

AN O/S PERSPECTIVE ON NETWORKS

Adem Efe Gencer¹

October 4th, 2012

¹Department of Computer Science, Cornell University

Papers

2

- **Active Messages: A Mechanism for Integrated Communication and Control**, Thorsten von Eicken, David E. Culler, Seth Copen Goldstein, and Klaus Erik Schauer. In Proceedings of the 19th Annual International Symposium on Computer Architecture, 1992.
- **U-Net: A User-Level Network Interface for Parallel and Distributed Computing**, Thorsten von Eicken, Anindya Basu, Vineet Buch and Werner Vogels. 15th SOSP, December 1995.

Authors - I

- Thorsten von Eicken
 - ▣ Founder of RightScale
 - ▣ Ph.D, CS at University of California, Berkeley
 - ▣ Assistant Prof. at Cornell (1993-1999)
- David E. Culler
 - ▣ Chair, Electrical Engineering and Computer Sciences at Berkeley
- Seth Copen Goldstein
 - ▣ Associate Professor at Carnegie Mellon University
- Klaus Erik Schauser
 - ▣ Associate professor at UCSB

Authors - II

4

- Thorsten von Eicken
- Anindya Basu
 - ▣ Advised by von Eicken
- Vineet Buch
 - ▣ Google
 - ▣ MS from Cornell
- Werner Vogels
 - ▣ Research Scientist at Cornell (1994-2004)
 - ▣ CTO and vice President of Amazon

Active Messages

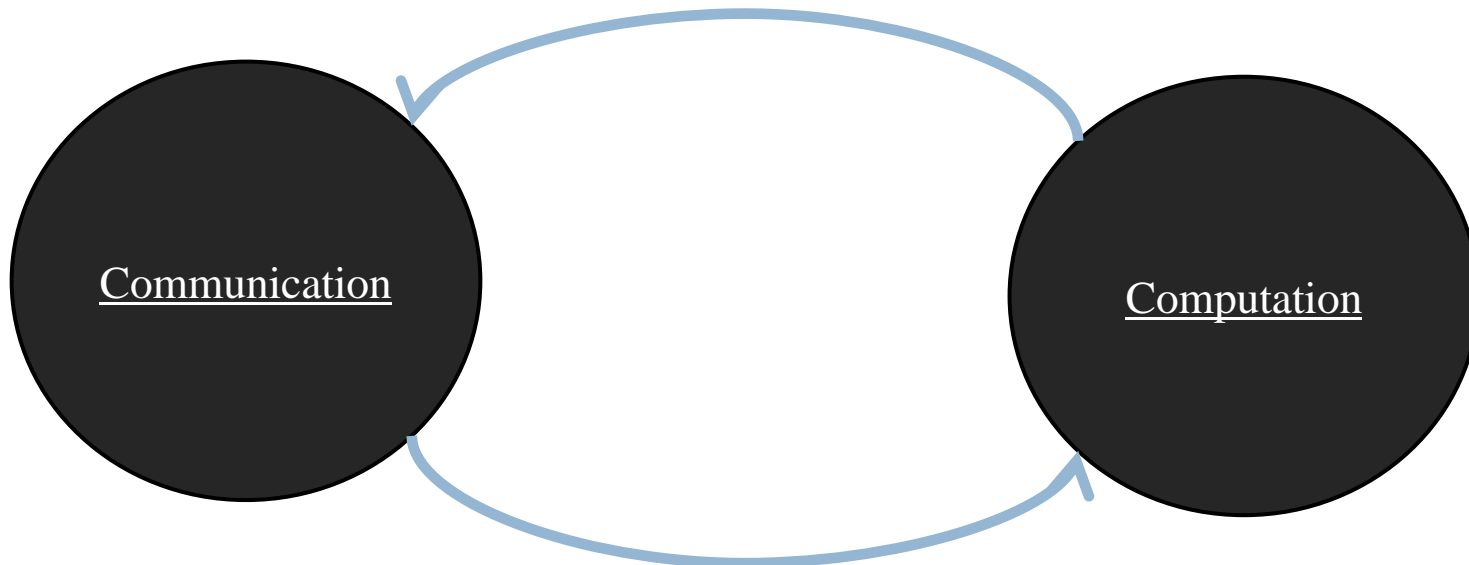
5

- Introduction
- Preview: Active Message
 - ▣ Traditional programming models
- Active Messages
- Environment
- Split-C
- Message Driven Model vs. Active Messages
- Final Notes

Introduction

6

- Inefficient use of underlying hardware
 - ▣ Poor overlap between communication and computation
 - ▣ High communication overhead



Introduction: Objective

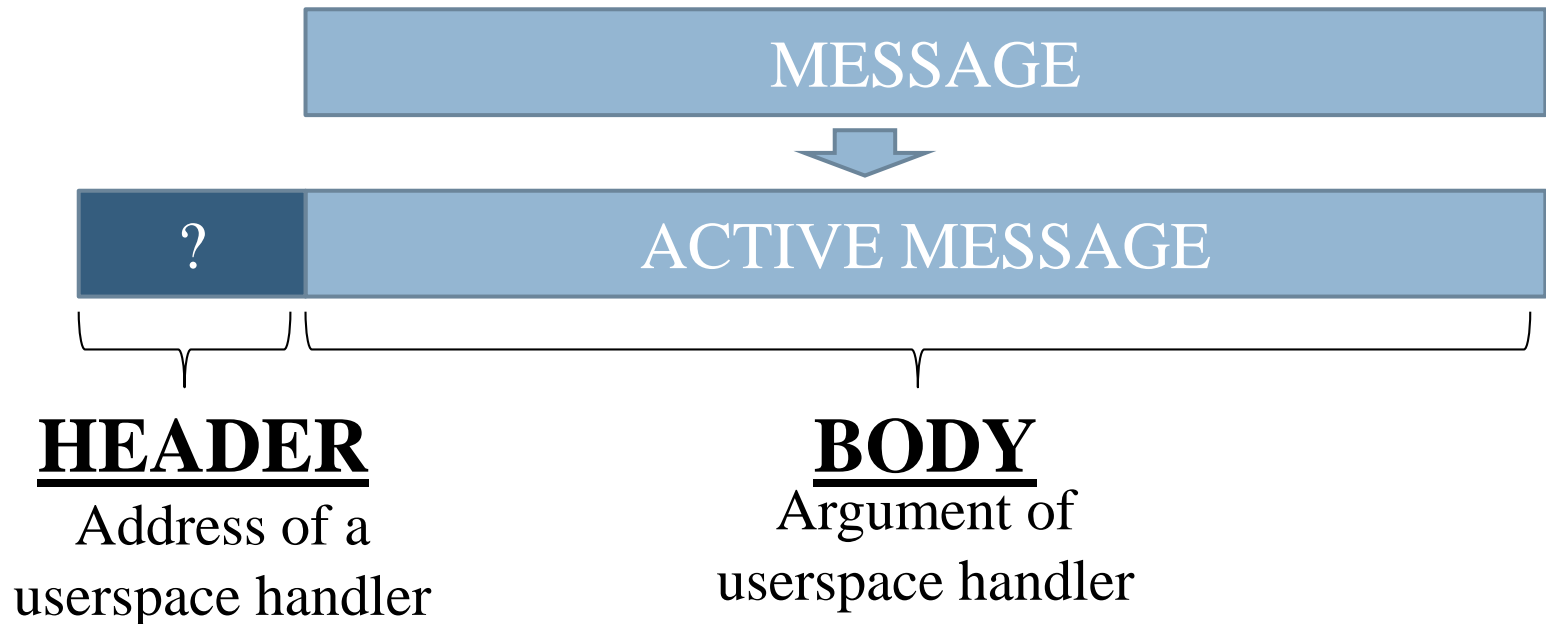
7

- Eliminate inefficiency by increasing the overlap between communication & computation!

Preview: Active Message

8

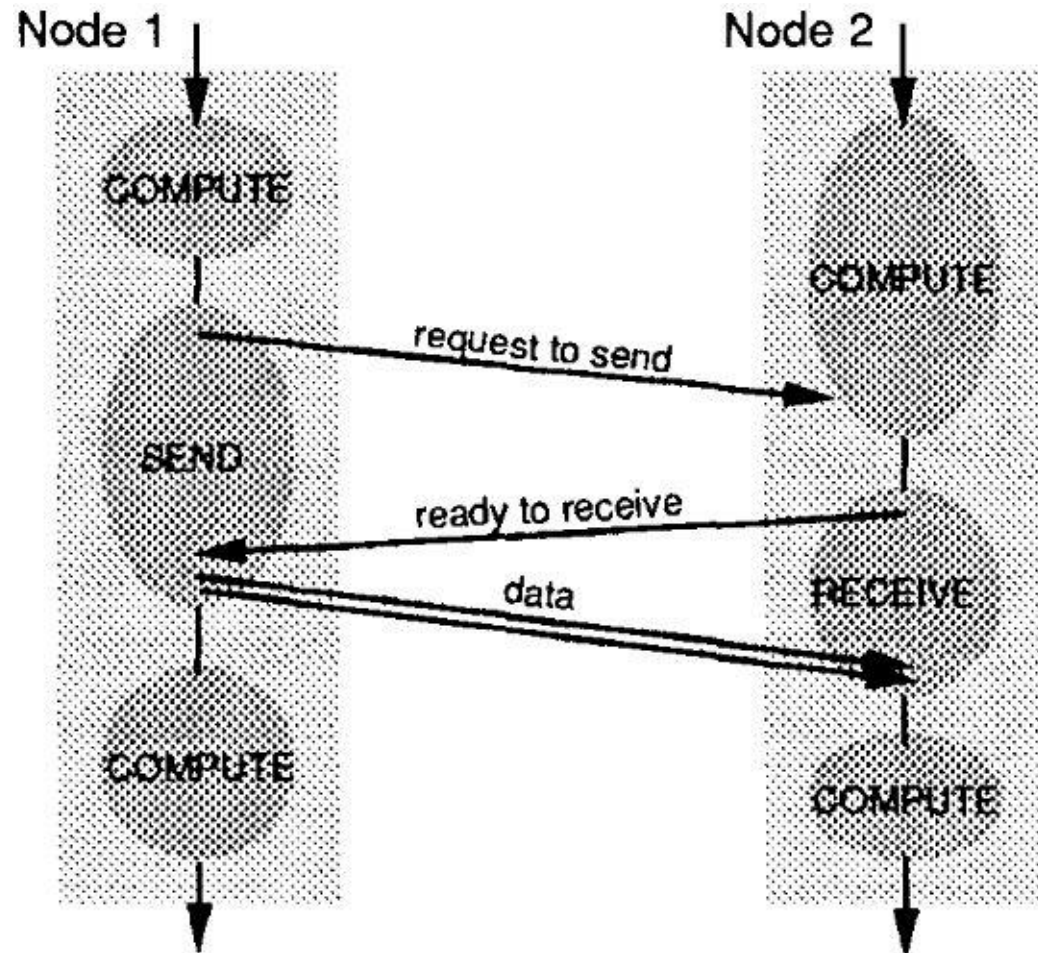
- Solution: A novel asynchronous communication mechanism! **Active messages**



Traditional Programming Models

9

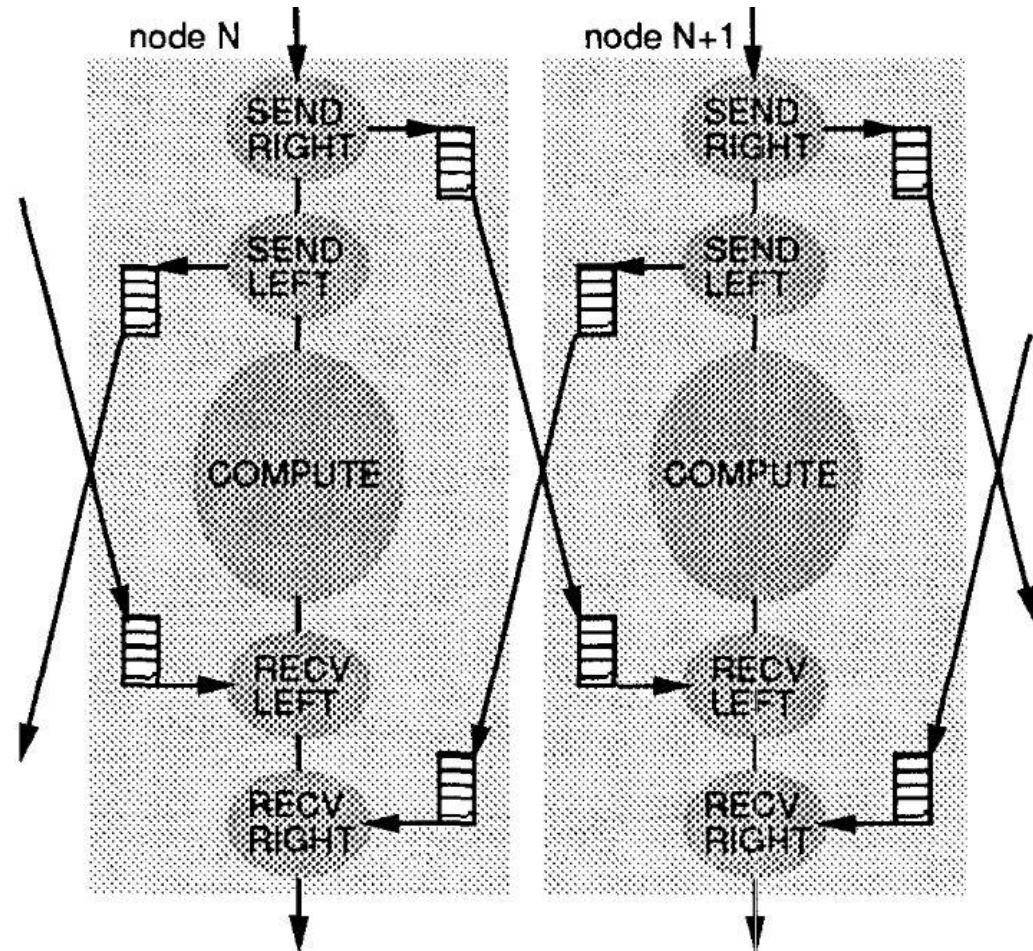
- ❑ Synchronous comm.
 - ❑ 3-phase protocol
 - ❑ Send / receive are blocking
- 👍 Simple
- 👍 No buffering
- 👎 High comm. latency
- 👎 Underutilized network bandwidth



Traditional Programming Models

10

- Asynchronous comm.
 - ▣ Send / receive are non-blocking
 - ▣ Message layer buffers the message
- ☞ Communication & Computation can overlap



Traditional Programming Models

11

	Synchronous	Asynchronous	
PROS	No Buffering	Overlap	} <i>Active Messages</i>
CONS	No Overlap High Comm. Latency	Buffering	

Active Messages

12

- Requires SPMD programming model
- Key optimization: not buffered
 - ▣ Except for network transport buffering
 - ▣ For large messages buffering is needed on sender side
 - ▣ Receiving side pre-allocates structures
- Handlers do not block - deadlock avoidance

Active Messages

13

- Sender
 - ▣ Packs the message with a header containing handler's address
 - ▣ Pass the message to the network
 - ▣ Continue computing
- Receiver
 - ▣ Message arrival invokes handler
 - ▣ Handler Extracts the data & interrupts computation
 - ▣ Passes it to the ongoing computation

- Implementation of active messages on different platforms is not the same (nCUBE/2 vs. CM-5).

Environment

14

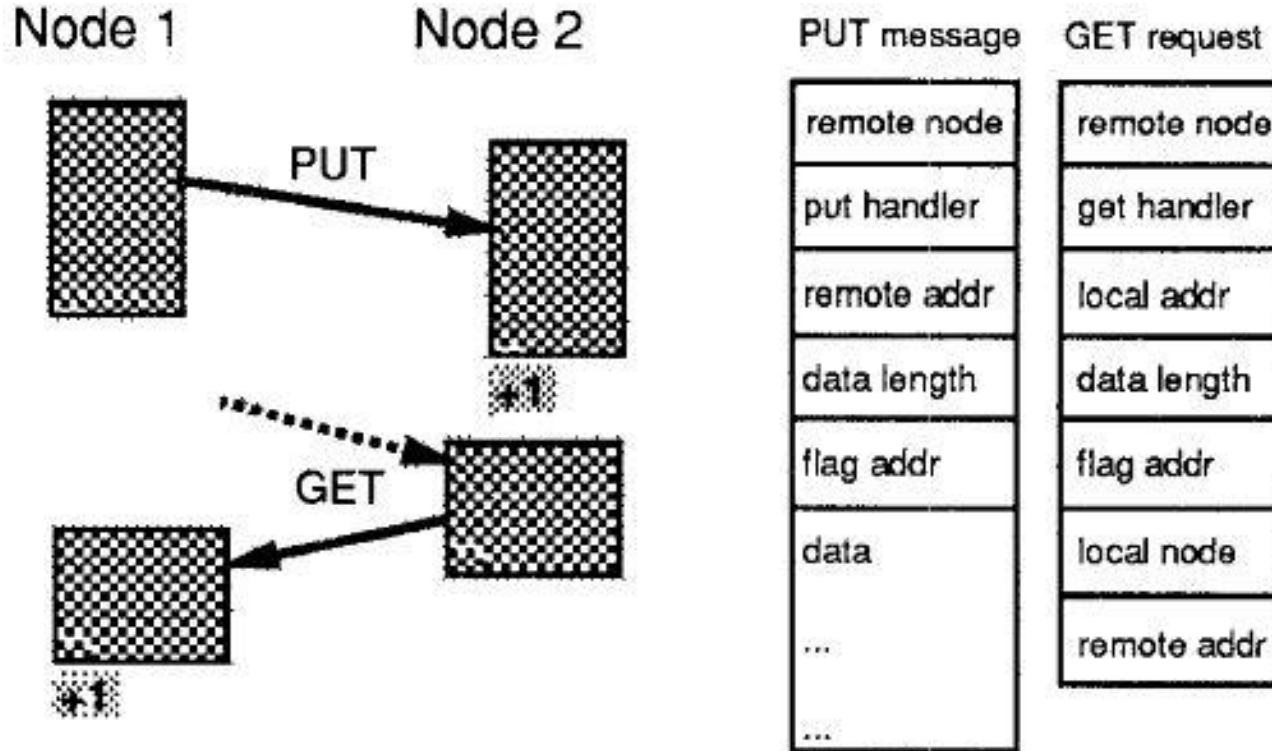
- Implementation on parallel supercomputers
 - *nCUBE/2*
 - Binary hypercube network
 - NI with 28 DMA channels
 - *CM-5 (connection machine-5)*
 - Hypertree network
- Split-C
 - A parallel extension of C



CM-5 machine of NSA

Split-C: PUT and GET

15



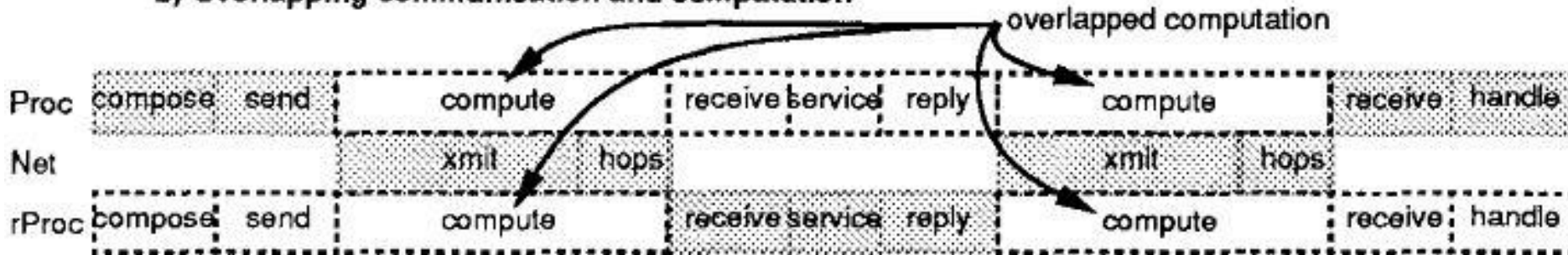
- Nonblocking
 - ▣ Message handler
 - ▣ Message formatter
- Used in matrix multiplication to show the performance of active message

Split-C: GET

16

- Matrix multiply achieves 95% performance in nCUBE/2 configuration.

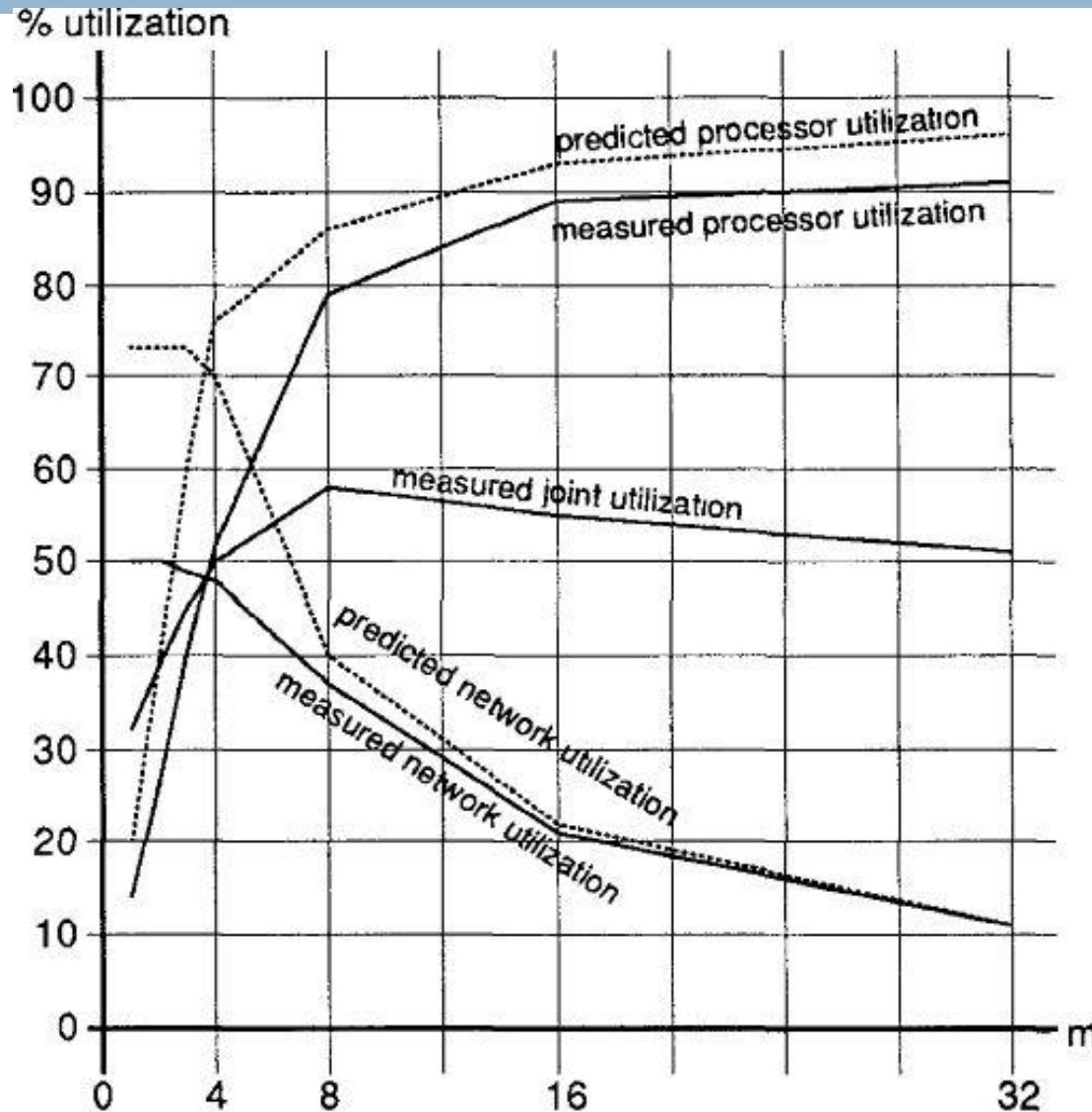
b) Overlapping communication and computation



- Xmit: time to inject the message into the network
- Hops: time for the network hops

Split-C: Matrix Multiply

17



- m : # of columns per processor
- # of nodes is constant ($N=128$)

Message Driven Model vs. Active Messages

18

- Message driven model
 - ▣ Computation in message handlers
 - ▣ Handler may suspend
 - ▣ Memory allocation & scheduling on message arrival
- Active Messages
 - ▣ Computation in background (handler extracts message)
 - ▣ Immediate execution

Final Notes

19

- ❑ Not a new parallel programming paradigm!
 - 👍 A primitive communication mechanism
 - 👍 Used to implement them efficiently
- 🔗 The need for userspace handler address leads to programming difficulty
- 🔗 Requires SPMD model (but not required on some modified versions)

U-Net

20

- Introduction
- U-Net
 - ▣ Remove the Kernel from Critical Path
- Communication in U-Net
- Protection
- Zero Copy
 - ▣ True zero copy vs. zero copy
- Environment
- Performance
- Final Notes

Introduction

21

- Low bandwidth and high communication latency
 - ▣ No longer hardware issues!
 - ▣ Problem is the message path through the *kernel*
- Small messages
 - ▣ Common in various applications
 - ▣ Requires low round-trip latency

Latency = processing overhead + network latency

Introduction: Objective

22

Make communication ***FASTER!***



U-Net

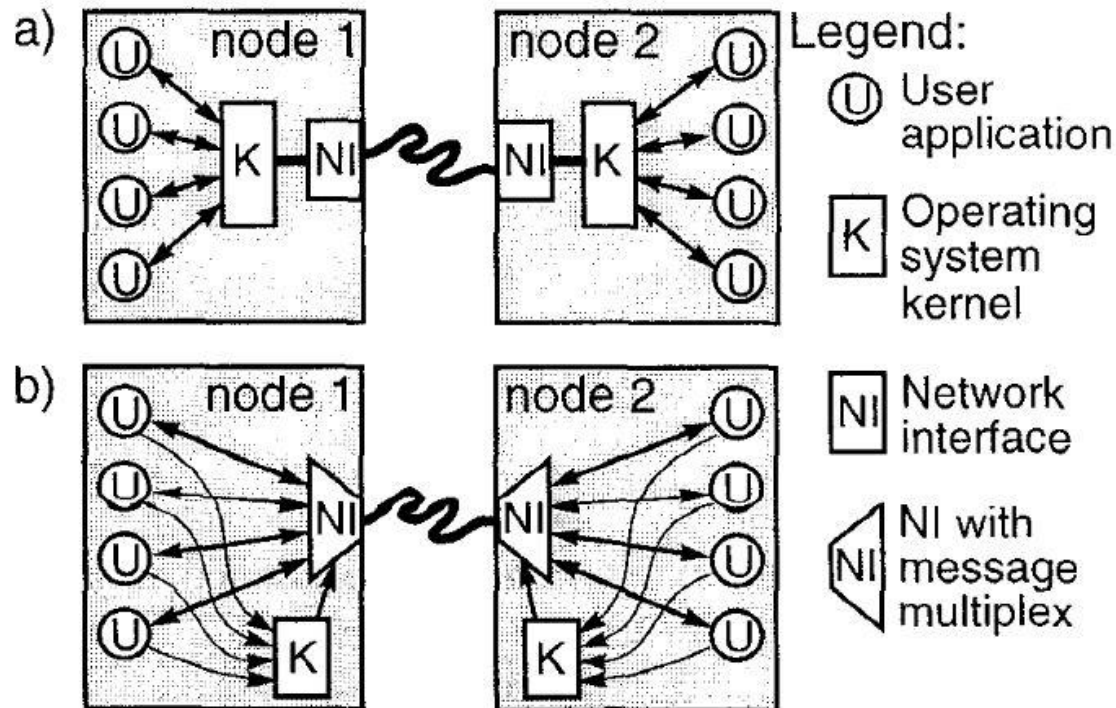
23

- ↳ Low communication latency
- ↳ High bandwidth
 - ↳ even with small messages!
- ↳ Support for workstations with off-the shelf network
- ↳ Keep providing full protection
 - ↳ Kernel controlled channel setup / tear-down
- ↳ TCP, UDP, Active Messages, ... can be implemented

Remove the kernel from critical path

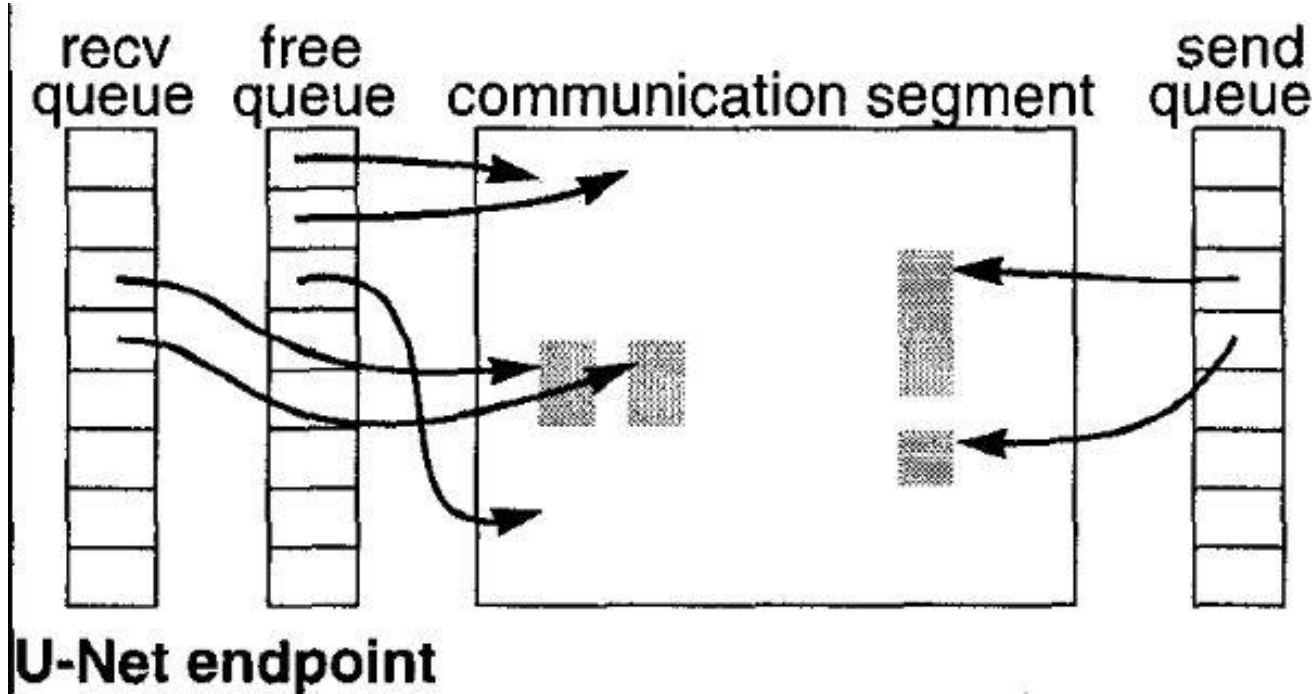
24

- Solution: Enable user-level access to NI to implement user-level communication protocols
 - ▣ Does not require OS modification or custom HW
 - ▣ Provides higher flexibility (app. specific protocols!)



Communication in U-Net

25

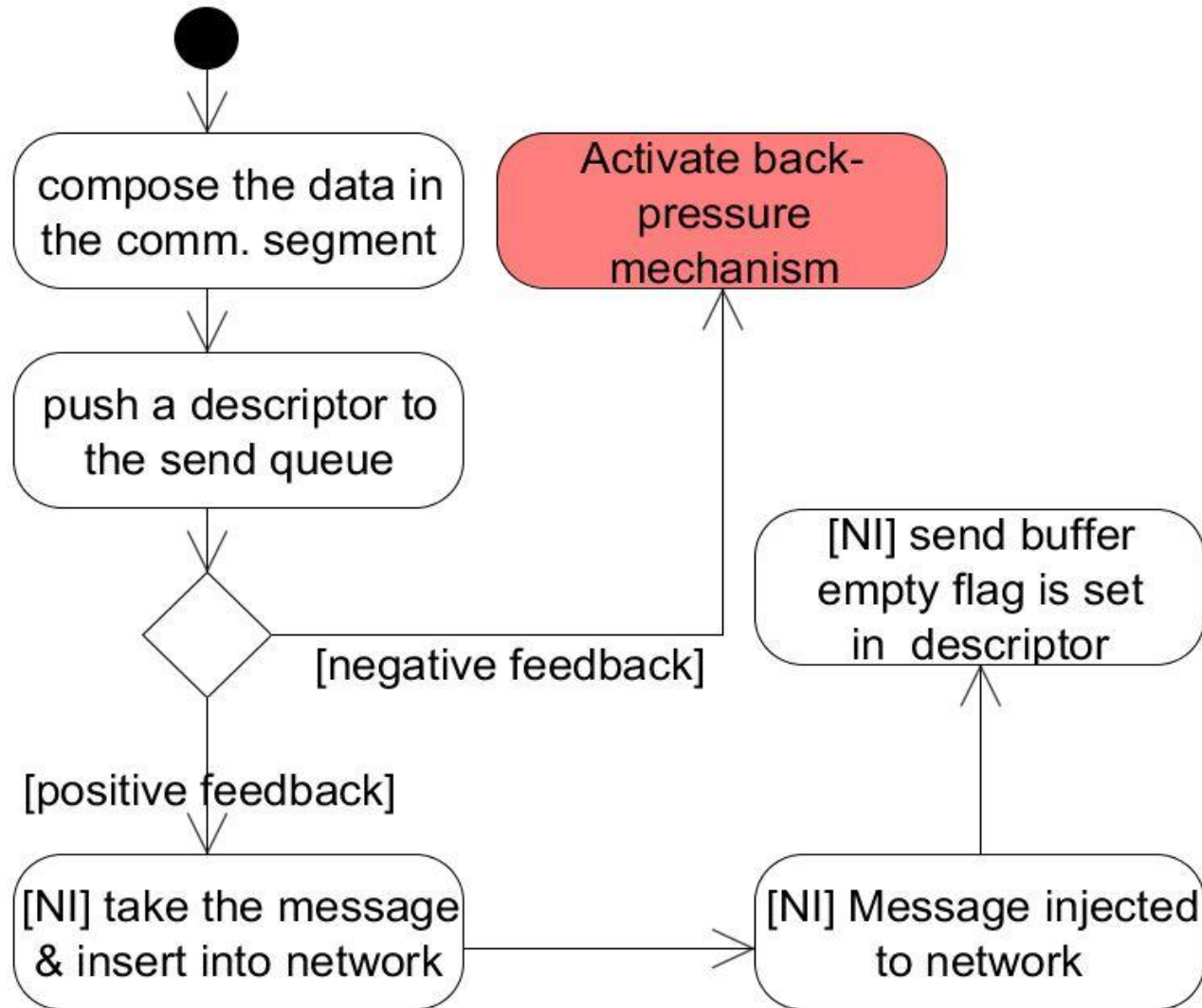


Each process creates endpoints to access network

- ❑ **U-net endpoint**: application handle into the network
- ❑ **Communication segment**: regions of memory that hold message data
- ❑ **Send/rcv/free queues**: hold message descriptors

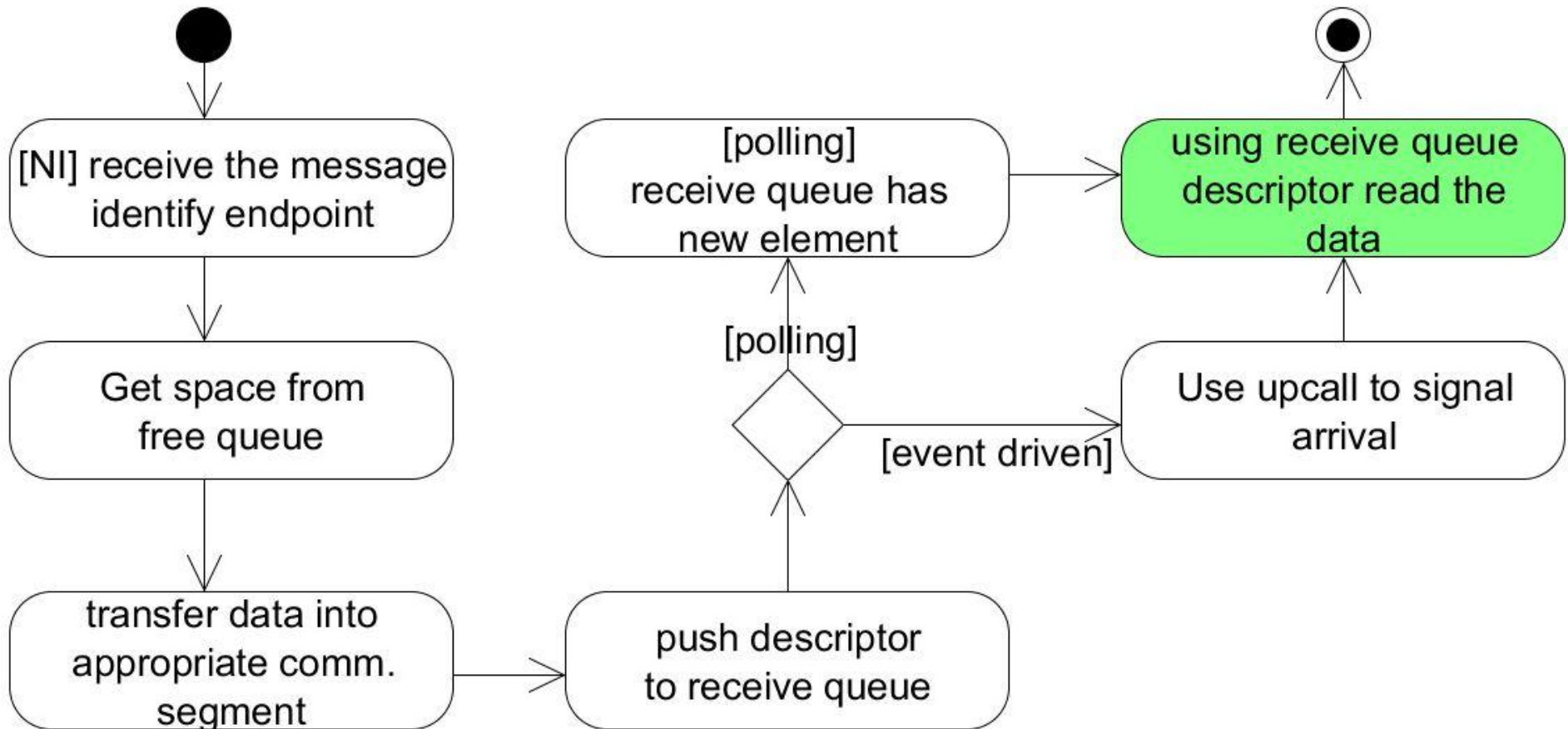
Communication in U-Net: Send

26



Communication in U-Net: Receive

27



Protection

- Owning process protection:

- Endpoints
- Communication Segments
- Send/Receive/Free Queues

Accessible by owning process only!

- Tag protection:

- Outgoing message
 - tagged with originating address
- Incoming message
 - delivered to correct destination point only!

Zero Copy

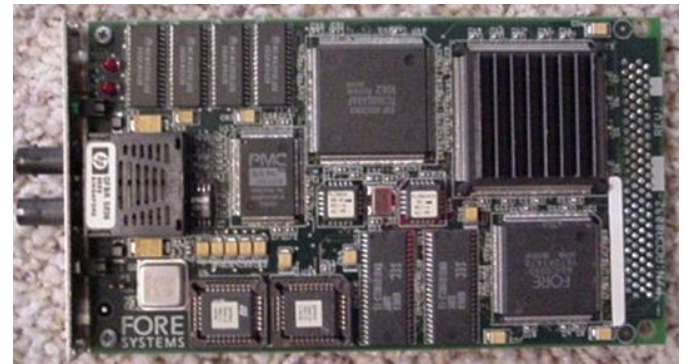
29

- Base level U-Net (zero copy)
 - ▣ Send/receive needs a buffer (not really “zero” copy!)
 - ▣ Requires a copy between app. data structures and a buffer in the comm. segment
- Direct Access U-Net (true zero copy)
 - ▣ Spans on entire address space
 - ▣ Has special hardware requirements

Environment

30

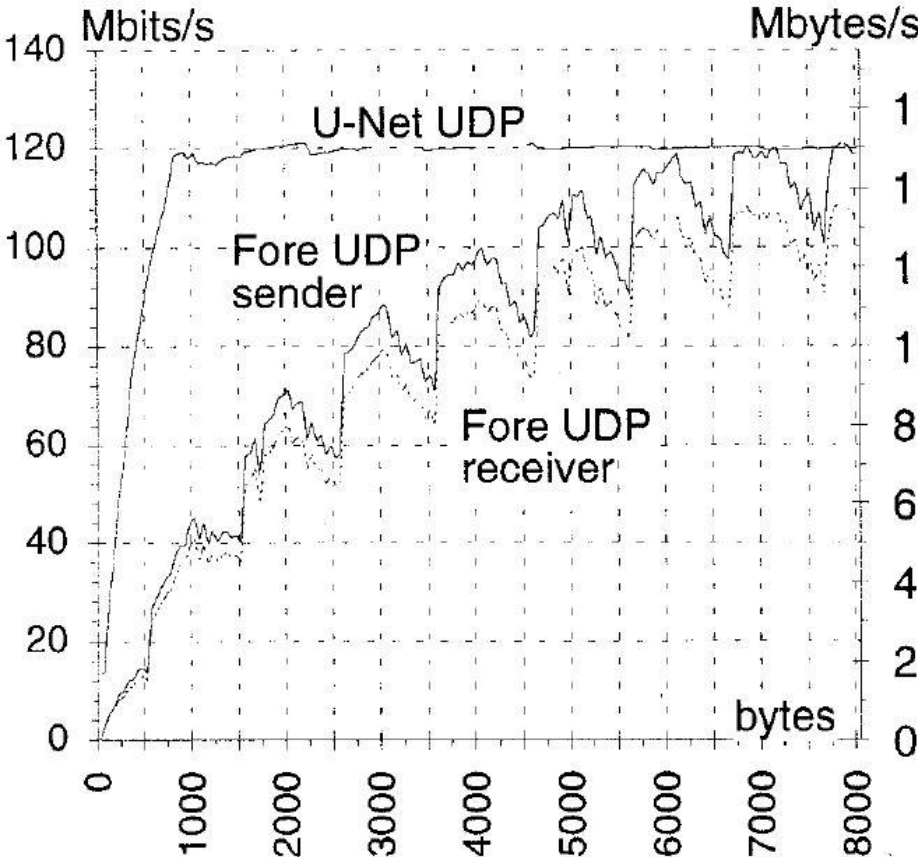
- Implementation on off-the-shelf hardware platform
 - *SPARCStations*
 - Fore SBA-100 NIC
 - Fore SBA-200 NIC
- SunOS 4.1.3
 - BSD-based Unix OS
 - Ancestor of Solaris (after version 5.0)
- Split-C
 - A parallel extension of C



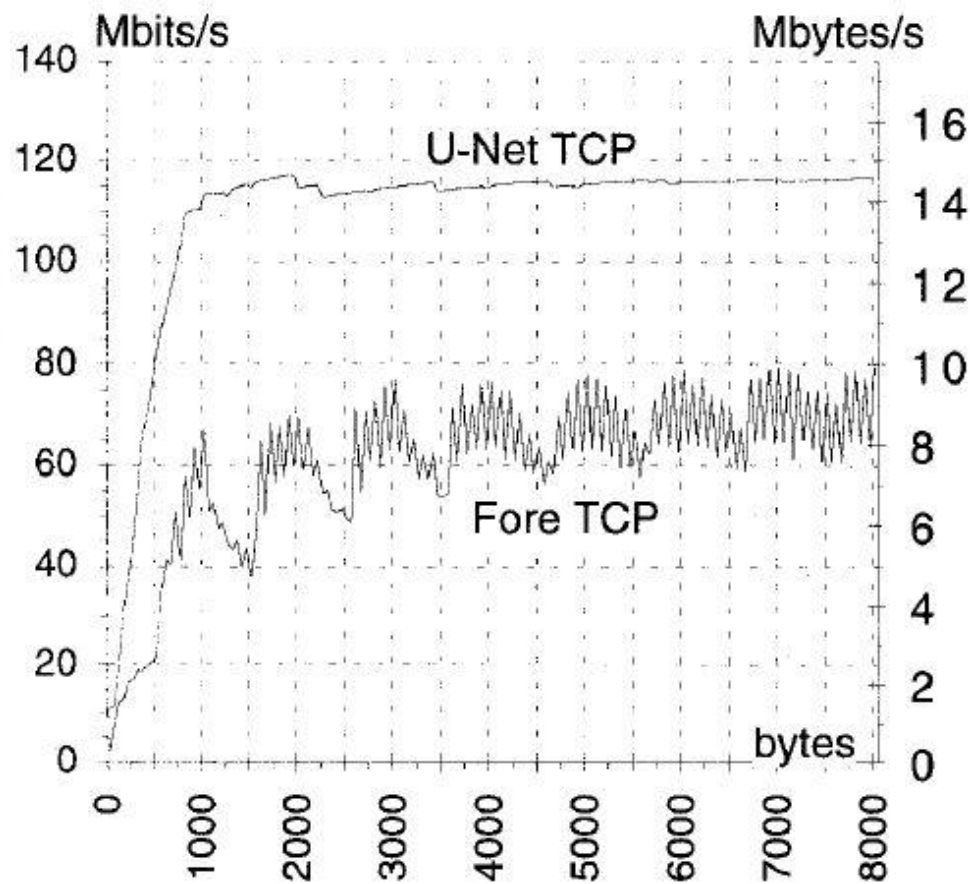
Fore SBA-200 NIC

Performance: UDP vs. TCP

31



UDP bandwidth as a function of message size

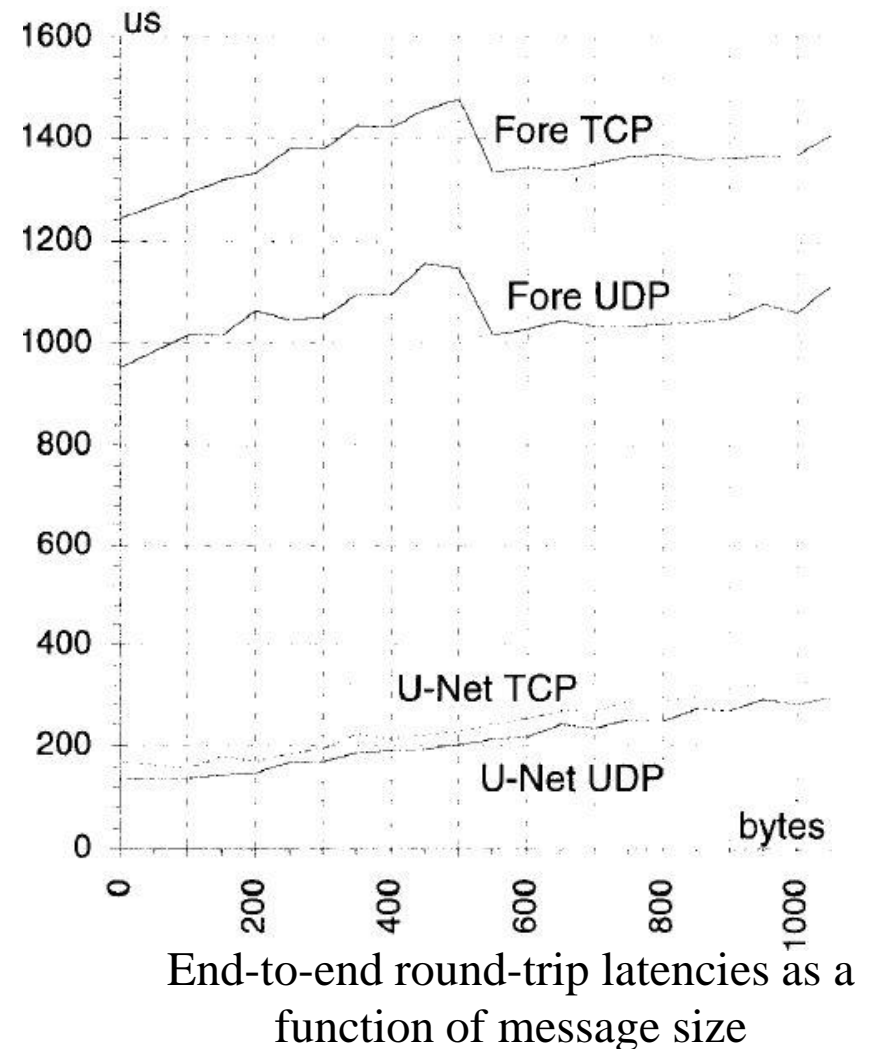


TCP bandwidth as a function of data generation by application

Performance: UDP vs. TCP

32

- Fast U-Net roundtrips:
 - ▣ Apprx. 7x speedup for TCP
 - ▣ Apprx. 5x speedup for UDP
- Fast U-Net TCP roundtrips let use of small window size.



Final Notes

33

- 👍 Low comm. latency
- 👍 High bandwidth
- 👍 Good flexibility

Protocol	Round-trip latency	Bandwidth 4K packets
Raw AAL5	65 μ s	120Mbits/s
Active Msgs	71 μ s	118Mbits/s
UDP	138 μ s	120Mbits/s
TCP	157 μ s	115Mbits/s
Split-C store	72 μ s	118Mbits/s

Table 3: U-Net latency and bandwidth Summary.

Cluster of workstations vs. supercomputers

- 👍 Comparable performance using U-Net
- 👎 Additional system resources (parallel process scheduler, file systems, ...) are needed