

Scalable Distributed Memory Multiprocessors

CS 418

Lectures 24-25

Outline

Scalability

- physical, bandwidth, latency and cost
- level of integration

Realizing Programming Models

- network transactions
- protocols
- safety
 - input buffer problem: N-1
 - fetch deadlock

Communication Architecture Design Space

- how much hardware interpretation of the network transaction?

- 2 -

CS 418 S'04

Limited Scaling of a Bus

Characteristic	Bus
Physical Length	~ 1 ft
Number of Connections	fixed
Maximum Bandwidth	fixed
Interface to Comm. medium	memory inf
Global Order	arbitration
Protection	Virt -> physical
Trust	total
OS	single
comm. abstraction	HW

Bus: each level of the system design is grounded in the scaling limits at the layers below and assumptions of close coupling between components

- 3 -

CS 418 S'04

Workstations in a LAN?

Characteristic	Bus	LAN
Physical Length	~ 1 ft	KM
Number of Connections	fixed	many
Maximum Bandwidth	fixed	???
Interface to Comm. medium	memory inf	peripheral
Global Order	arbitration	???
Protection	Virt -> physical	OS
Trust	total	none
OS	single	independent
comm. abstraction	HW	SW

No clear limit to physical scaling, little trust, no global order, consensus difficult to achieve.
Independent failure and restart

- 4 -

CS 418 S'04

Scalable Machines

What are the design trade-offs for the spectrum of machines between?

- specialize or commodity nodes?
- capability of node-to-network interface
- supporting programming models?

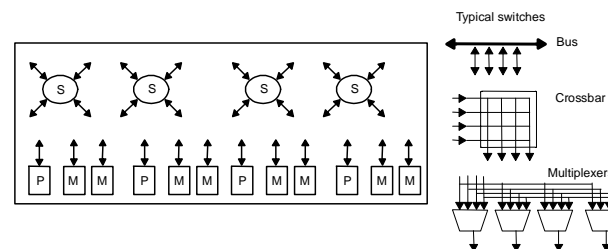
What does scalability mean?

- avoids inherent design limits on resources
- **bandwidth** increases with P
- **latency** does not
- **cost** increases slowly with P

- 5 -

CS 418 S'04

Bandwidth Scalability



What fundamentally limits bandwidth?

- single set of wires

Must have many independent wires

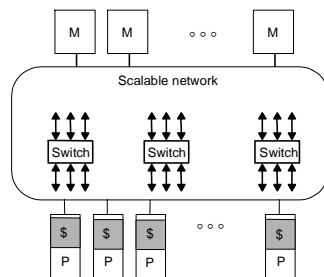
Connect modules through switches

Bus vs Network Switch?

- 6 -

CS 418 S'04

Dancehall MP Organization



Network bandwidth?

Bandwidth demand?

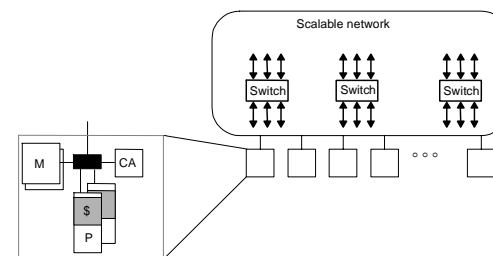
- independent processes?
- communicating processes?

Latency?

- 7 -

CS 418 S'04

Generic Distributed Memory Org.



Network bandwidth?

Bandwidth demand?

- independent processes?
- communicating processes?

Latency?

- 8 -

CS 418 S'04

Key Property

Large # of independent communication paths between nodes

- allow a **large # of concurrent transactions using different wires**

Initiated independently

No global arbitration

Effect of a transaction only visible to the nodes involved

- effects propagated through additional transactions

- 9 -

CS 418 S'04

Latency Scaling

$T(n) = \text{Overhead} + \text{Channel Time} + \text{Routing Delay}$

Overhead?

Channel Time(n) = n/B

- Bandwidth at bottleneck

RoutingDelay(h,n)

- 10 -

CS 418 S'04

Typical Example

max distance: $\log n$

number of switches: $\alpha n \log n$

overhead = 1 us, BW = 64 MB/s, 200 ns per hop

Pipelined

$T_{64}^{p}(128) = 1.0 \text{ us} + 2.0 \text{ us} + 6 \text{ hops} * 0.2 \text{ us/hop} = 4.2$

$T_{1024}^{p}(128) = 1.0 \text{ us} + 2.0 \text{ us} + 10 \text{ hops} * 0.2 \text{ us/hop} = 5.0$

Store and Forward

$T_{64}^{sf}(128) = 1.0 \text{ us} + 6 \text{ hops} * (2.0 + 0.2) \text{ us/hop} = 14.2 \text{ us}$

$T_{1024}^{sf}(128) = 1.0 \text{ us} + 10 \text{ hops} * (2.0 + 0.2) \text{ us/hop} = 23 \text{ us}$

- 11 -

CS 418 S'04

Cost Scaling

$\text{cost}(p,m) = \text{fixed cost} + \text{incremental cost}(p,m)$

Bus Based SMP?

Ratio of processors : memory : network : I/O ?

Parallel efficiency(p) = $\text{Speedup}(P) / P$

Costup(p) = $\text{Cost}(p) / \text{Cost}(1)$

Cost-effective: $\text{speedup}(p) > \text{costup}(p)$

- 12 -

CS 418 S'04

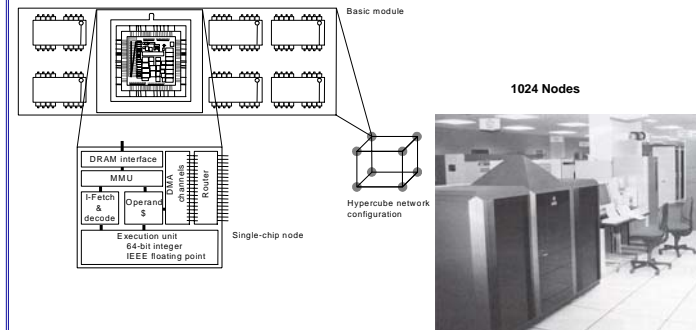
Physical Scaling

Chip-level integration
Board-level integration
System-level integration

- 13 -

CS 418 S'04

nCUBE/2 Machine Organization

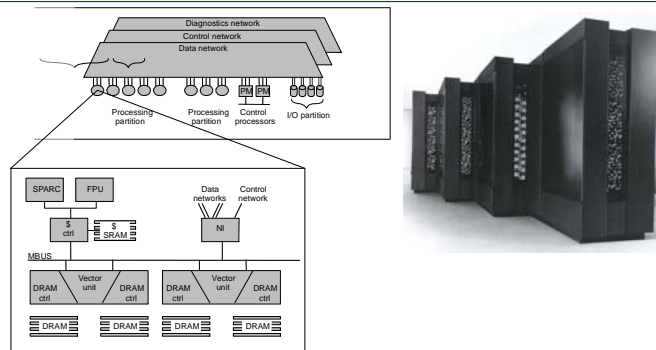


Entire machine synchronous at 40 MHz

- 14 -

CS 418 S'04

CM-5 Machine Organization



Board-level integration

- 15 -

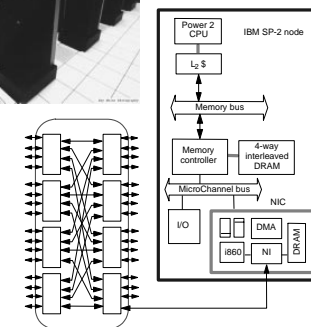
CS 418 S'04

System Level Integration

IBM SP-2



General interconnection network formed from 8-port switches



- 16 -

CS 418 S'04

Outline

Scalability

- physical, bandwidth, latency and cost
- level of integration

Realizing Programming Models

- network transactions
- protocols
- safety
 - input buffer problem: N-1
 - fetch deadlock

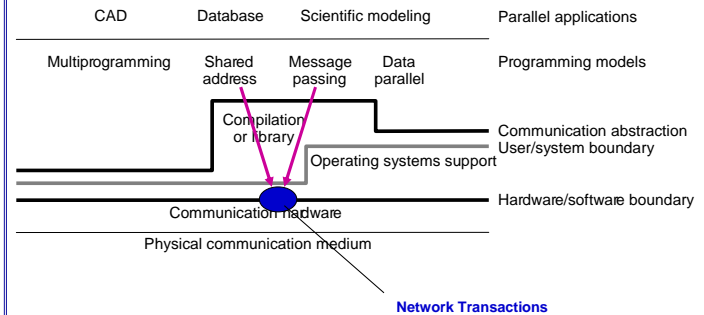
Communication Architecture Design Space

- how much hardware interpretation of the network transaction?

- 17 -

CS 418 S'04

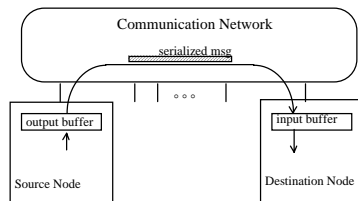
Programming Models Realized by Protocols



- 18 -

CS 418 S'04

Network Transaction Primitive



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

- 19 -

CS 418 S'04

Bus Transactions vs Net Transactions

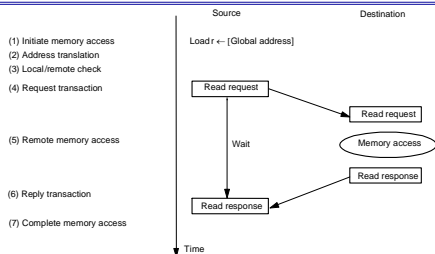
Issues:

protection check	V->P	??
format	wires	flexible
output buffering	reg, FIFO	??
media arbitration	global	local
destination naming and routing		
input buffering	limited	many source
action		
completion detection		

- 20 -

CS 418 S'04

Shared Address Space Abstraction



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?

- 21 -

CS 418 S'04

The Fetch Deadlock Problem

Even if a node cannot issue a request, it must sink network transactions.

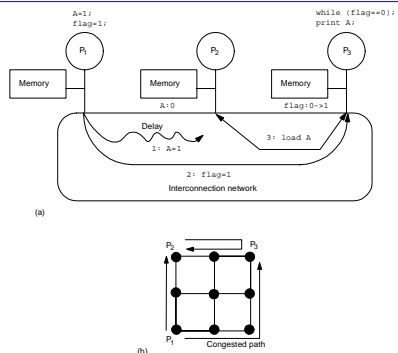
Incoming transaction may be a request, which will generate a response.

Closed system (finite buffering)

- 22 -

CS 418 S'04

Consistency



write-atomicity violated without caching

- 23 -

CS 418 S'04

Key Properties of SAS Abstraction

Source and destination data addresses are specified by the source of the request

- a degree of logical coupling and trust

No storage logically "outside the application address space(s)"

- may employ temporary buffers for transport

Operations are fundamentally request-response

Remote operation can be performed on remote memory

- logically does not require intervention of the remote processor

- 24 -

CS 418 S'04

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

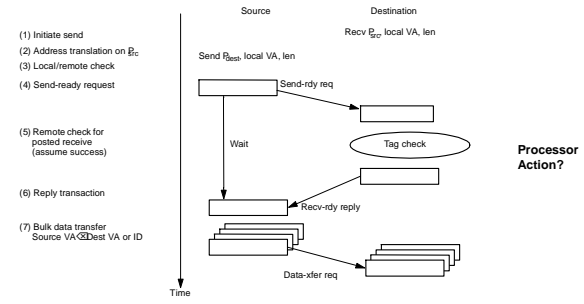
Asynchronous

- Send completes after send buffer may be reused

- 25 -

CS 418 S'04

Synchronous Message Passing

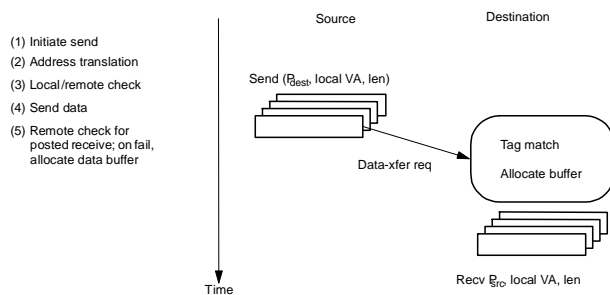


Constrained programming model.
 Deterministic! What happens when threads added?
 Destination contention very limited.
 User/System boundary?

- 26 -

CS 418 S'04

Asynch. Message Passing: Optimistic

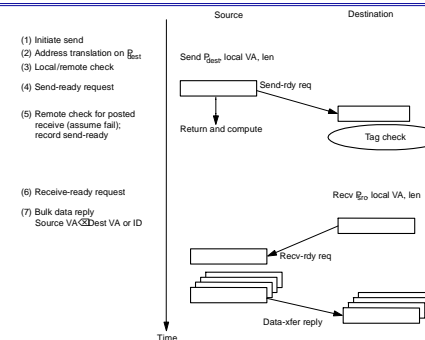


More powerful programming model
 Wildcard receive => non-deterministic
 Storage required within msg layer?

- 27 -

CS 418 S'04

Asynch. Msg Passing: Conservative



Where is the buffering?
 Contention control? Receiver initiated protocol?
 Short message optimizations

- 28 -

CS 418 S'04

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

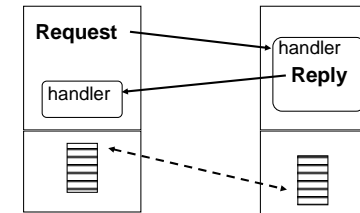
Fundamentally a 3-phase transaction

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

- 29 -

CS 418 S'04

Active Messages



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

- 30 -

CS 418 S'04

Common Challenges

Input buffer overflow

- N-1 queue over-commitment => must slow sources
- Reserve space per source (credit)
 - when available for reuse?
 - » Ack or Higher level
- Refuse input when full
 - backpressure in reliable network
 - tree saturation
 - deadlock free
 - what happens to traffic not bound for congested dest?
- Reserve ack back channel
- Drop packets
- Utilize higher-level semantics of programming model

- 31 -

CS 418 S'04

Challenges (cont)

Fetch Deadlock

- For network to remain deadlock free, nodes must continue accepting messages, even when cannot source msgs
- what if incoming transaction is a request?
 - Each may generate a response, which cannot be sent!
 - What happens when internal buffering is full?

Logically independent request/reply networks

- physical networks
- virtual channels with separate input/output queues

Bound requests and reserve input buffer space

- $K(P-1)$ requests + K responses per node
- service discipline to avoid fetch deadlock?

NACK on input buffer full

- NACK delivery?

- 32 -

CS 418 S'04

Challenges in Realizing Programming Models in the Large

One-way transfer of information

No global knowledge, nor global control

- barriers, scans, reduce, global-OR give fuzzy global state

Very large number of concurrent transactions

Management of input buffer resources

- many sources can issue a request and over-commit destination before any see the effect

Latency is large enough that you are tempted to "take risks"

- optimistic protocols
- large transfers
- dynamic allocation

Many many more degrees of freedom in design and engineering of these system

- 33 -

CS 418 S'04

Summary

Scalability

- physical, bandwidth, latency and cost
- level of integration

Realizing Programming Models

- network transactions
- protocols
- safety
 - input buffer problem: N-1
 - fetch deadlock

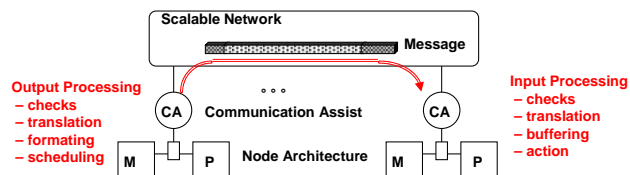
Communication Architecture Design Space

- how much hardware interpretation of the network transaction?

- 34 -

CS 418 S'04

Network Transaction Processing



Key Design Issues:

- How much interpretation of the message?
- How much dedicated processing in the Comm. Assist?

- 35 -

CS 418 S'04

Spectrum of Designs

Increasing HW Support, Specialization, Intrusiveness, Performance (???)

None: Physical bit stream

- blind, physical DMA

nCUBE, iPSC, . . .

User/System

- User-level port
- User-level handler

CM-5, *T
J-Machine, Monsoon, . . .

Remote virtual address

- Processing, translation

Paragon, Meiko CS-2

Global physical address

- Proc + Memory controller

RP3, BBN, T3D

Cache-to-cache

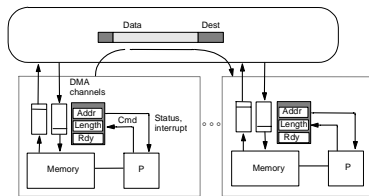
- Cache controller

Dash, KSR, Flash

- 36 -

CS 418 S'04

Net Transactions: Physical DMA



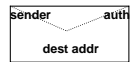
DMA controlled by regs, generates interrupts
Physical => OS initiates transfers

Send-side

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

Receive

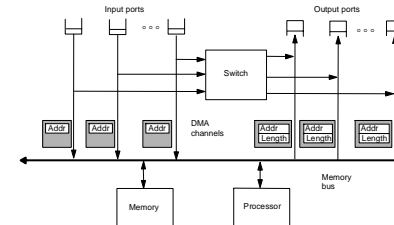
- must receive into system buffer, since no interpretation in CA



- 37 -

CS 418 S'04

nCUBE Network Interface



independent DMA channel per link direction

- leave input buffers always open
- segmented messages

routing interprets envelope

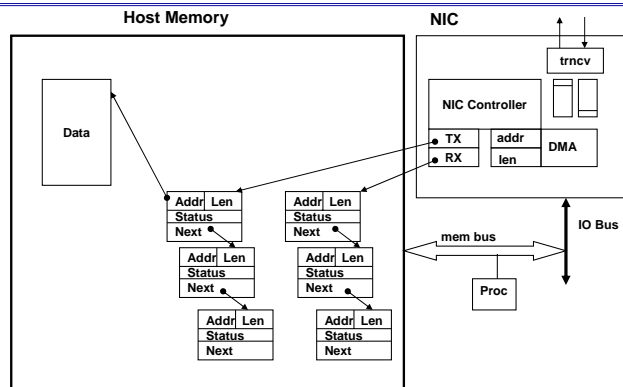
- dimension-order routing on hypercube
- bit-serial with 36 bit cut-through

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

- 38 -

CS 418 S'04

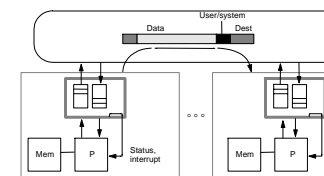
Conventional LAN Network Interface



- 39 -

CS 418 S'04

User Level Ports



initiate transaction at user level

deliver to user without OS intervention

network port in user space

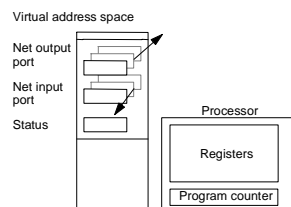
User/system flag in envelope

- protection check, translation, routing, media access in src CA
- user/sys check in dest CA, interrupt on system

- 40 -

CS 418 S'04

User Level Network ports



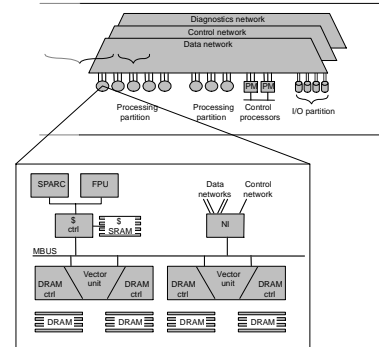
Appears to user as logical message queues plus status
 What happens if no user pop?

- 41 -

CS 418 S'04

Example: CM-5

Input and output FIFO for each network
 2 data networks
 tag per message
 · index NI mapping table
 context switching?



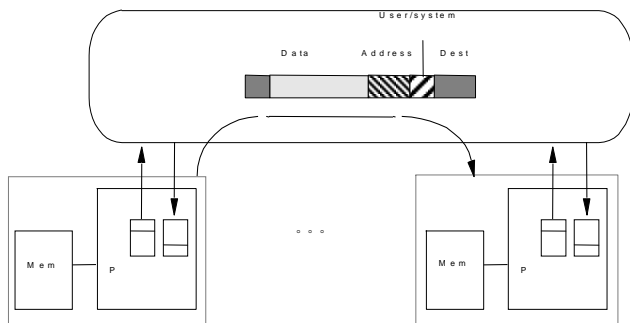
*T integrated NI on chip
 iWARP also

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

- 42 -

CS 418 S'04

User Level Handlers

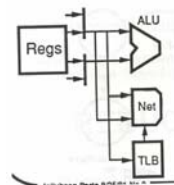


Hardware support to vector to address specified in message
 · message ports in registers

- 43 -

CS 418 S'04

J-Machine



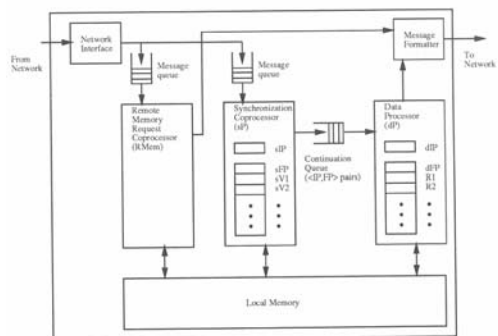
Each node a small msg driven processor
 HW support to queue msgs and dispatch to msg handler task



- 44 -

CS 418 S'04

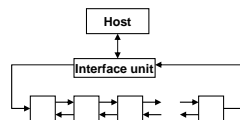
*T



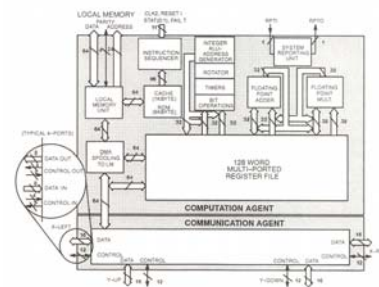
- 45 -

CS 418 S'04

iWARP



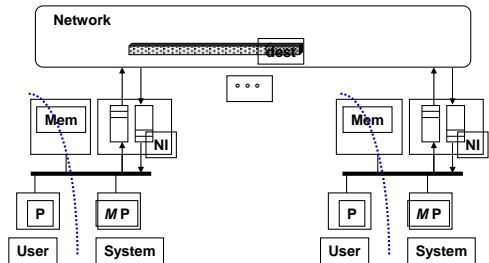
Nodes integrate communication with computation on systolic basis
 Msg data direct to register
 Stream into memory



- 46 -

CS 418 S'04

Dedicated Message Processing Without Specialized Hardware Design

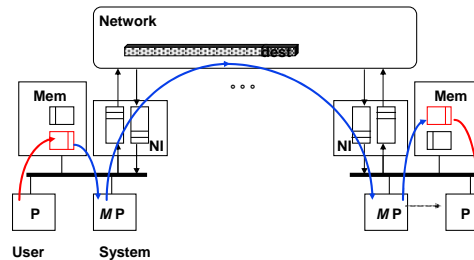


General Purpose processor performs arbitrary output processing (at system level)
 General Purpose processor interprets incoming network transactions (at system level)
 User Processor <-> Msg Processor via shared memory
 Msg Processor <-> Msg Processor via system network transaction

- 47 -

CS 418 S'04

Levels of Network Transaction

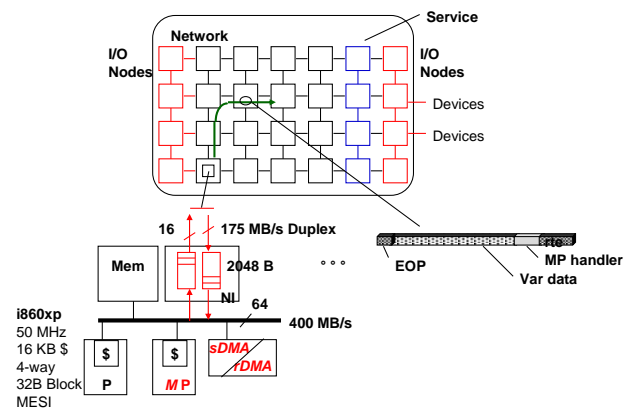


User Processor stores cmd / msg / data into shared output queue
 • must still check for output queue full (or make elastic)
 Communication assists make transaction happen
 • checking, translation, scheduling, transport, interpretation
 Effect observed on destination address space and/or events
 Protocol divided between two layers

- 48 -

CS 418 S'04

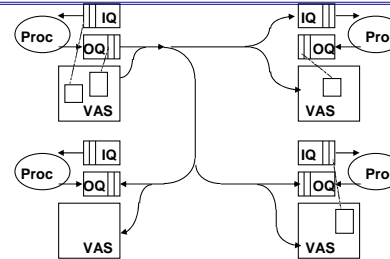
Example: Intel Paragon



- 49 -

CS 418 S'04

User Level Abstraction



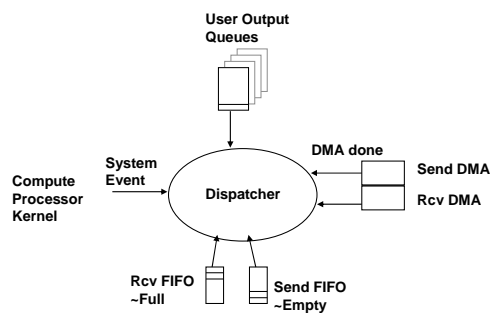
Any user process can post a transaction for any other in protection domain

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

- 50 -

CS 418 S'04

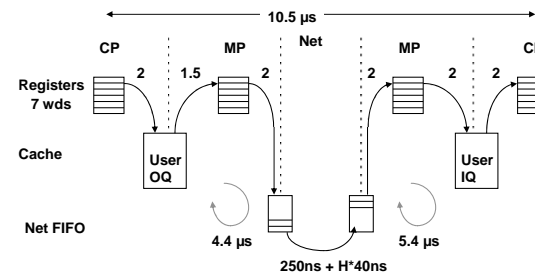
Msg Processor Events



- 51 -

CS 418 S'04

Basic Implementation Costs: Scalar



Cache-to-cache transfer (two 32B lines, quad word ops)

- producer: read(miss,S), chk, write(S,WT), write(I,WT), write(S,WT)
- consumer: read(miss,S), chk, read(H), read(miss,S), read(H), write(S,WT)

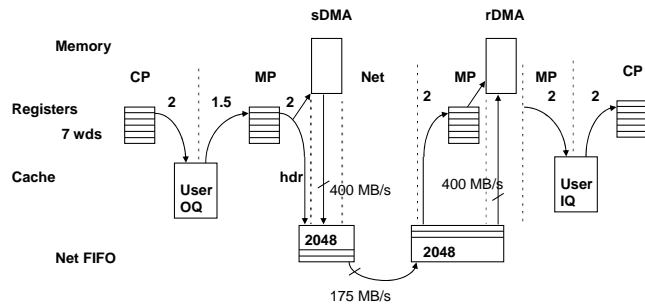
to NI FIFO: read status, chk, write, . . .

from NI FIFO: read status, chk, dispatch, read, read, . . .

- 52 -

CS 418 S'04

Virtual DMA -> Virtual DMA

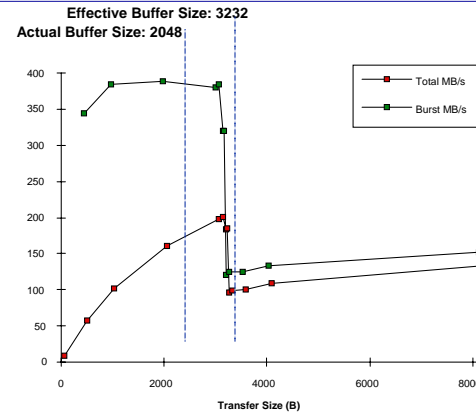


Send MP segments into 8K pages and does VA -> PA
 Rcv MP reassembles, does dispatch and VA -> PA per page

- 53 -

CS 418 S'04

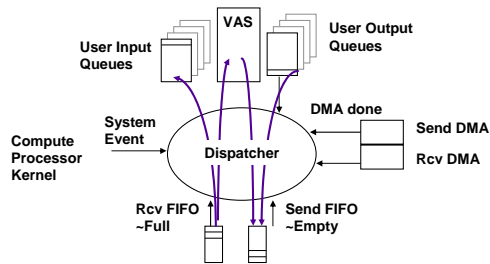
Single Page Transfer Rate



- 54 -

CS 418 S'04

Msg Processor Assessment



Concurrency Intensive

- Need to keep inbound flows moving while outbound flows stalled
- Large transfers segmented

Reduces overhead but adds latency

- 55 -

CS 418 S'04