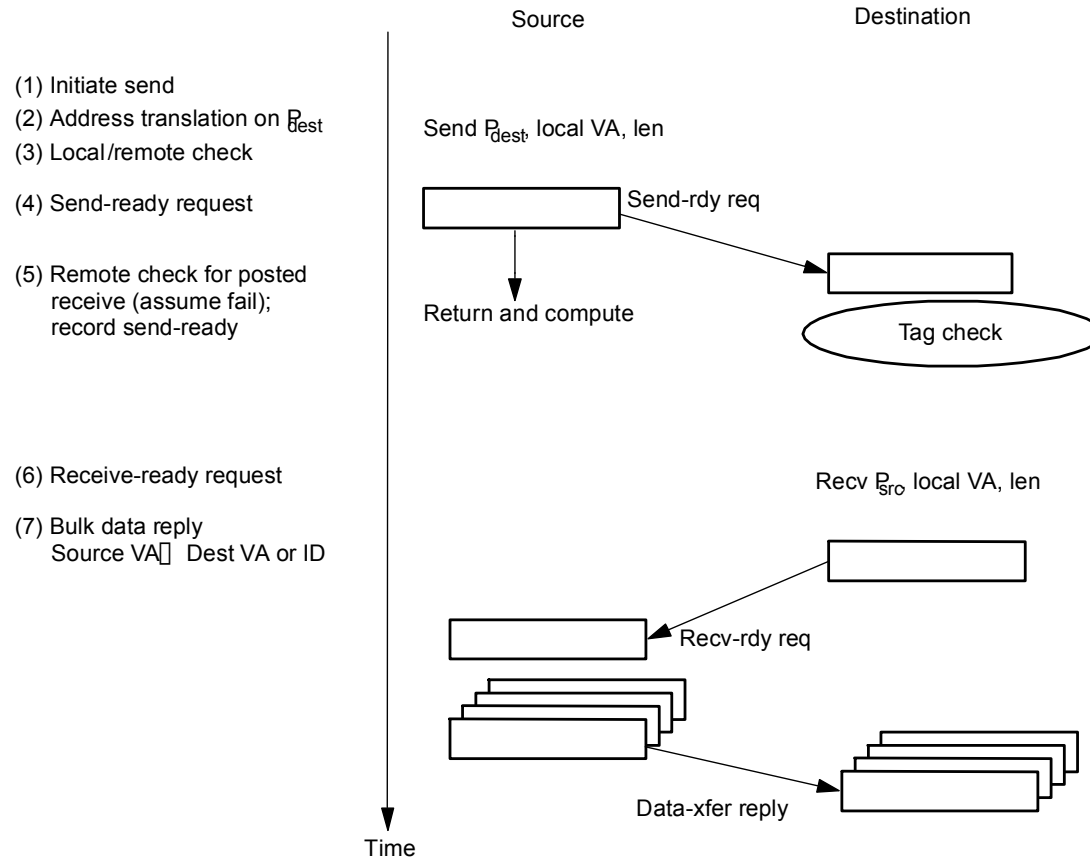Data-xfer req

Tag match

Allocate buffer

Recv $P_{src}$, local VA, len

Time

More powerful programming model

Wildcard receive => non-deterministic

Storage required within msg layer?

# Asynch. Msg Passing: Conservative

Source                                    Destination

(1) Initiate send
(2) Address translation on $P_{dest}$          Send $P_{dest}$, local VA, len
(3) Local/remote check

(4) Send-ready request                         Send-rdy req

(5) Remote check for posted
    receive (assume fail);          Return and compute
    record send-ready                              Tag check

(6) Receive-ready request                      Recv $P_{src}$ local VA, len

(7) Bulk data reply
    Source VA □ Dest VA or ID

                                               Recv-rdy req

                                               Data-xfer reply

                                Time

Where is the buffering?

Contention control?  Receiver initiated protocol?

Short message optimizations

# Key Features of Msg Passing Abstraction

Source knows send data address, dest. knows receive data address

- after handshake they both know both

Arbitrary storage "outside the local address spaces"

- may post many sends before any receives
- non-blocking asynchronous sends reduces the requirement to an arbitrary number of descriptors
  - fine print says these are limited too

Fundamentally a 3-phase transaction

- includes a request / response
- can use optimisitic 1-phase in limited "Safe" cases
  - credit scheme

# Active Messages

Request

handler

Reply

handler

User-level analog of network transaction

- transfer data packet and invoke handler to extract it from the network and integrate with on-going computation

Request/Reply

Event notification: interrupts, polling, events?

May also perform memory-to-memory transfer

# Network Transaction Processing



Scalable Network

Message

**Output Processing**
– checks
– translation
– formating
– scheduling

CA    Communication Assist    CA

Node Architecture

M    P

M    P

**Input Processing**
– checks
– translation
– buffering
– action

Key Design Issue:

How much interpretation of the message?

How much dedicated processing in the Comm. Assist?

# Hardware Capabilities of Network Interface

- Just DMA Support
  - E.g. Normal LAN interfaces, nCUBE/2, …

- Direct CPU Access, with fine-grained communication
  - CM-5, Monsoon, J-machine, …

- Dedicated Hardware
  - E.g. Intel Paragon, Myrinet, DEC PCI Memory Channel, …

# Net Transactions: Physical DMA



DMA controlled by regs, generates interrupts

Physical => OS initiates transfers

Send-side

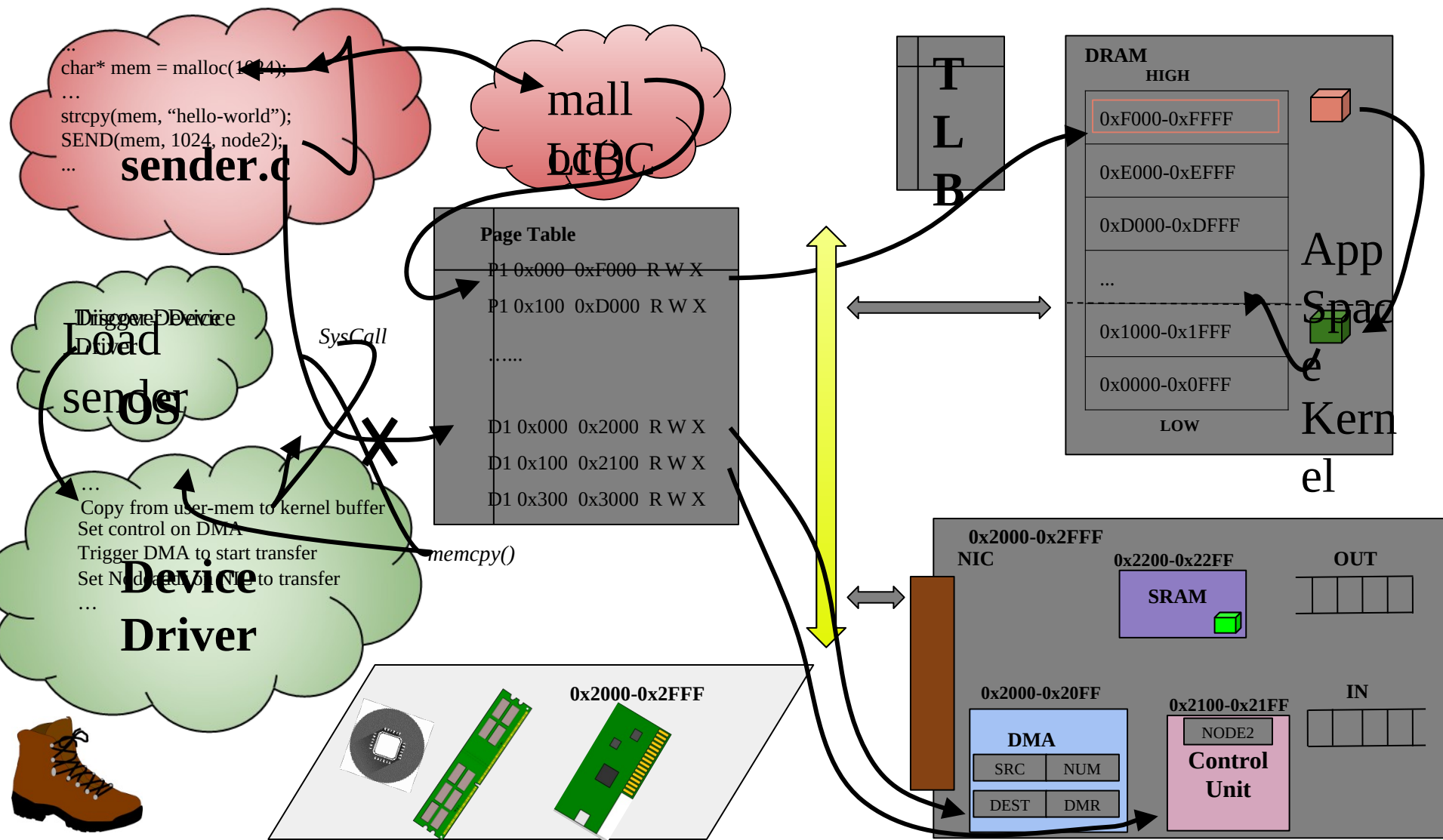· construct system "envelope" around user data in kernel area

Receive

· must receive into system buffer, since no interpretation inCA

# Conventional LAN Network Interface

**Host Memory**

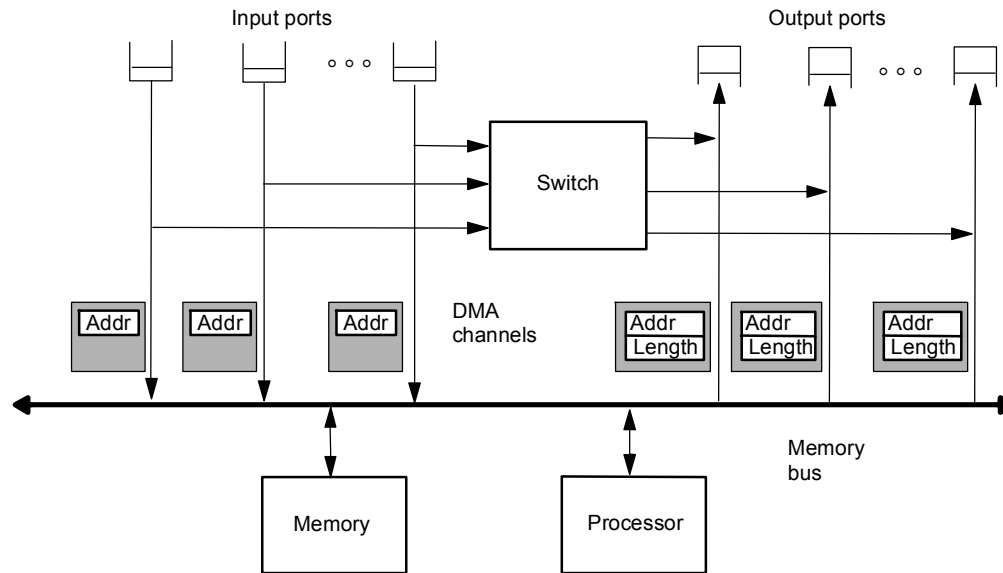**NIC**

**Data**

**NIC Controller**

**trncv**

**TX**
**RX**

**addr**
**len**

**DMA**

| Addr | Len |
|------|-----|
| Status | |
| Next | |

| Addr | Len |
|------|-----|
| Status | |
| Next | |

| Addr | Len |
|------|-----|
| Status | |
| Next | |

| Addr | Len |
|------|-----|
| Status | |
| Next | |

| Addr | Len |
|------|-----|
| Status | |
| Next | |

| Addr | Len |
|------|-----|
| Status | |
| Next | |

**mem bus**

**IO Bus**

**Proc**

..
char* mem = malloc(1024);
…
strcpy(mem, "hello-world");
SEND(mem, 1024, node2);
…

**sender.c**

mall
oc()

**T L B**

**DRAM**
HIGH

0xF000-0xFFFF

0xE000-0xEFFF

0xD000-0xDFFF

...

0x1000-0x1FFF

0x0000-0x0FFF

LOW

App
Spac
e
Kern
el

**Page Table**

P1 0x000  0xF000  R W X

P1 0x100  0xD000  R W X

…...

D1 0x000  0x2000  R W X

D1 0x100  0x2100  R W X

D1 0x300  0x3000  R W X

*SysCall*

Trigger Device
Load
sender
OS

Device
Driver

…
Copy from user-mem to kernel buffer
Set control on DMA
Trigger DMA to start transfer
Set Node2 via NIC to transfer
…

*memcpy()*

0x2000-0x2FFF

0x2000-0x2FFF
NIC

0x2200-0x22FF

**SRAM**

OUT

0x2000-0x20FF

**DMA**

| SRC | NUM |
| DEST | DMR |

0x2100-0x21FF

NODE2

**Control Unit**

IN

15

..
char* mem = malloc(1024);
…
RECV(mem, 1024, node1);
printf("%s\n", mem);
...

**recv.c**

mall
oc()
LIBC

**T
L
B**

**DRAM**
HIGH

0xF000-0xFFFF

0xE000-0xEFFF

0xD000-0xDFFF

...

0x1000-0x1FFF

0x0000-0x0FFF

LOW

App
Spac
e
Kern
el

**Page Table**

P1 0x000  0xF000  R W X

P1 0x100  0xD000  R W X

…...

D1 0x000  0x2000  R W X

D1 0x100  0x2100  R W X

D1 0x300  0x3000  R W X

Trigger Device
Driver on
memc
service arrival
pyOS

*SysCall*

*Interrupt*

...
Interrupt OS to notify message arrival

Set control on DMA
Trigger DMA to trf to Kernel
Buffers
...

**Device
Driver**

*memcpy()*

0x2000-0x2FFF

0x2000-0x2FFF
NIC

0x2200-0x22FF

**SRAM**

OUT

0x2000-0x20FF

**DMA**

SRC | NUM

DEST | DMR

0x2100-0x21FF

**Control
Unit**

NODE2

IN

16

# **Software Stack**

Sending

- Trap to OS    - Why????

- Copy data to Pinned Memory    - Why???

- Wait for DMA to become available

- Program DMA (addr, length, ….)

Receiving

- If you don't know where to put it, interrupt host CPU

- OS programs DMA to copy to pinned memory    - Why???

- On subsequent application RECEIVE call, copy from system to user buffer

Overhead can run to dozens/hundreds of Microseconds

# nCUBE Network Interface



independent DMA channel per link direction
- leave input buffers always open
- segmented messages

routing interprets envelope
- dimension-order routing on hypercube
- bit-serial with 36 bit cut-through

| Os | 16 ins | 260 cy |
|----|--------|--------|
|    | 13 us  |        |
| Or | 18     | 200 cy |
|    | 15 us  |        |
| - includes interrupt | | |

18

# Kernel-based Communication

150-200 usecs

Send

Recv

Copy to OS Buffer

Copy to Appln.

NI

NI

Program NI

Program NI

Interrupt

# User-level Communication

15-20 usecs

Send

Recv

NI

NI

# How do we achieve User-Level Access?

Smart Scheduling (CM-5)

Additional Hardware (Paragon, Myrinet, DEC Memory Channel, …)

# Example: CM-5

NI sits on Memory Bus

Input and output FIFO for each network

NI registers are mapped into user-address space

Processor can load/store into these registers

NI can also be made to interrupts

2 data networks

Other Examples: *T integrated NI on chip, J-Machine, iWARP



| Os | 50 cy | 1.5 us |
|----|-------|--------|
| Or | 53 cy | 1.6 us |
| interrupt | | 10us |

# Coordinated Scheduling (Co-Scheduling)

We need to ensure you can only send to or receive from whoever you are authorized to do so.

Source and Destination should have counterparts of same application scheduled in their respective nodes. i.e. coordinated context-switching

But that does not suffice!!!

- We also need coordinated Network context switching

- CM-5 does network context switching through "All Fall Down"

# User-level Messaging with Additional H/W (But NOT specialized H/W)



General Purpose processor performs arbitrary output processing (at system level)

General Purpose processor interprets incoming network transactions (at system level)

User Processor <–> Msg Processor share memory

Msg Processor <–> Msg Processor via system network transaction

# Levels of Network Transaction



User Processor stores cmd / msg / data into shared output queue
- must still check for output queue full (or make elastic)

Communication assists make transaction happen
- checking, translation, scheduling, transport, interpretation

Effect observed on destination address space and/or events

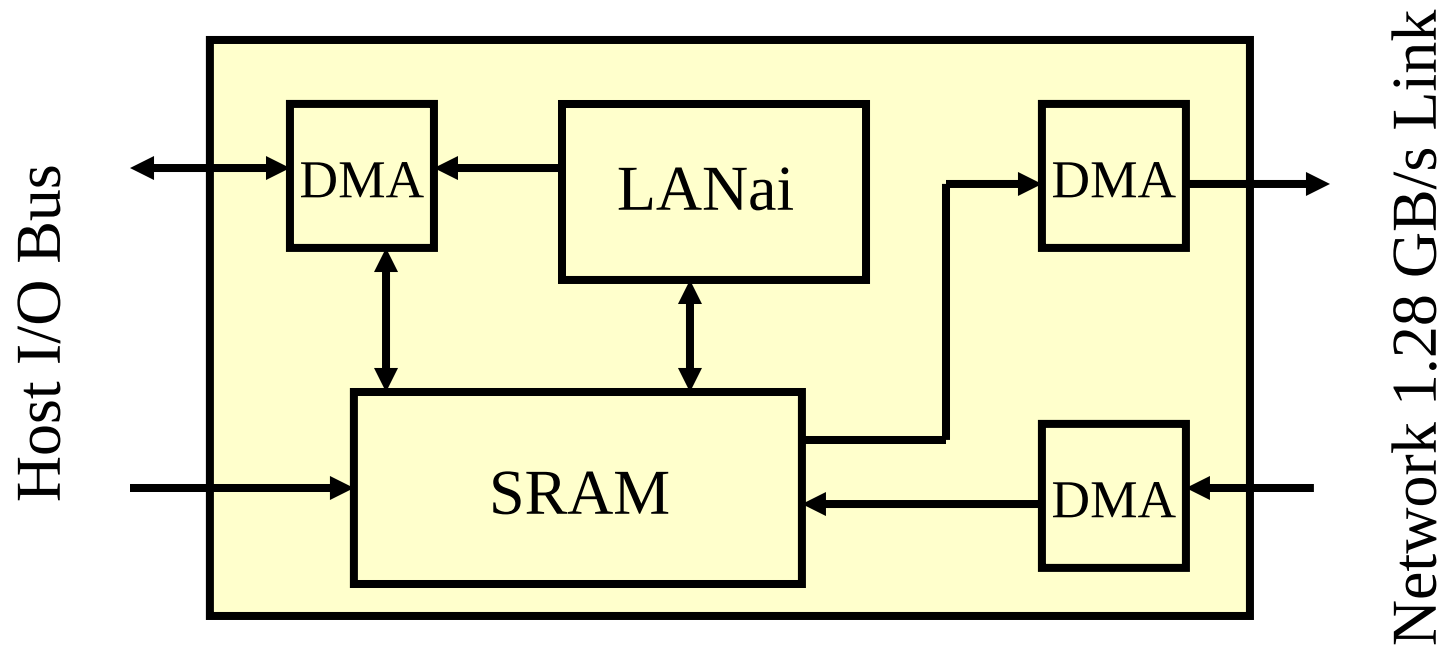Protocol divided between two layers

# Example: Intel Paragon



**Service**

**Network**

**I/O Nodes**

**I/O Nodes**

Devices

Devices

**16**   **175 MB/s Duplex**

**Mem**   **2048 B**   ° ° °   **EOP**

**NI**

**rte**

**MP handler**

**Var data**

**i860xp**
50 MHz
16 KB $
4-way
32B Block
MESI

**64**   **400 MB/s**

**$**   **$**   **sDMA**

**P**   **M P**   **rDMA**

# User Level Abstraction



Any user process can post a transaction for any other in protection domain

- communication layer moves $OQ_{src} \rightarrow IQ_{dest}$
- may involve indirection: $VAS_{src} \rightarrow VAS_{dest}$

# Specialized NICs (Myrinet)
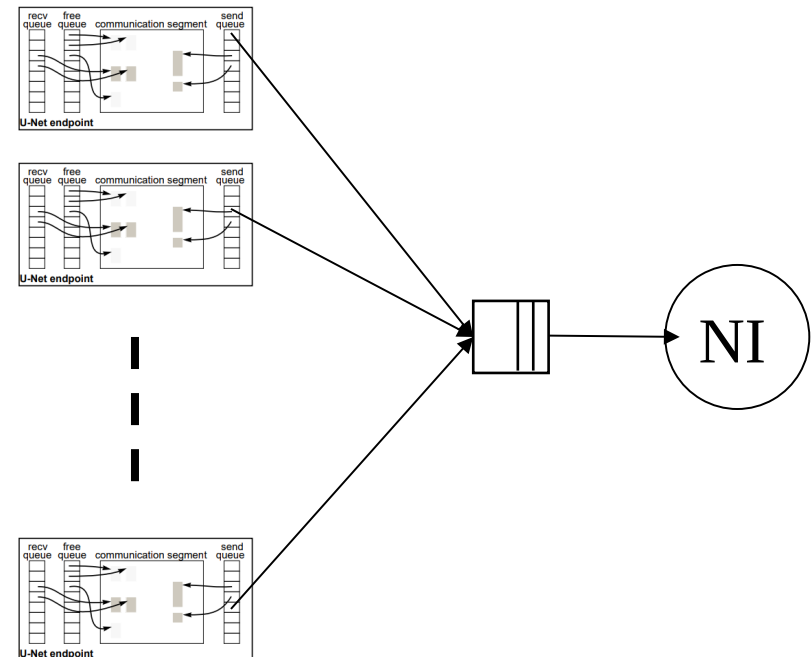
# Data Structures (End Point)

Where should each of the queues be maintained?

- Send – On the NI
- Receive – On the Host
- Free – On the Host



Need to use memory locations for synchronization

Scalability with # of endpoints (Doorbell mechanism + Hardware Combining Queue)

# DEC Memory Channel

"Push-Only" Memory-based message passing

Map part of the Virtual Address Space to the NI
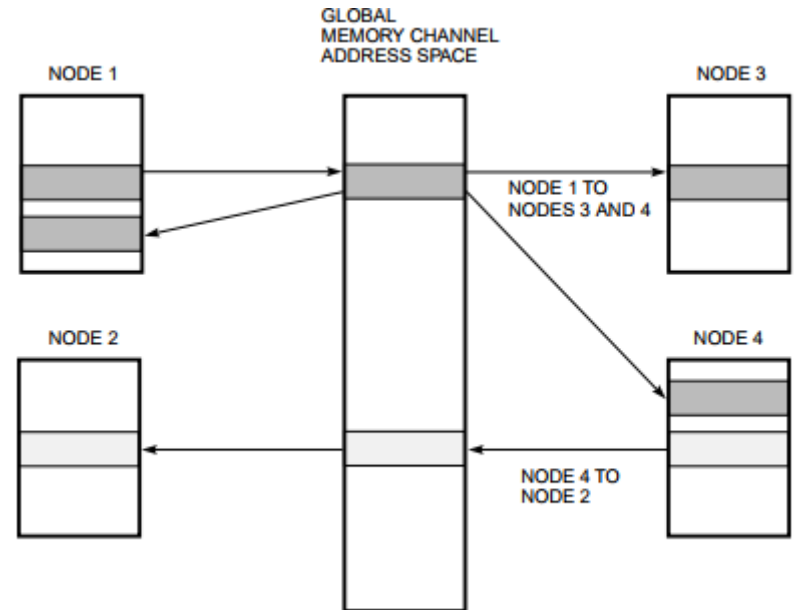
Open as "Outgoing" or "Incoming" channels

Incoming channels allocate a physical frame that is pinned

No frames allocated for outgoing channels.

One-to-Many channels possible

Send is simply a sequence of stores and Receives are a sequence of loads.

Completely avoids additional copies!



GLOBAL MEMORY CHANNEL ADDRESS SPACE

NODE 1

NODE 3

NODE 1 TO NODES 3 AND 4

NODE 2

NODE 4

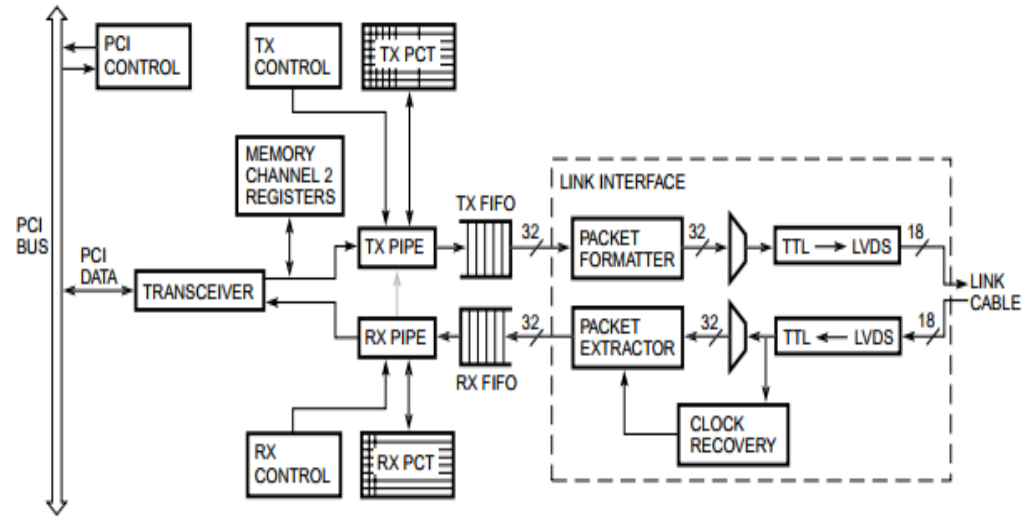NODE 4 TO NODE 2

# DEC Memory Channel

**Send**:

CPU Stores words into outgoing channel (page in VA)

Write-Buffer in CPU coalesces words (resulting in a Cache block write)

NIC uses Physical Frame # to Index PCT

Get Descriptor for Destination  (route, control bits, Frame # at receiver)

Compose Packet and Send



**Receive**:

Extract Receive Frame #

Use that to index PCT

Verify validity

DMA to Physical Frame

Receive Frames are cacheable (cache coherent transfer)