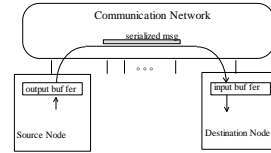


Realizing Programming Models

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

Network Transaction Primitive



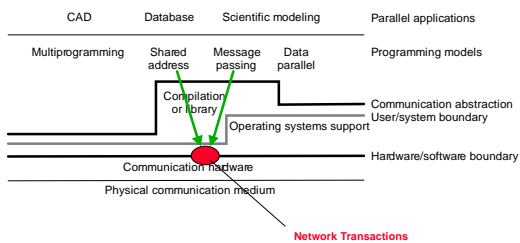
- one-way transfer of information from a source output buffer to a dest. input buffer
 - causes some action at the destination
 - occurrence is not directly visible at source
- deposit data, state change, reply

3/5/99

CS258 S99

2

Programming Models Realized by Protocols

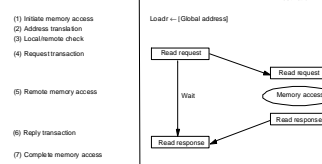


3/5/99

CS258 S99

3

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?

3/5/99

CS258 S99

4

The Fetch Deadlock Problem

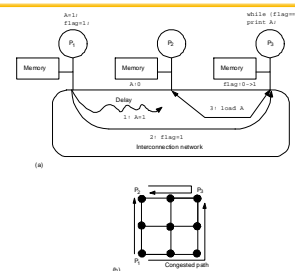
- Even if a node cannot issue a request, it must sink network transactions.
- Incoming transaction may be a request, which will generate a response.
- Closed system (finite buffering)

3/5/99

CS258 S99

5

Consistency



- write-atomicity violated without caching

3/5/99

CS258 S99

6

Key Properties of Shared Address 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 address space”
 - » 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

3/5/99

CS258 S99

7

Message passing

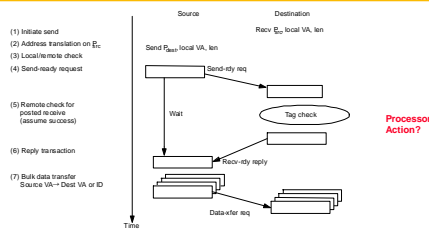
- Bulk transfers
- Complex synchronization semantics
 - more complex protocols
 - More complex action
- Synchronous
 - Send completes after matching rcv and source data sent
 - Receive completes after data transfer complete from matching send
- Asynchronous
 - Send completes after send buffer may be reused

3/5/99

CS258 S99

8

Synchronous Message Passing



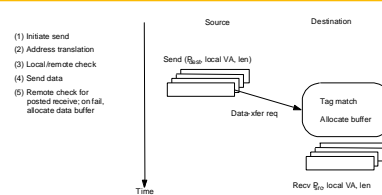
- Constrained programming model.
- Deterministic! What happens when threads added?
- Destination contention very limited.
- User/System boundary?

3/5/99

CS258 S99

9

Asynch. Message Passing: Optimistic



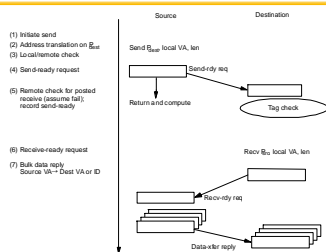
- More powerful programming model
- Wildcard receive => non-deterministic
- Storage required within msg layer?

3/5/99

CS258 S99

10

Asynch. Msg Passing: Conservative



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

3/5/99

CS258 S99

11

Key Features of Msg Passing Abstraction

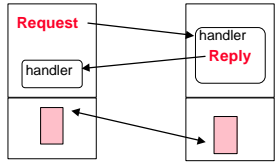
- Source knows send data address, dest. knows receive data address
 - after handshake they both know both
- Arbitrary storage “outside the local address spaces”
 - may post many sends before any receives
 - non-blocking asynchronous sends reduces the requirement to an arbitrary number of descriptors
 - » 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

3/5/99

CS258 S99

12

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**

3/5/99 CS258 S99 13

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

3/5/99 CS258 S99 14

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?

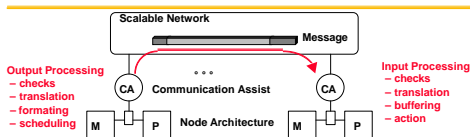
3/5/99 CS258 S99 15

Challenges in Realizing Prog. 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**

3/5/99 CS258 S99 16

Network Transaction Processing



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

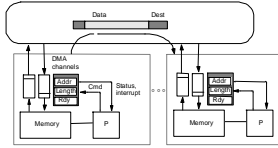
3/5/99 CS258 S99 17

Spectrum of Designs

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

3/5/99 CS258 S99 18

Net Transactions: Physical DMA

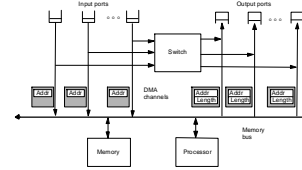


- DMA controlled by regs, generates interrupts
- Physical => OS initiates transfers
- Send-side
 - construct system "envelope" around user data in kernel area
- Receive
 - must receive into system buffer, since no interpretation in CA



3/5/99 CS258 S99 19

nCUBE Network Interface

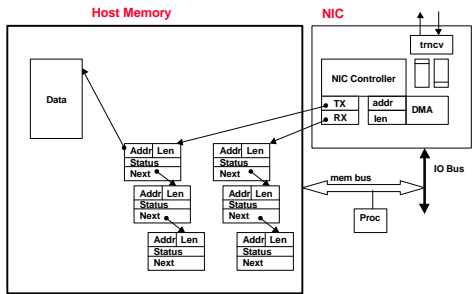


- independent DMA channel per link direction
 - leave input buffers always open
 - segmented messages
- routing interprets envelope
 - dimension-order routing on hypercube
 - bit-serial with 36 bit cut-through

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

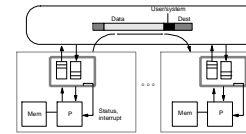
3/5/99 CS258 S99 20

Conventional LAN NI



3/5/99 CS258 S99 21

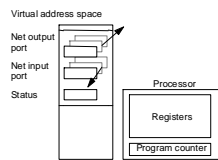
User Level Ports



- 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

3/5/99 CS258 S99 22

User Level Network ports

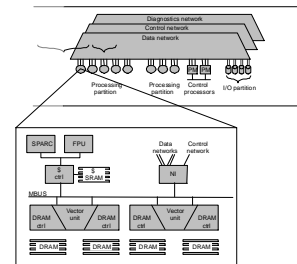


- Appears to user as logical message queues plus status
- What happens if no user pop?

3/5/99 CS258 S99 23

Example: CM-5

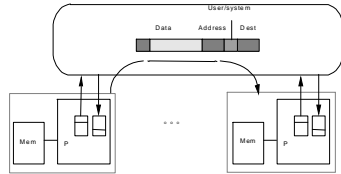
- Input and output FIFO for each network
- 2 data networks
- tag per message
 - index NI mapping table
- context switching?
- *T integrated NI on chip
- iWARP also



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

3/5/99 CS258 S99 24

User Level Handlers



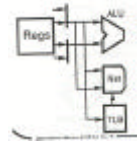
- Hardware support to vector to address specified in message
 - message ports in registers

3/5/99

CS258 S99

25

J-Machine



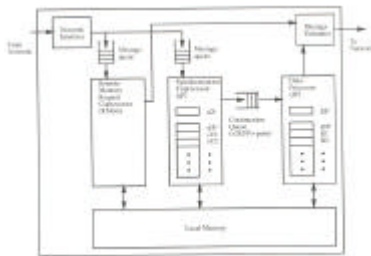
- Each node a small mdg driven processor
- HW support to queue msgs and dispatch to msg handler task

3/5/99

CS258 S99

26

*T

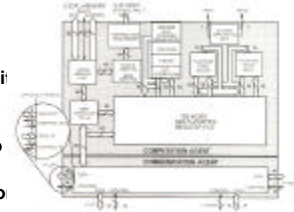
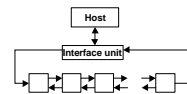


3/5/99

CS258 S99

27

iWARP



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

3/5/99

CS258 S99

28

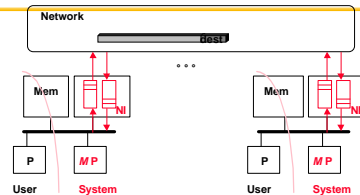
Dedicated processing without dedicated hardware design

3/5/99

CS258 S99

29

Dedicated Message Processor



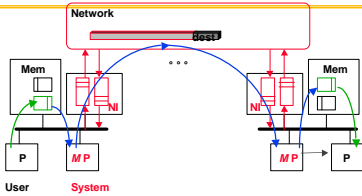
- General Purpose processor performs arbitrary output processing (at system level)
- General Purpose processor interprets incoming network transactions (at system level)
- User Processor <-> Msg Processor share memory
- Msg Processor <-> Msg Processor via system network transaction

3/5/99

CS258 S99

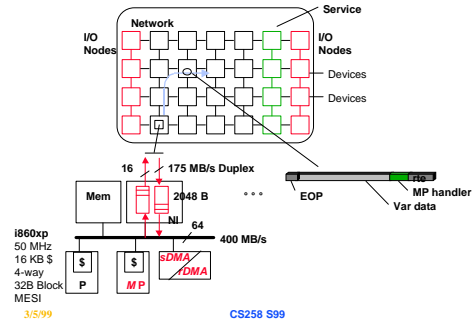
30

Levels of Network Transaction



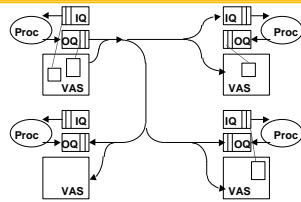
- User Processor stores cmd / msg / data into shared output queue
 - must still check for output queue full (or make elastic)
 - Communication assists make transaction happen
 - checking, translation, scheduling, transport, interpretation
 - Effect observed on destination address space and/or events
- 3/5/99 CS258 S99 31

Example: Intel Paragon



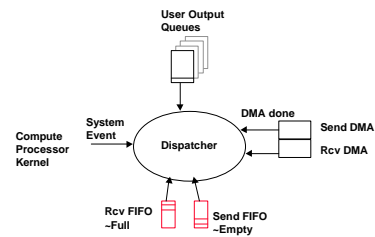
3/5/99 CS258 S99 32

User Level Abstraction (Lok Liu)



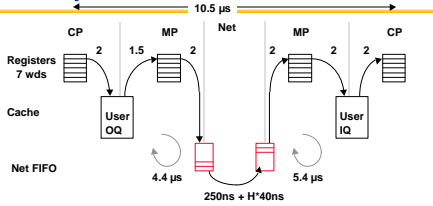
- Any user process can post a transaction for any other in protection domain
 - communication layer moves $OQ_{src} \rightarrow IQ_{dest}$
 - may involve indirection: $VAS_{src} \rightarrow VAS_{dest}$
- 3/5/99 CS258 S99 33

Msg Processor Events



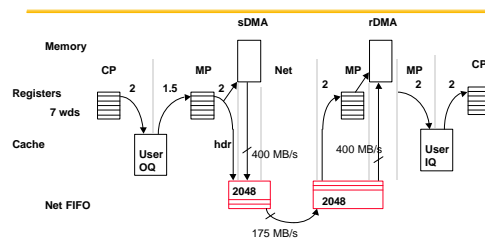
3/5/99 CS258 S99 34

Basic Implementation Costs: Scalar



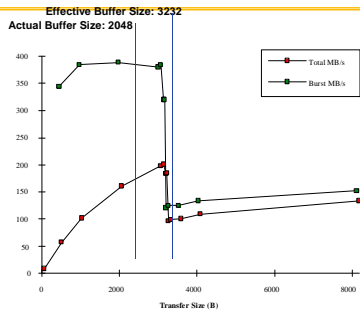
- Cache-to-cache transfer (two 32B lines, quad word ops)
 - producer: read(miss,S), chk, write(S,WT), write(I,WT), write(S,WT)
 - consumer: read(miss,S), chk, read(H), read(miss,S), read(H), write(S,WT)
 - to NI FIFO: read status, chk, write, ...
 - from NI FIFO: read status, chk, dispatch, read, read, ...
- 3/5/99 CS258 S99 35

Virtual DMA -> Virtual DMA



- Send MP segments into 8K pages and does VA -> PA
 - Rcv MP reassembles, does dispatch and VA -> PA per page
- 3/5/99 CS258 S99 36

Single Page Transfer Rate

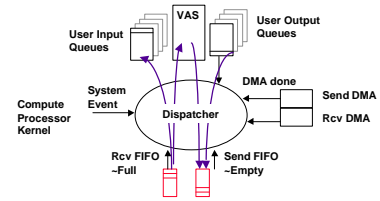


3/5/99

CS258 S99

37

Msg Processor Assessment



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

3/5/99

CS258 S99

38