

## Node-to-Network Interface in Scalable Multiprocessors

CS 258, Spring 99  
David E. Culler  
Computer Science Division  
U.C. Berkeley

## Racap: 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

3/10/99

CS258 S99

2

## Racap: Challenges (cont)

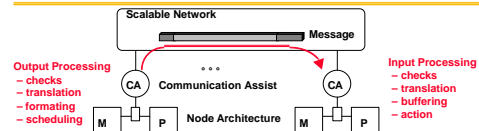
- **Fetch Deadlock**
  - For network to remain deadlock free, nodes must continue accepting messages, even when cannot source msgs
  - what if incoming transaction is a request?
    - » Each may generate a response, which cannot be sent!
    - » What happens when internal buffering is full?
- **logically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

3/10/99

CS258 S99

3

## Network Transaction Processing



- **Key Design Issue:**
- How much interpretation of the message?
- How much dedicated processing in the Comm. Assist?

3/10/99

CS258 S99

4

## Spectrum of Designs

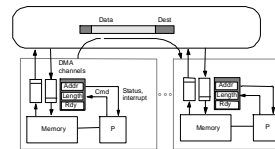
- **None: Physical bit stream**
    - blind, physical DMA nCUBE, iPSC, ...
  - **User/System**
    - User-level port CM-5, \*T
    - User-level handler J-Machine, Monsoon, .
    - ..
  - **Remote virtual address**
    - Processing, translation Paragon, Meiko CS-2
  - **Global physical address**
    - Proc + Memory controller RP3, BBN, T3D
  - **Cache-to-cache**
    - Cache controller Dash, KSR, Flash
- Increasing HW Support, Specialization, Intrusiveness, Performance (???)

3/10/99

CS258 S99

5

## Net Transactions: Physical DMA



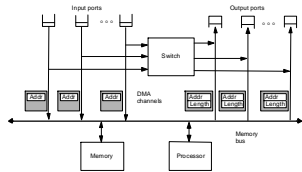
- **DMA controlled by regs, generates interrupts**
- **Physical => OS initiates transfers**
- **Send-side**
  - construct system "envelope" around user data in kernel area
- **Receive**
  - must receive into system buffer, since no interpretation inCA

3/10/99

CS258 S99

6

## 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

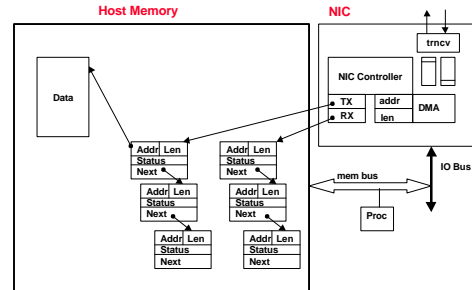
Os	16 ins	260 cy	13 us
Or	18	200 cy	15 us
			- includes interrupt

3/10/99

CS258 S99

7

## Conventional LAN NI

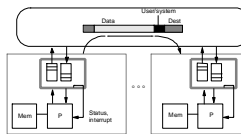


3/10/99

CS258 S99

8

## User Level Ports



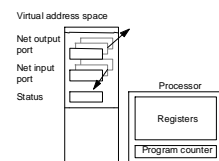
- initiate transaction at user level
- deliver to user without OS intervention
- network port in user space
- User/system flag in envelope
  - protection check, translation, routing, media access in src CA
  - user/sys check in dest CA, interrupt on system

3/10/99

CS258 S99

9

## User Level Network ports



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

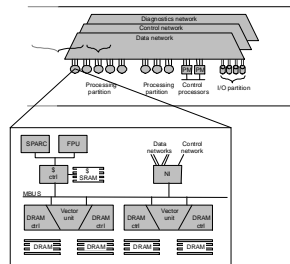
3/10/99

CS258 S99

10

## Example: CM-5

- Input and output FIFO for each network
- 2 data networks
  - index NI mapping table
- tag per message
- context switching?
- \*T integrated NI on chip
- iWARP also



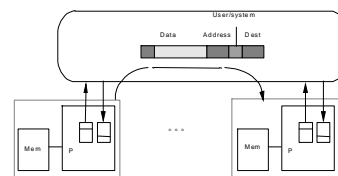
Os	50 cy	1.5 us
Or	53 cy	1.6 us
		interrupt 10us

3/10/99

CS258 S99

11

## User Level Handlers



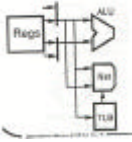
- Hardware support to vector to address specified in message
  - message ports in registers

3/10/99

CS258 S99

12

## J-Machine: Msg-Driven Processor



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

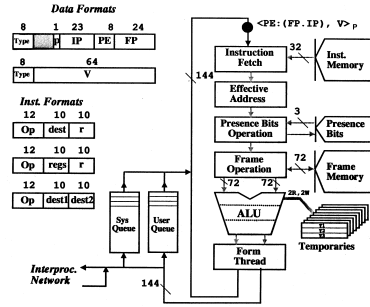


3/10/99

CS258 S99

13

## Monsoon Explicit Token-Store

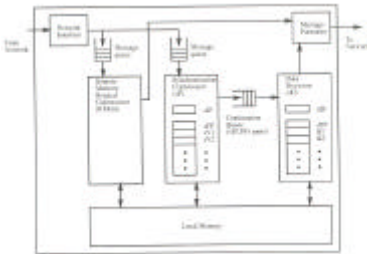


3/10/99

CS258 S99

14

## \*T: Network Co-Processor

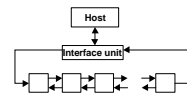


3/10/99

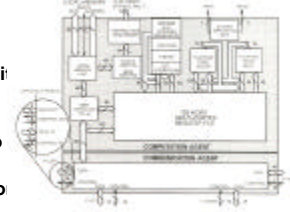
CS258 S99

15

## iWARP: Systolic Computation



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



3/10/99

CS258 S99

16

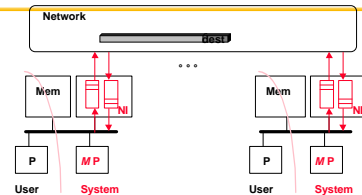
## Dedicated processing without dedicated hardware design

3/10/99

CS258 S99

17

## Dedicated Message Processor



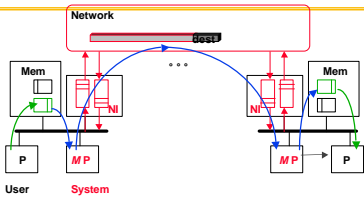
- General Purpose processor performs arbitrary output processing (at system level)
- General Purpose processor interprets incoming network transactions (at system level)
- User Processor <-> Msg Processor share memory
- Msg Processor <-> Msg Processor via system network transaction

3/10/99

CS258 S99

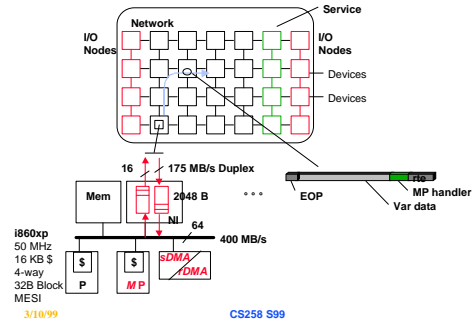
18

## Levels of Network Transaction



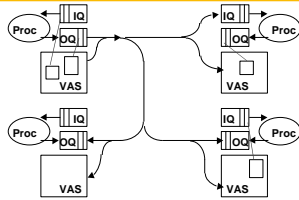
- User Processor stores cmd / msg / data into shared output queue
    - must still check for output queue full (or make elastic)
  - Communication assists make transaction happen
    - checking, translation, scheduling, transport, interpretation
  - Effect observed on destination address space and/or events
- 3/10/99 CS258 S99 19

## Example: Intel Paragon



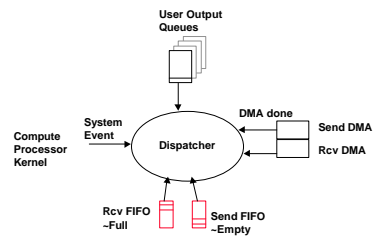
3/10/99 CS258 S99 20

## User Level Abstraction (Lok Liu)



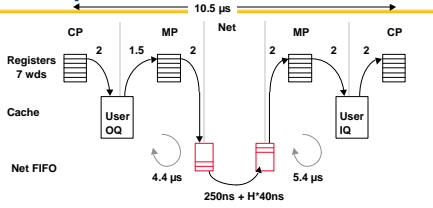
- Any user process can post a transaction for any other in protection domain
    - communication layer moves  $OQ_{src} \rightarrow IQ_{dest}$
    - may involve indirection:  $VAS_{src} \rightarrow VAS_{dest}$
- 3/10/99 CS258 S99 21

## Msg Processor Events



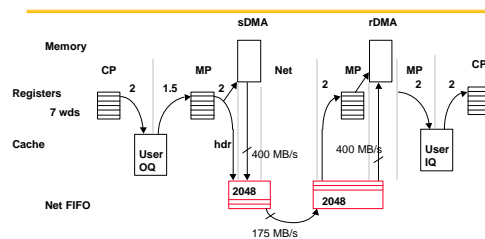
3/10/99 CS258 S99 22

## Basic Implementation Costs: Scalar



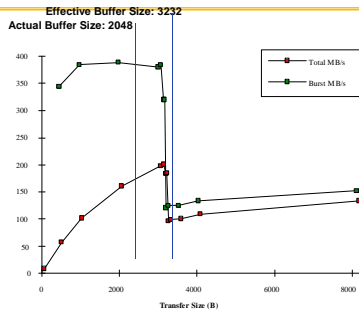
- Cache-to-cache transfer (two 32B lines, quad word ops)
    - producer: read(miss,S), chk, write(S,WT), write(I,WT), write(S,WT)
    - consumer: read(miss,S), chk, read(H), read(miss,S), read(H), write(S,WT)
  - to NI FIFO: read status, chk, write, ...
  - from NI FIFO: read status, chk, dispatch, read, read, ...
- 3/10/99 CS258 S99 23

## Virtual DMA -> Virtual DMA



- Send MP segments into 8K pages and does VA -> PA
  - Rcv MP reassembles, does dispatch and VA -> PA per page
- 3/10/99 CS258 S99 24

## Single Page Transfer Rate

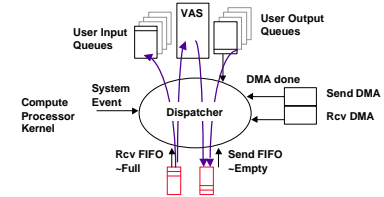


3/10/99

CS258 S99

25

## Msg Processor Assessment



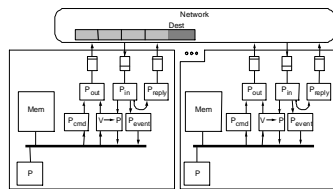
- **Concurrency Intensive**
  - Need to keep inbound flows moving while outbound flows stalled
  - Large transfers segmented
- **Reduces overhead but adds latency**

3/10/99

CS258 S99

26

## Case Study: Meiko CS2 Concept



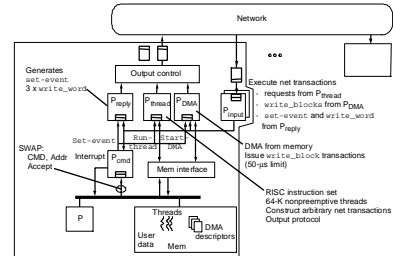
- **Circuit-switched Network Transaction**
  - source-dest circuit held open for request response
  - limited cmd set executed directly on NI
- **Dedicated communication processor for each step in flow**

3/10/99

CS258 S99

27

## Case Study: Meiko CS2 Organization

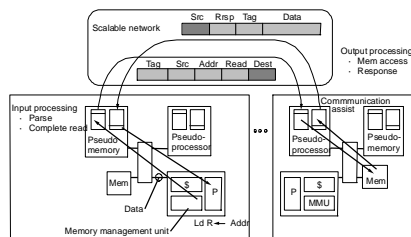


3/10/99

CS258 S99

28

## Shared Physical Address Space



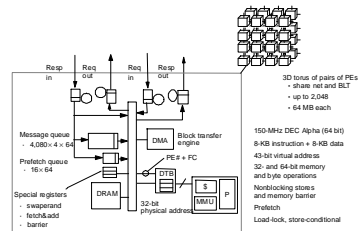
- **NI emulates memory controller at source**
- **NI emulates processor at dest**

3/10/99

CS258 S99

29

## Case Study: Cray T3D



- **Build up info in 'shell'**
- **Remote memory operations encoded in address**

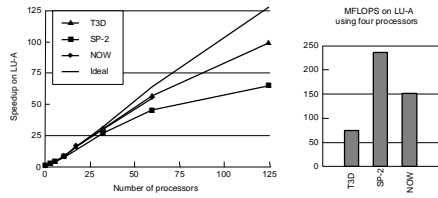
3/10/99

CS258 S99

30



## Application Performance on LU

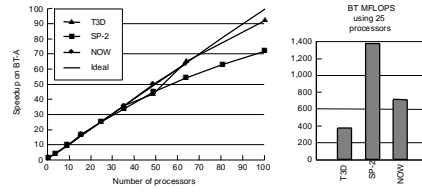


3/10/99

CS258 S99

37

## Application Performance on BT

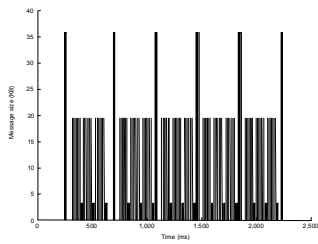


3/10/99

CS258 S99

38

## Message Profile on BT

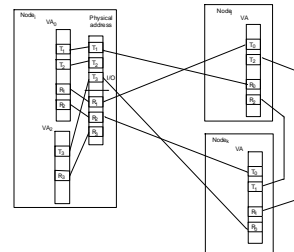


3/10/99

CS258 S99

39

## Reflective Memory



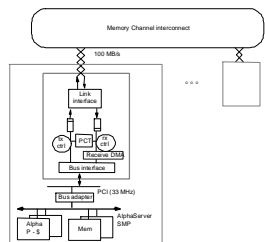
- Writes to local region reflected to remote

3/10/99

CS258 S99

40

## Case Study: DEC Memory Channel



- See also Shrimp

3/10/99

CS258 S99

41