

file ↓

cc: Holloway,
Retberg,
Barker
bug: me

To: Steve Levy
Dan Sullivan
Dave Walden
Alex McKenzie
From: Steve Blumenthal

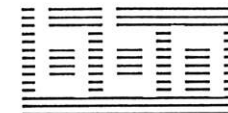
A copy of slides
we presented to
DARPA and NSA

MONARCH

BBN Systems and Technologies

Steve Blumenthal
Phil Carvey
Gregg Bromley

January 23, 1990



Monarch

Agenda

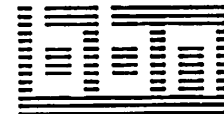
- **Monarch Introduction** Blumenthal
- **Monarch Hardware** Carvey
- **High Density Packaging**
- **Programming Model and Software** Bromley
- **Program Plan and Costs** Blumenthal



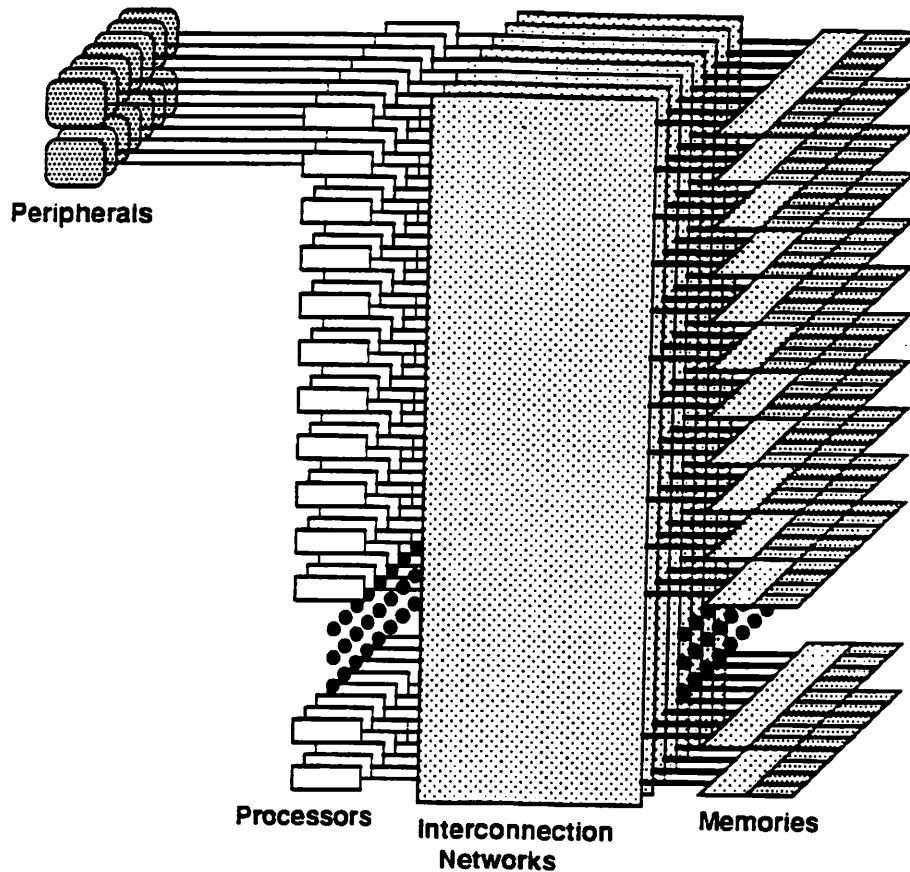
Monarch

Overview

- **Uniform Access Shared Memory Machine**
- **Single Program/Multiple Thread**
Fine Grain Parallel MIMD Programming Model
- **Scalable to Teraflop Machine**
- **Low Latency Interconnection Scheme**



Monarch System Design



The Switch is the Key

- Minimize Effective Latency
- Provides Very High Processor-Memory Bandwidth
- Economically Scalable System



Monarch

Processors

- **Commercial Processors can Meet Multiprocessor Requirements**
 - Cost/Performance
 - Lower Development Risk
 - Evolutionary Family Tracks Technology
 - Community of Uniprocessor Users and Software
- **with Latency Masking Techniques**
 - Parallel Instruction Execution & Memory Access
 - Operand Prefetch
 - Multiple Operand Fetch



Monarch

Scalability

<u>Procs</u>	<u>GIPS</u>	<u>GFLOPS</u>	<u>Mbytes</u>
65K	6,500	6,500	17,180
8K	800	800	2,150
1024	100	100	270
128	13	13	67

- **Cost Effective**

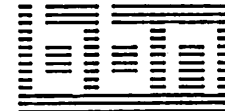
- Target raw parts cost Cost \$2000 per Processor
- Target raw parts cost \$20 per MIP
- Target raw parts cost \$20 per MFLOP



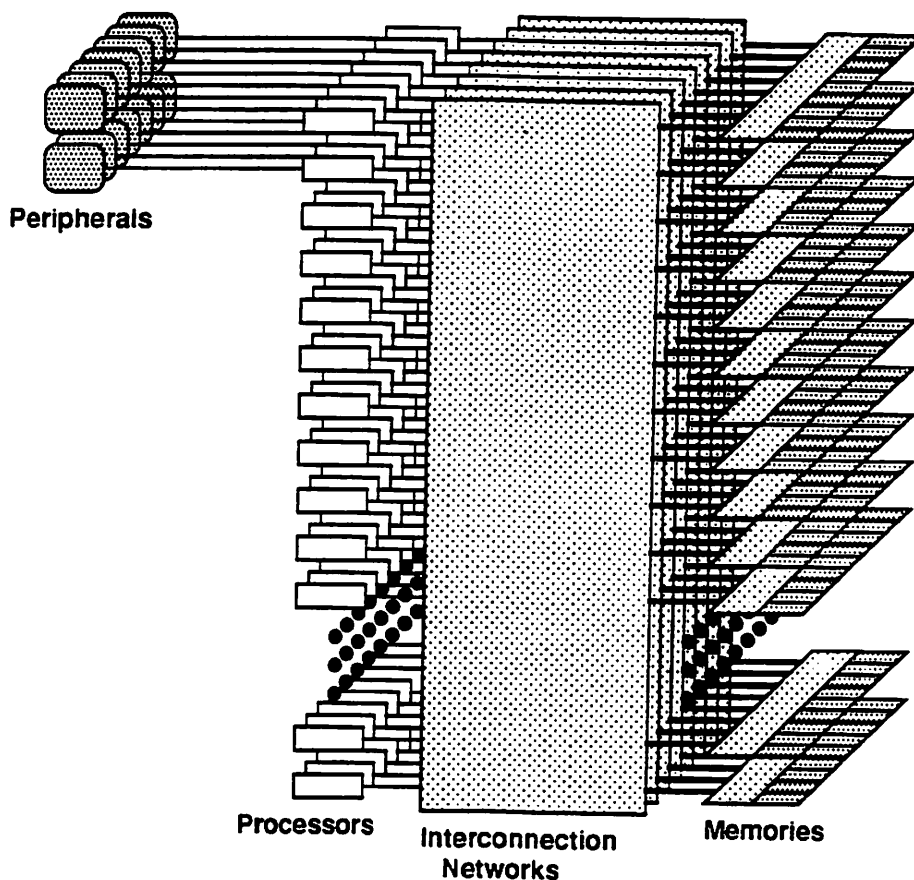
Monarch

Recent Developments

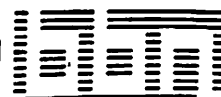
- **BBN ACI introduced the TC/2000 multiprocessor**
 - Software - Significant Development Investment
 - MACH Operating System
 - Programming Languages
 - Parallel Runtime & Synchronization Libraries
 - Debuggers and Tools
- **BBN STC has continued to investigate highly parallel VLSI multiprocessor systems**
 - Updated Switch Design
 - Use of Commercial Microprocessor
 - Latency Masking Research
 - New High Density Packaging Concepts
- **BBN STC IR & D study to adapt switch technology to multigigabit networks**



Monarch System Design



- * Memory physically partitioned into equal elements accessed via spatial switching
- * When memory elements have sufficient bandwidth, two or more processor - interconnect sets can be supported per memory elements via time multiplexing

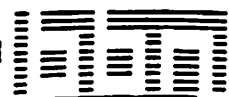


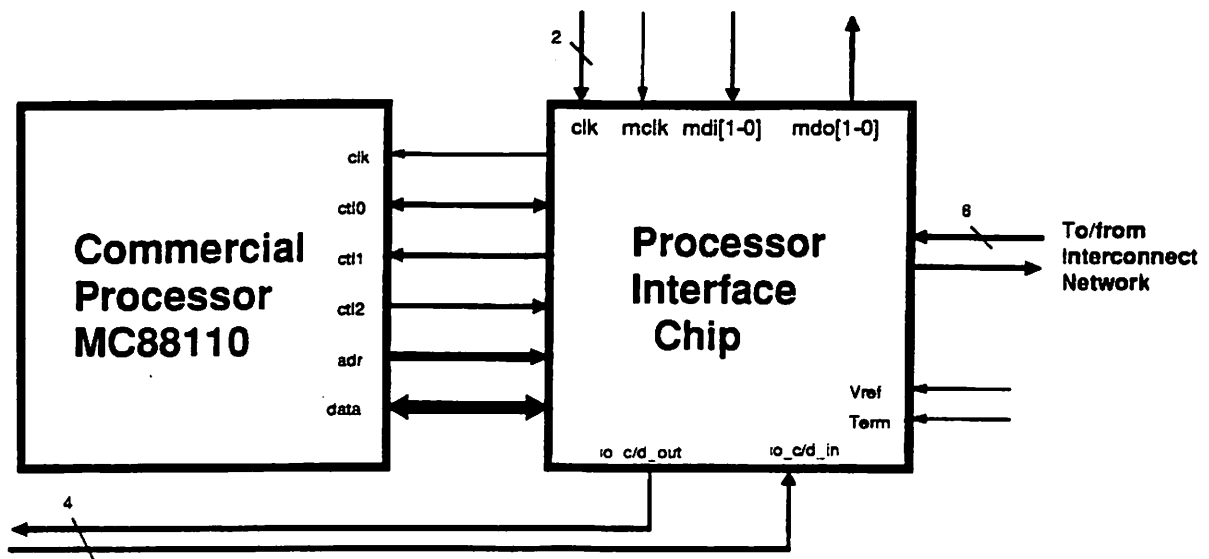
Monarch Characteristics

- * Shared memory machine, all memory is equally accessible to all processors with constant access time**

- * System consists of:**
 - Processor modules = Motorola MC88110 plus custom coprocessor interface to interconnection network**
 - Interconnection network = multiple stages built from just two custom switch elements**
 - Memory modules = two identical custom network interface chips and 34 SRAM chips**
 - I/O subsystems, VME bus compatible**
 - Utility modules**

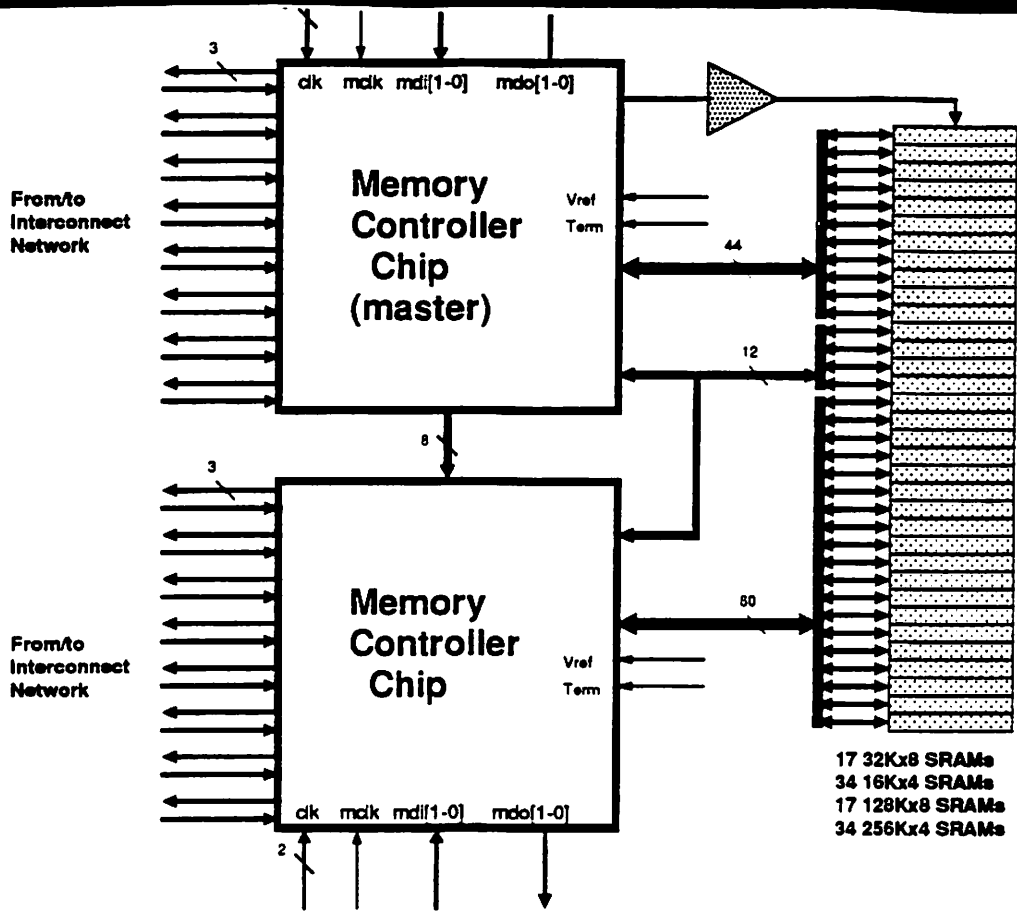
- * Interconnection network supports**
 - very low latency transactions**
 - pipelined operations for high wire utilization**
 - combining reads and writes**
 - large plurality of paths from any processor to a memory module, paths selected via random number generator**
 - M>N path redundancy for fault tolerance**
 - continuous link monitoring**
 - link testability during machine operation**
 - process, temperature, & link delay compensation of all links during machine operation**
 - extensible to very high bandwidths**
 - memory module to processor communication path during each frame for very fast processor to processor interrupts**
 - very low blocking probability**





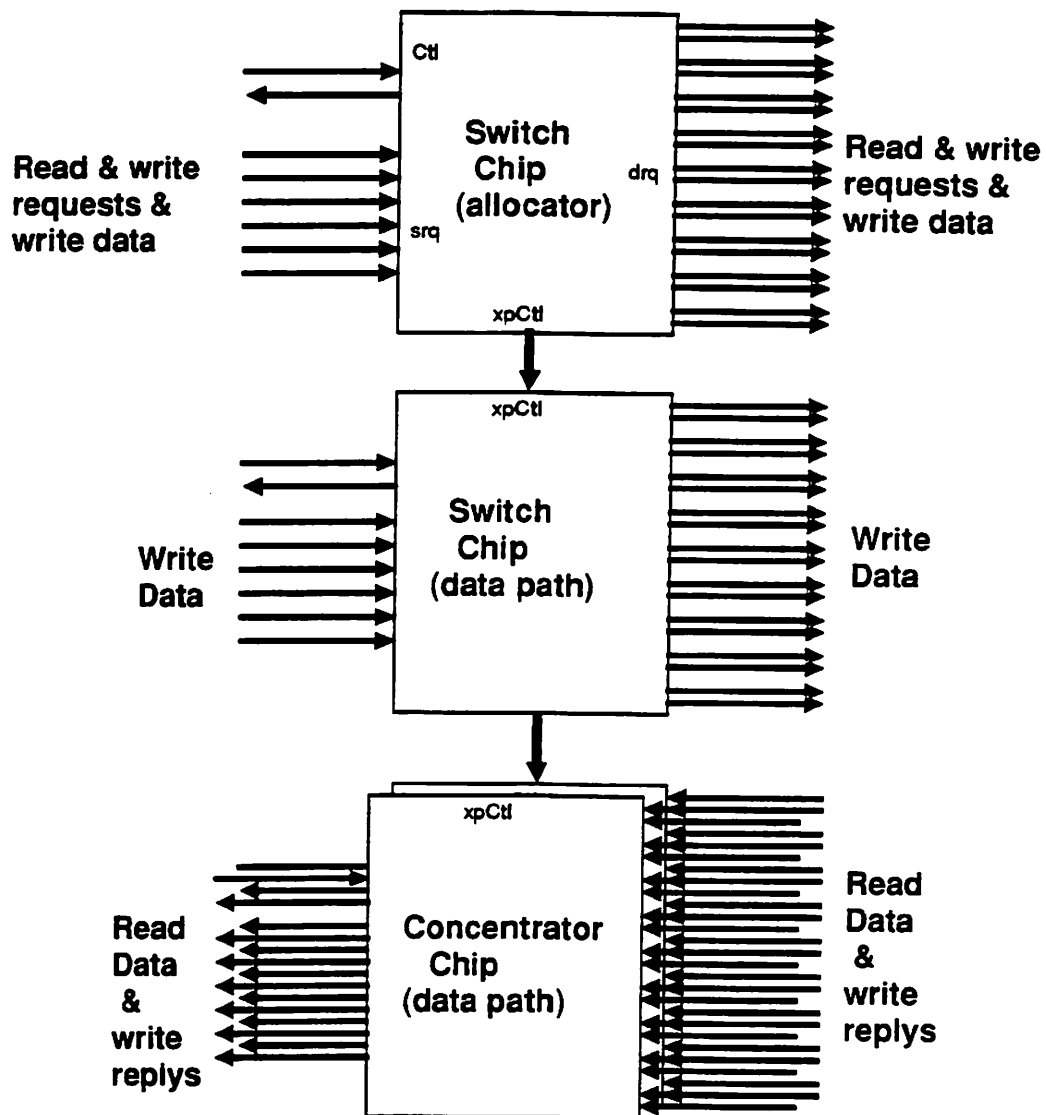
- * Management of instruction fetches & operand loads and stores
- * Small cache of recently fetched operand strips
- * Address hashing
- * Normal & low priority transactions
- * Steal, write if not stolen, read & write combining
- * Real time clock interrupts
- * Processor to processor transactions
- * Configuration dependent request formatting
- * Random path assignment
- * Reset, interrupt vector, bootstrap code





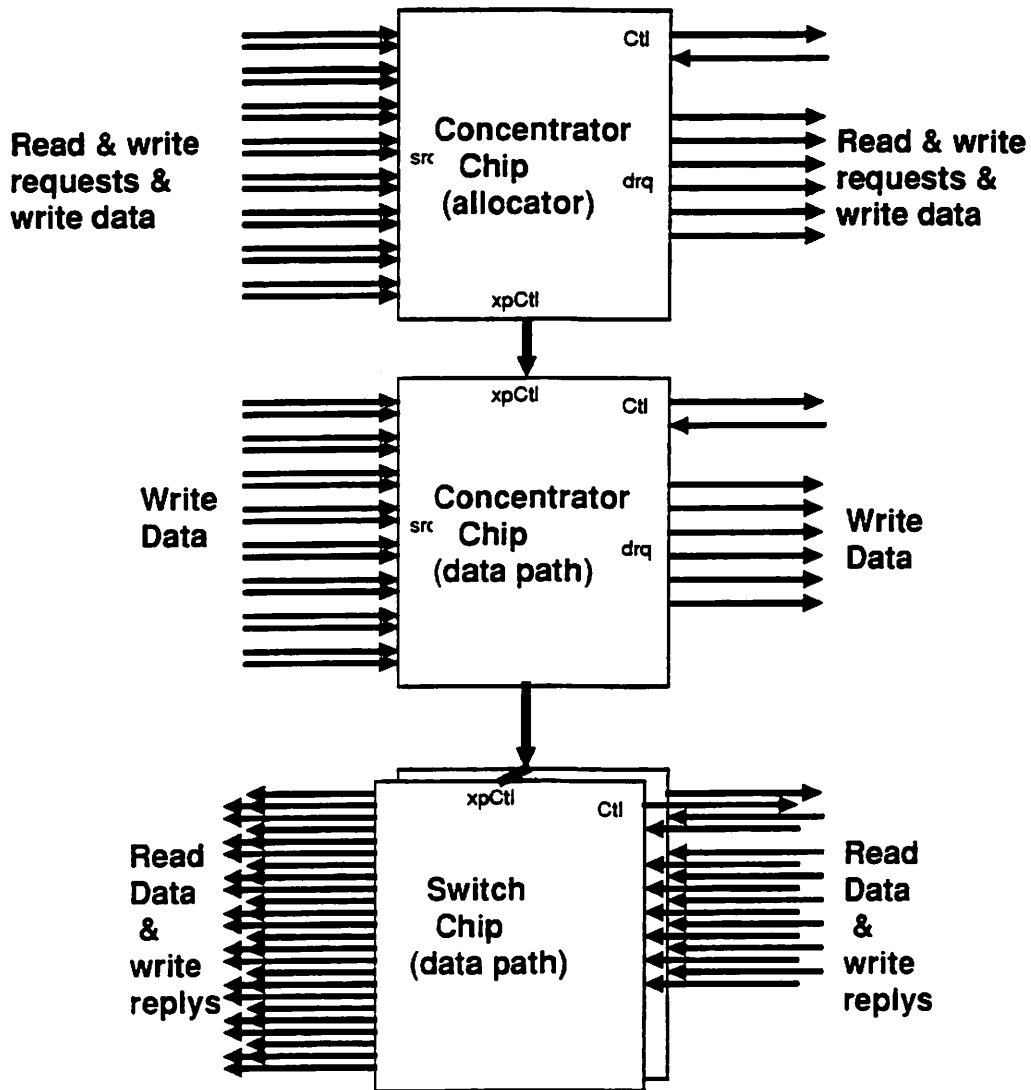
- * 8 channels with three input & output wires
- * MC supports two modes: master and slave
- * Three reads & five writes per frame imply blocking probability under 6%
- * Small number of write queues per channel
- * In address mode, MC captures request addresses, arbitrates for three SRAM access slots, controls the SRAM, & broadcasts control info to companion chip
- * In slave mode, MC inputs control info & transfers data from SRAM buffers to switch channel
- * MC supports read & write combining, write if not stolen, and steal operations
- * Double error detection, single error correction
- * Mailbox in each MC transmitted to corresponding processor each frame





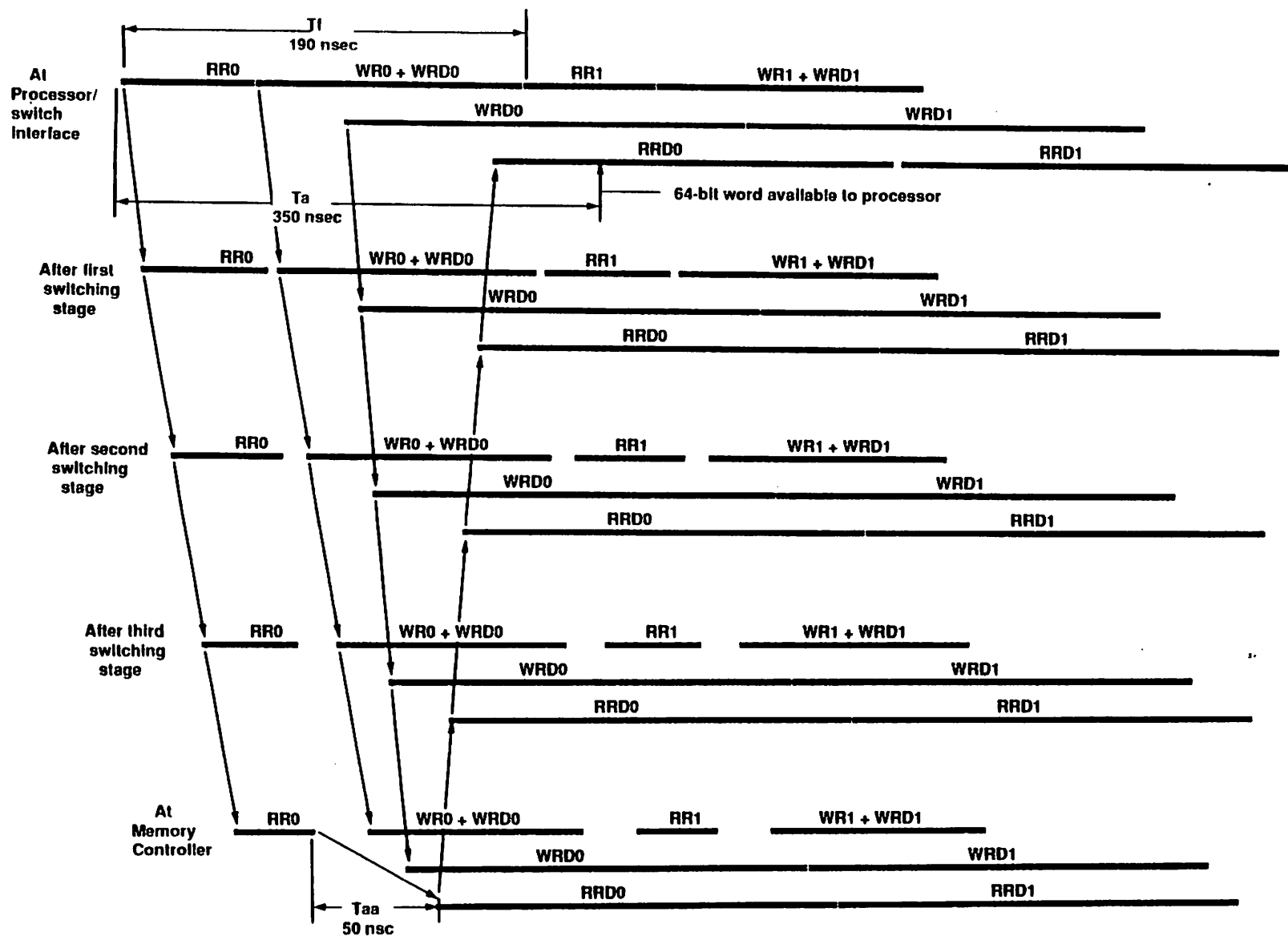
- * 6 inputs, 8 pair of outputs, 8-way switching
- * Redundant pair of outputs
- * 12 wires per channel
- * 256 bit reads and 256 bit writes
- * Net aggregate bandwidth 337 Mbytes/sec
- * Latency <20 nsec
- * Frame period 190 nsec



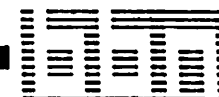


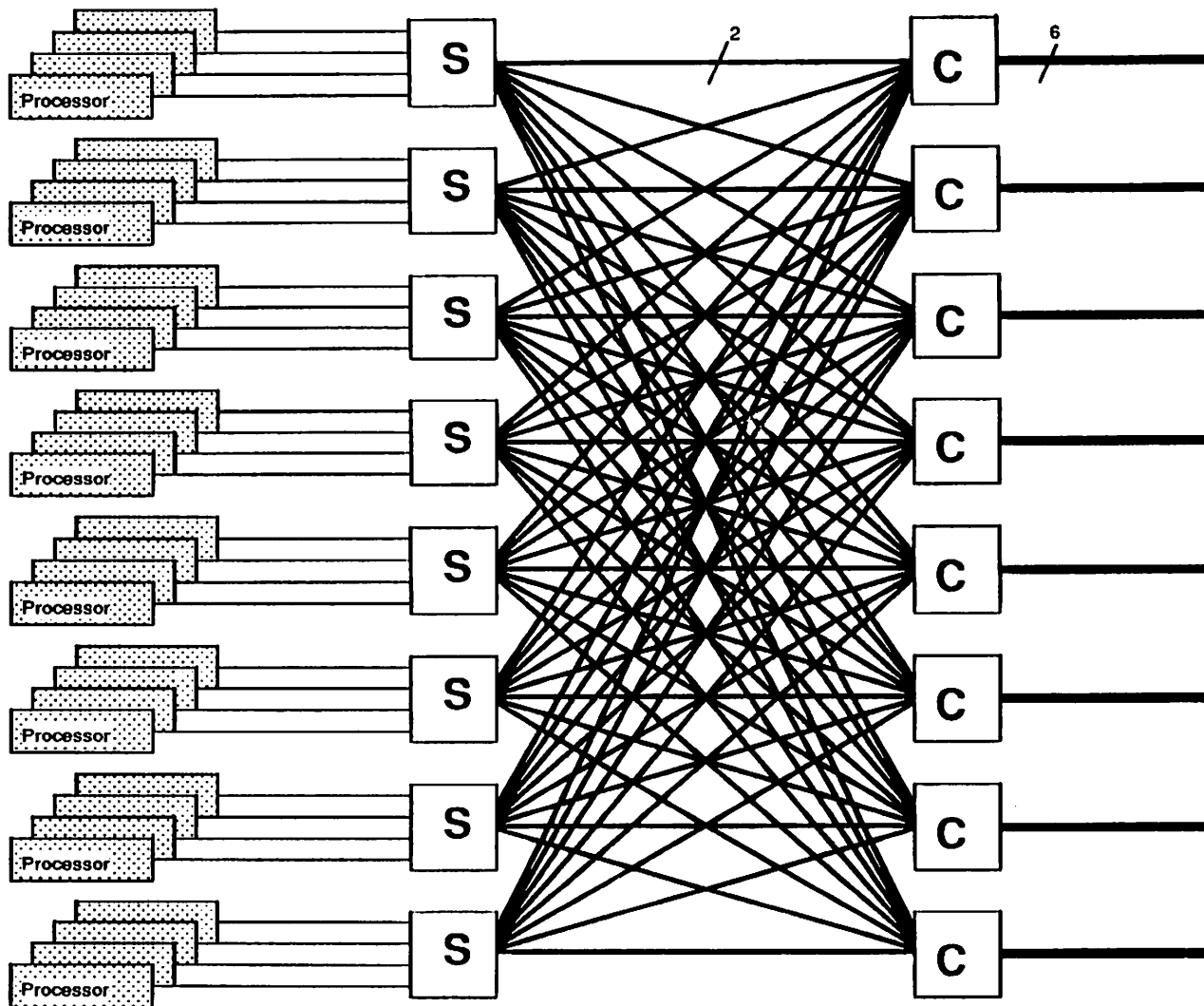
- * 16 inputs, 6 outputs, one way switching
- * 12 wires per channel
- * 256 bit reads and 256 bit writes
- * Net aggregate bandwidth 337 Mbytes/sec
- * Latency < 20 nsec
- * Frame period 190 nsec





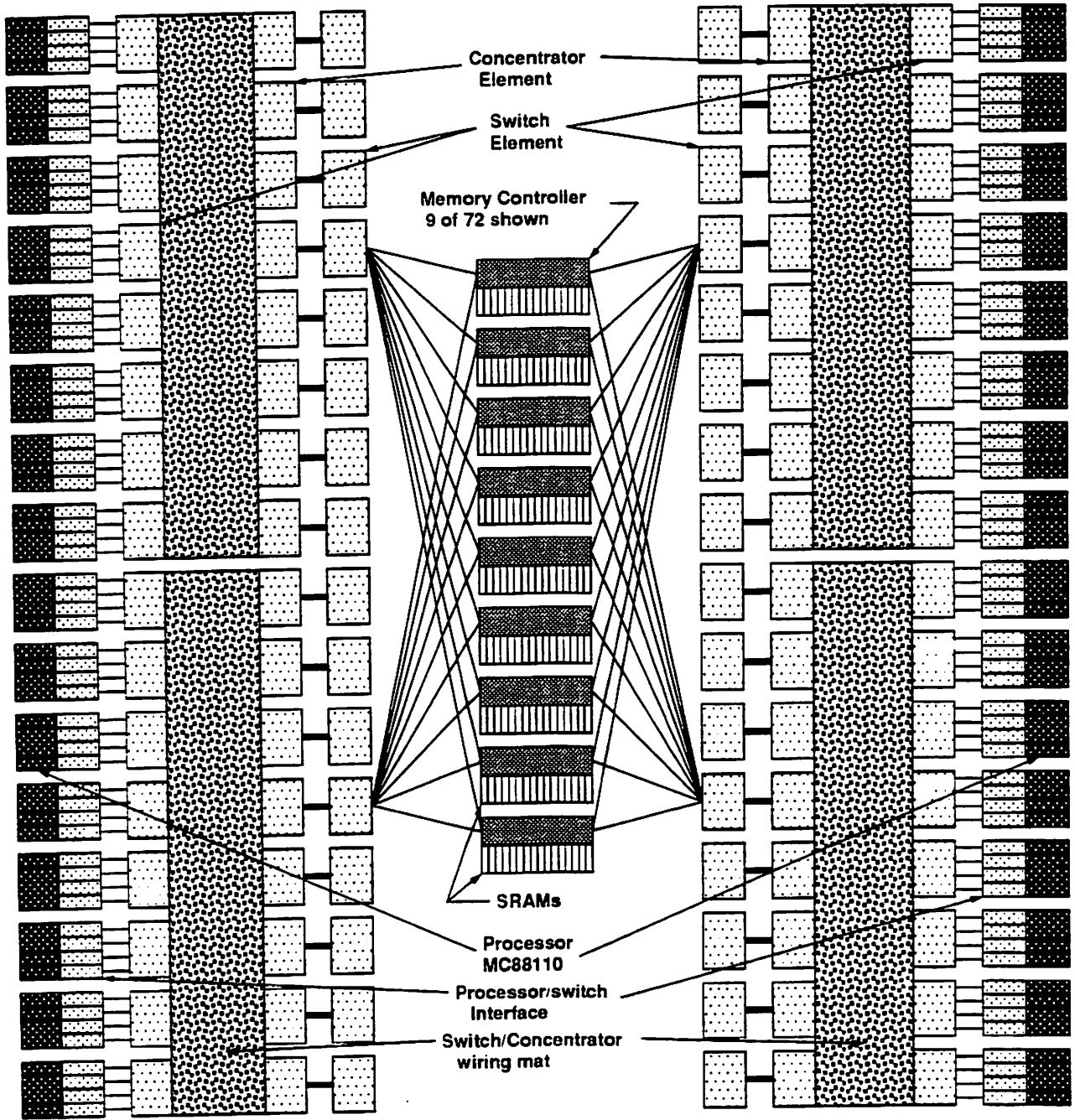
Request and reply timing for two consecutive reads and writes





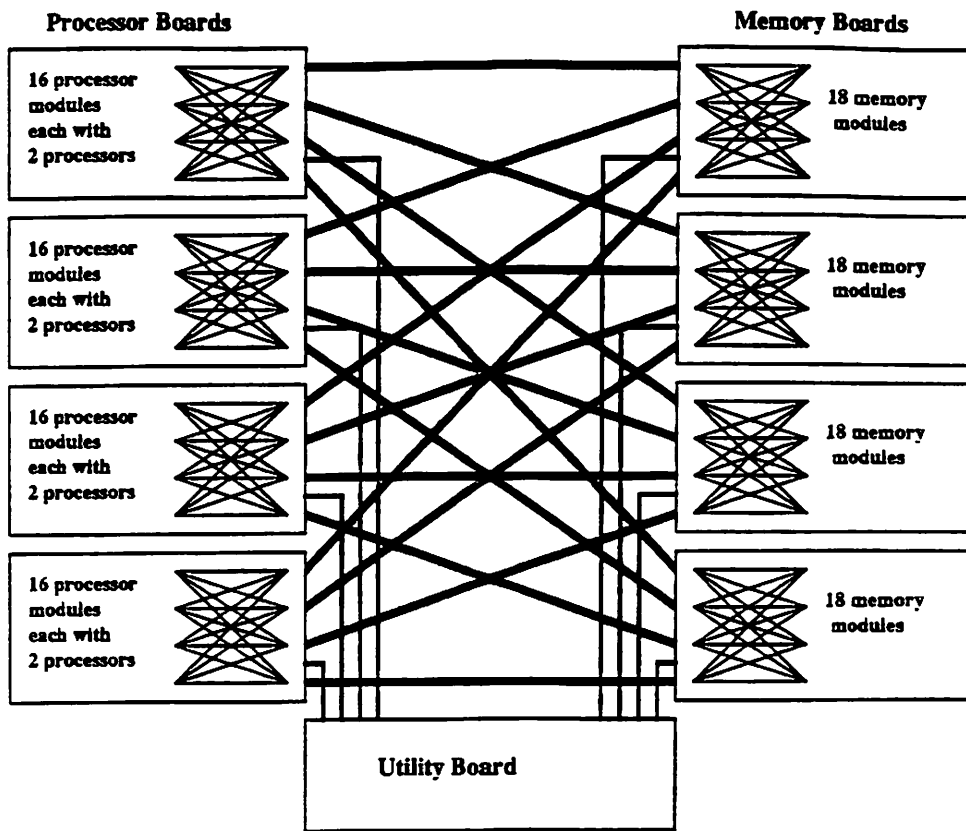
Monarch Processors and first stage switching



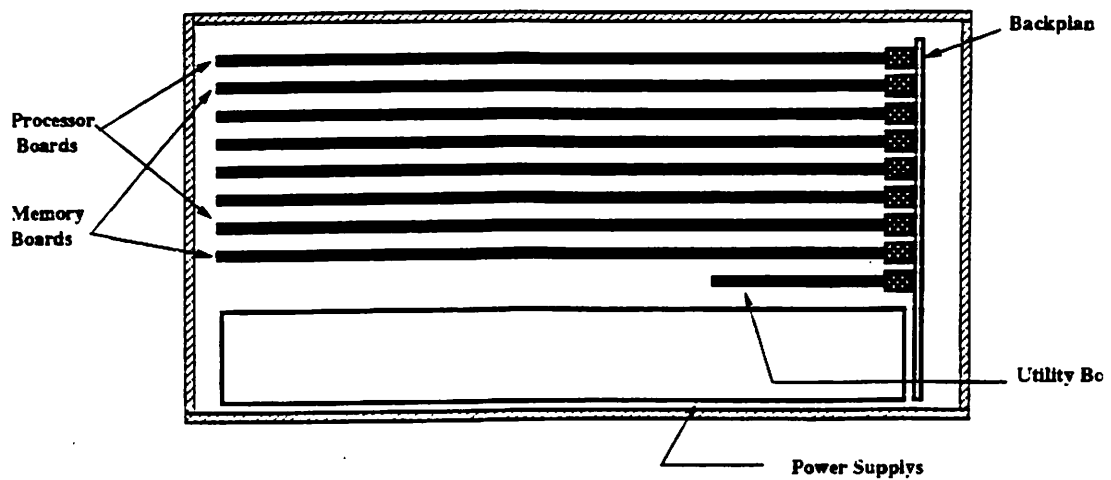


Monarch 128 Processor Block Diagram



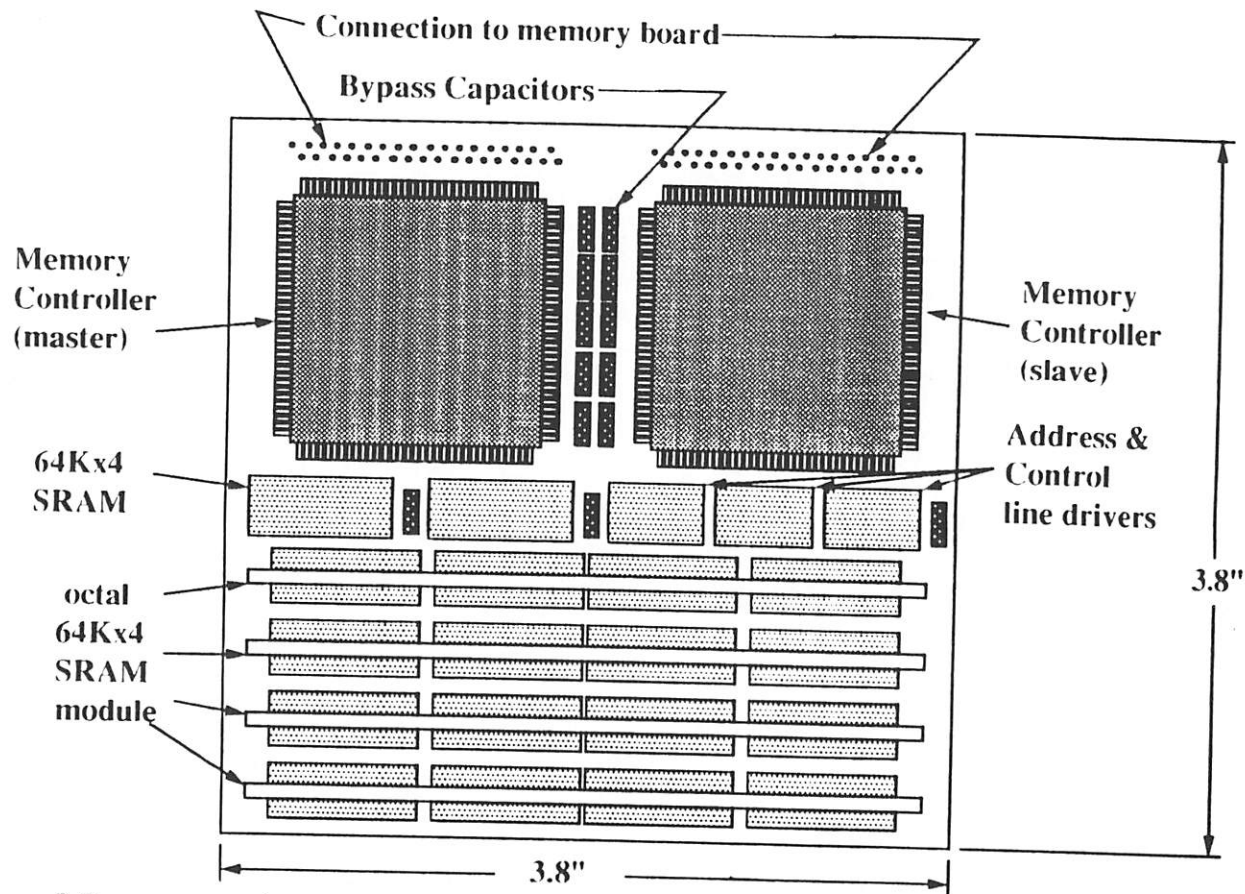


128 Processor Monarch as a set of interconnected boards

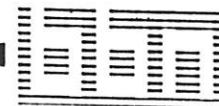


128 Processor Monarch card cage

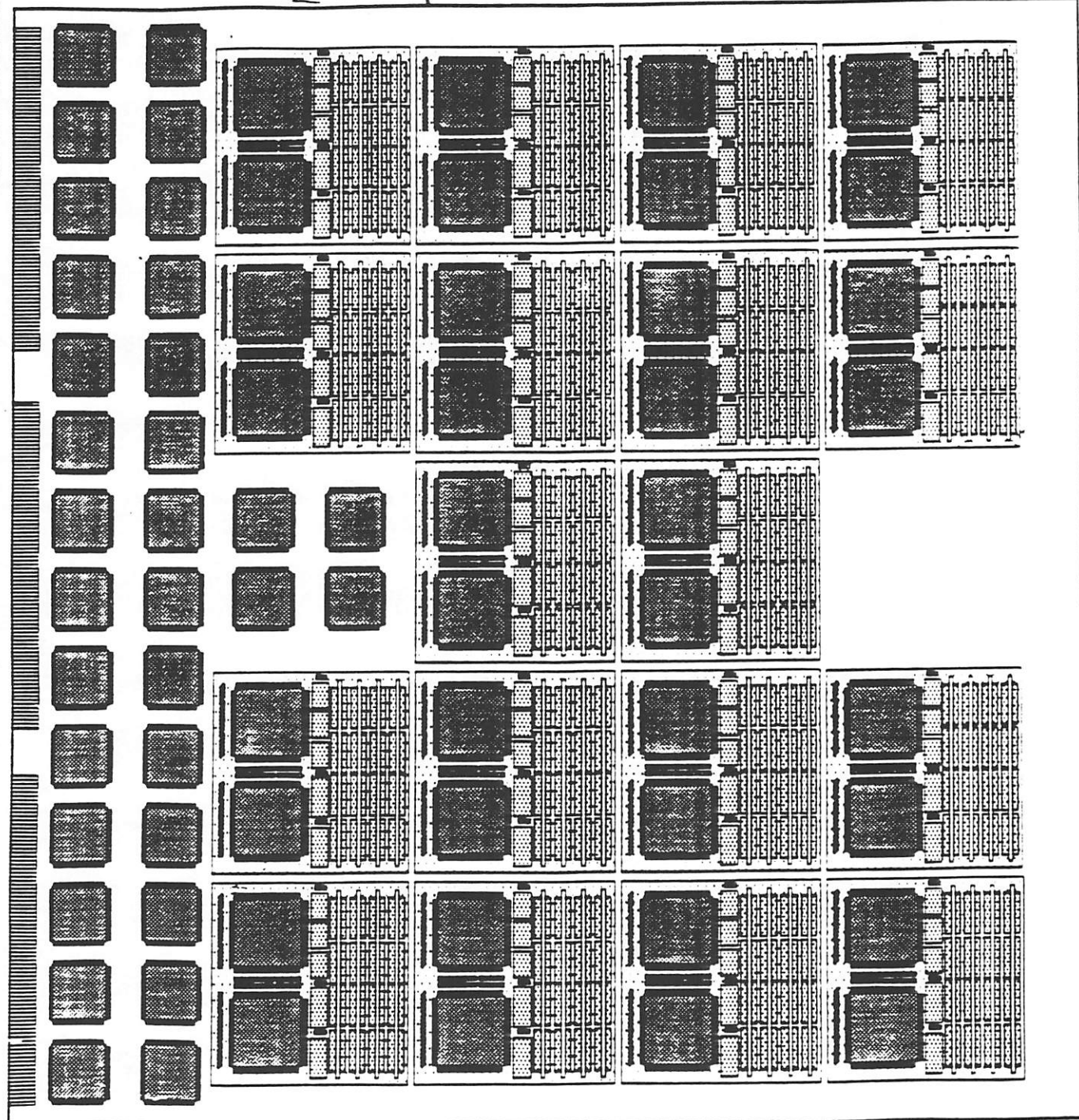


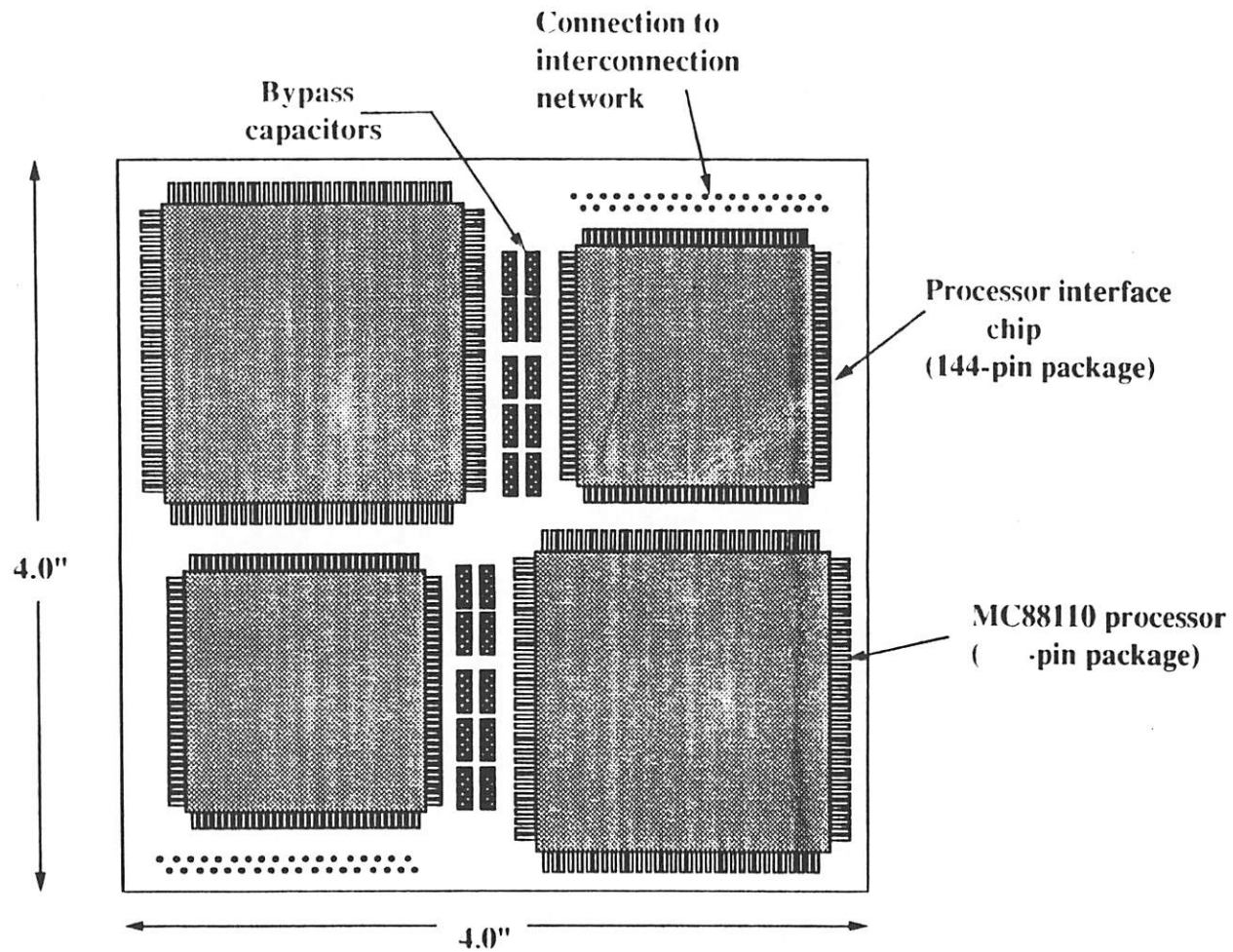


Memory element 136-bit accesses, 28.6-nsec cycle

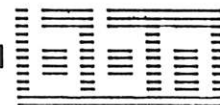


MONARCH
MEMORY
BOARD





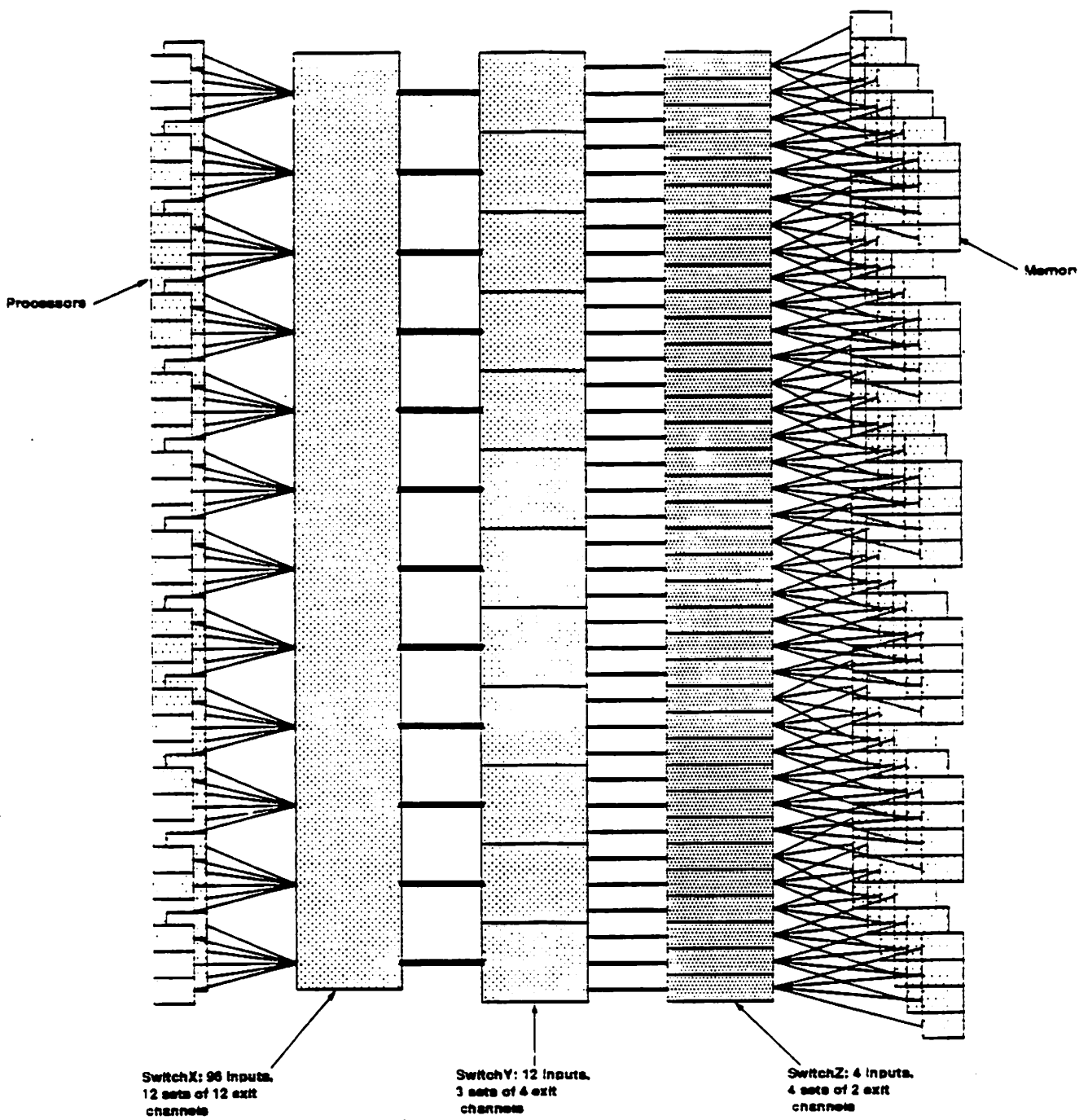
**Processor Module consisting of two MC88110s
and two processor interface chips**



High Density Monarch

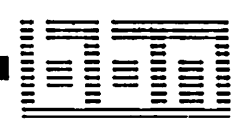
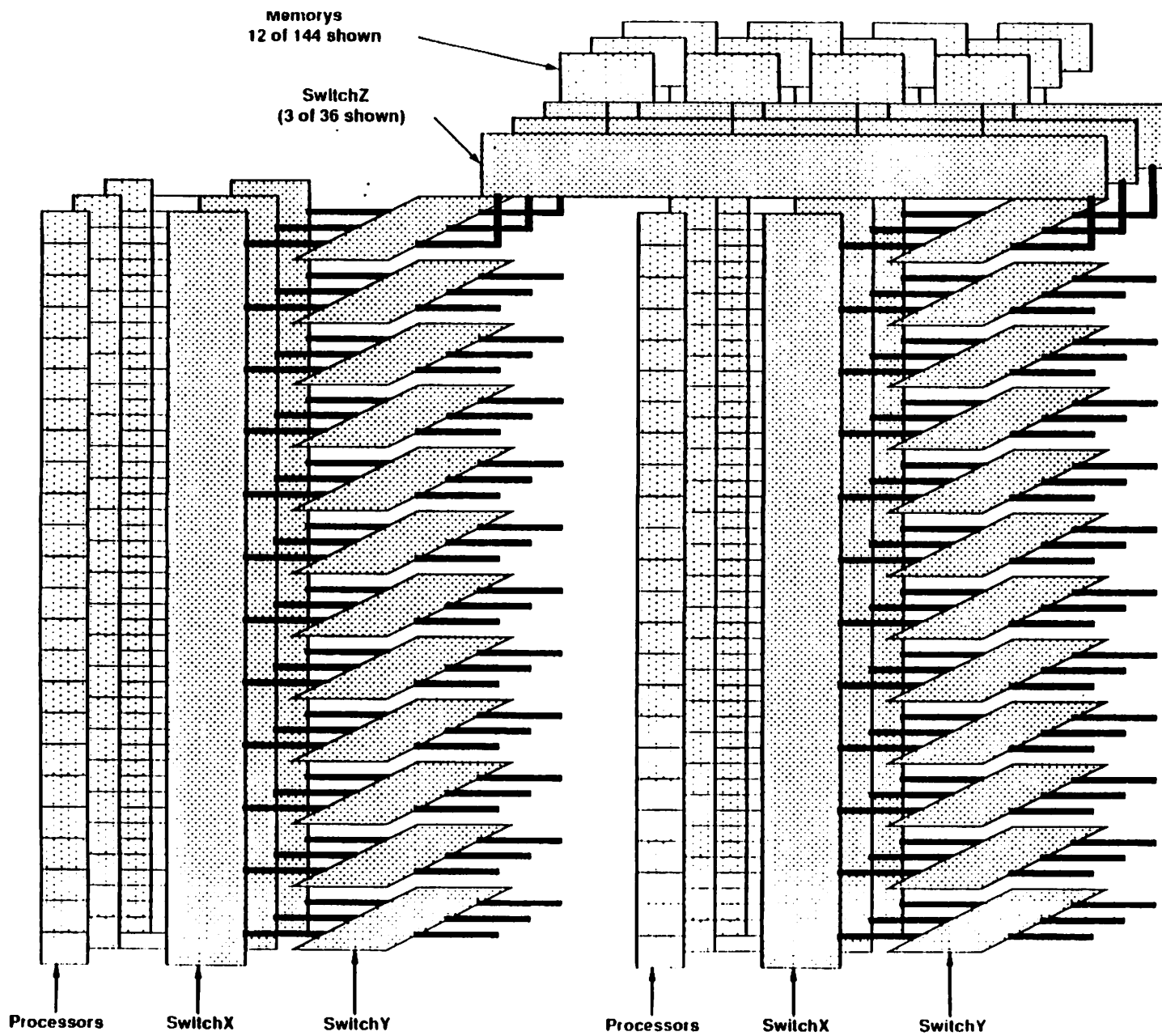
- * Standard PC board and backplane packaging technology for smaller machines, high density packaging for large machines**
- * High density packaging desired because of mechanical limitations, backplane-cable-connector constraints**
- * Packaging very difficult problem for large machines**
- * Requirements:**
 - small size to reduce speed of light latency and skin effect losses in transmission lines**
 - low temperature rise to decrease failure rates**
 - serviceability without major machine disassembly**
 - control of intermodule signal noise coupling**
 - fault tolerance of power distribution**
 - fault detection without invasive probing**
 - small number of different mechanical assemblies**
- * Ideal machine would use no cables interconnecting assemblies - optical interconnect**
- * Liquid cooling required above 2-3 KW per cubic foot, ideal machine should not use immersion cooling**
- * Ideal machine would be volumetric homogeneous over some minimum sized volume**

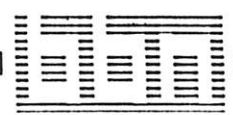
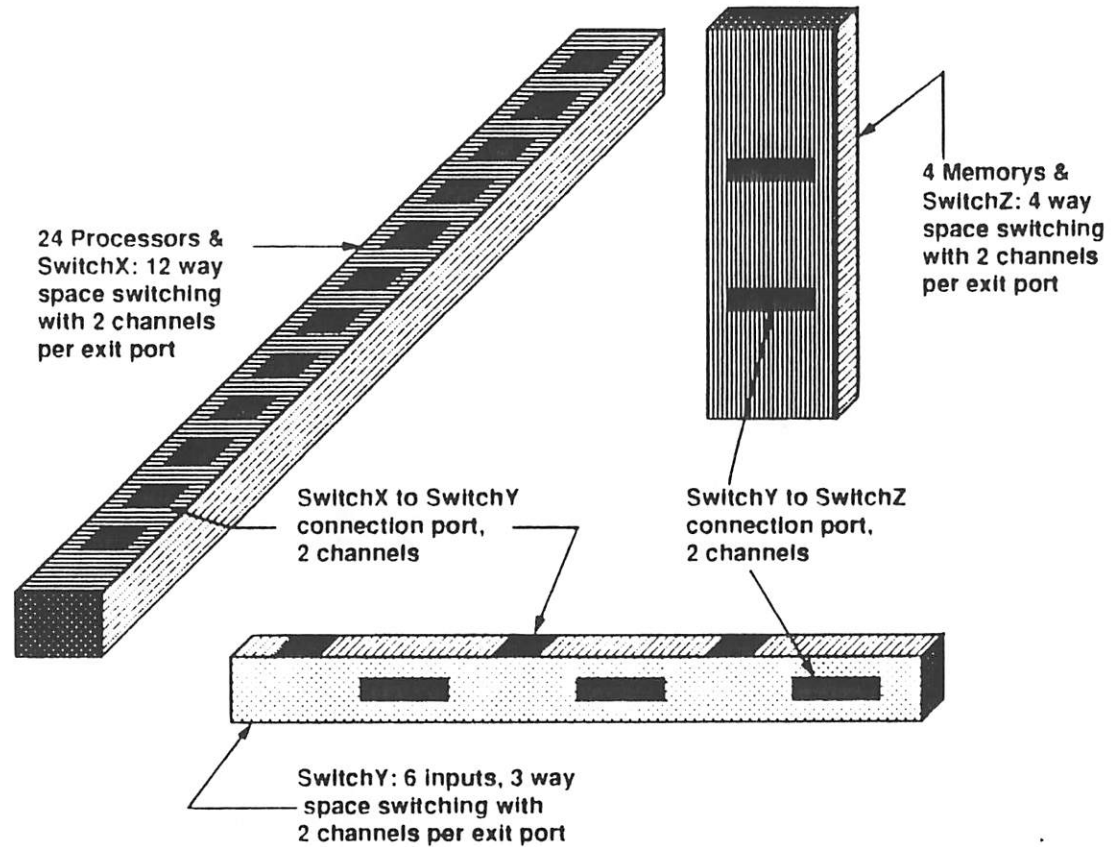


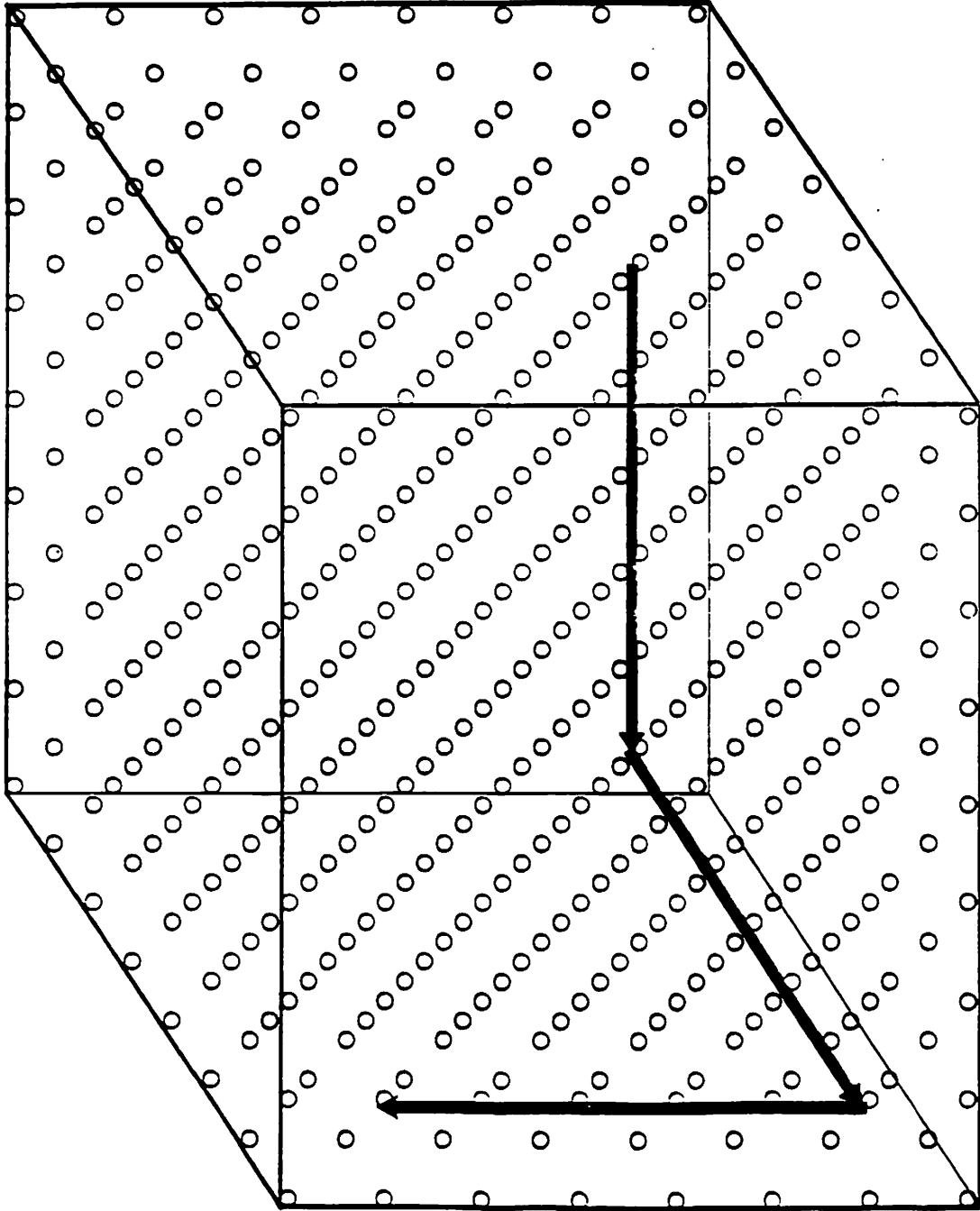


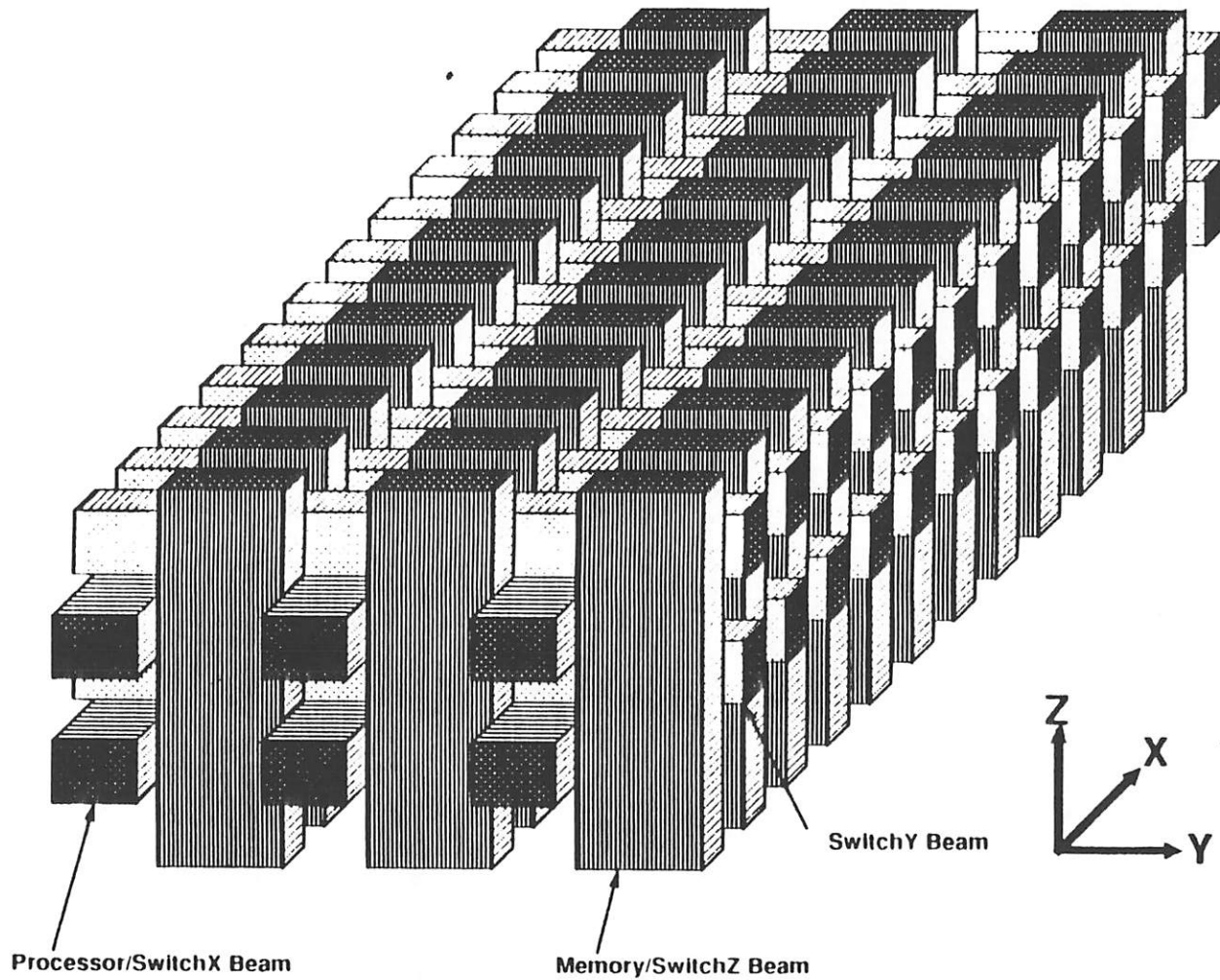
96 Processor/144 Memory system interconnected via a three stage indirect network S12:3:4



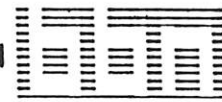


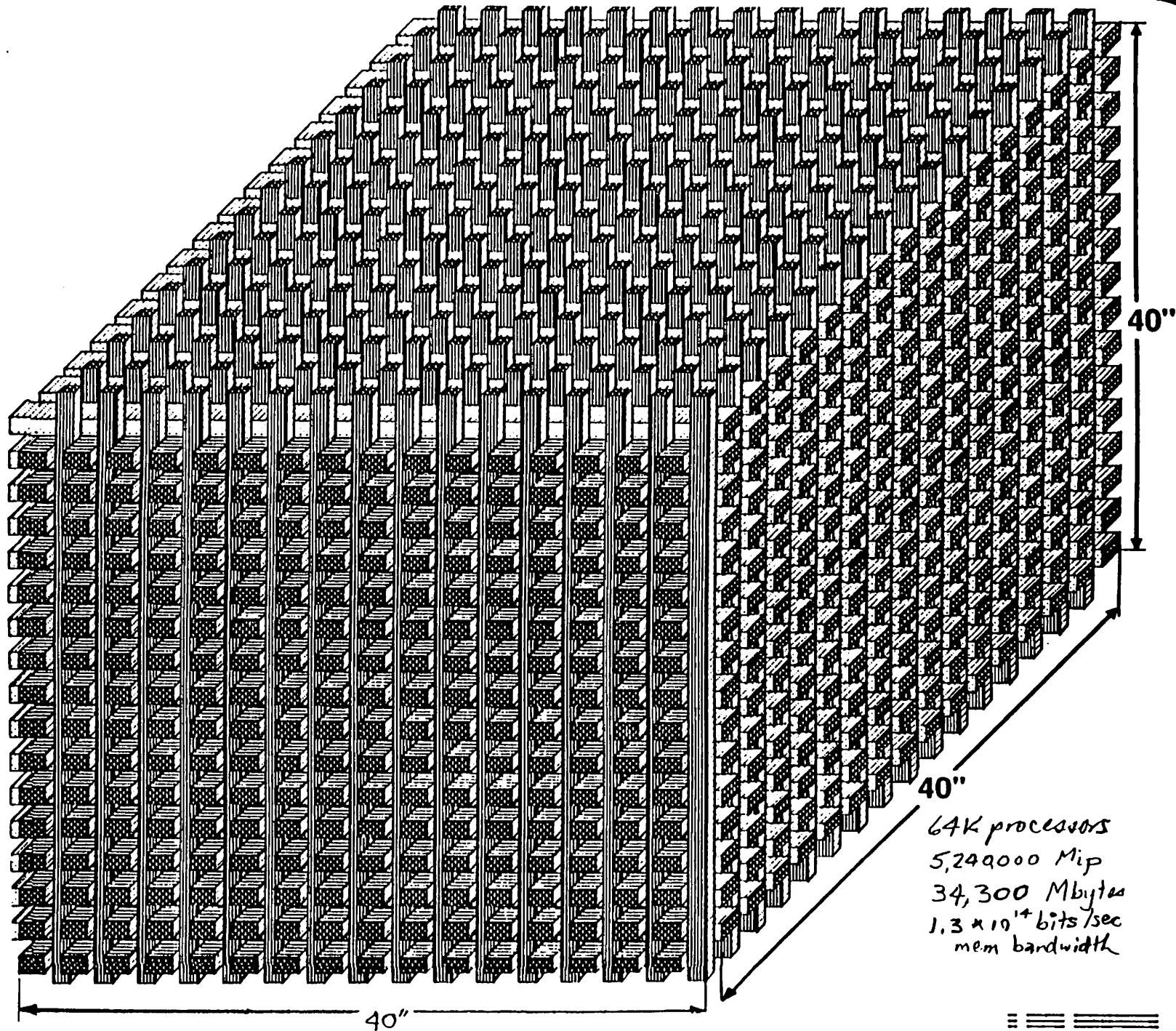




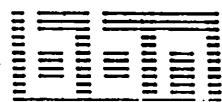


96 Processor/144 Memory system with 12:3:4 switching

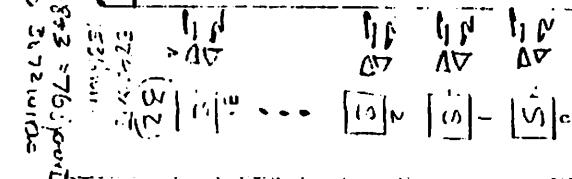
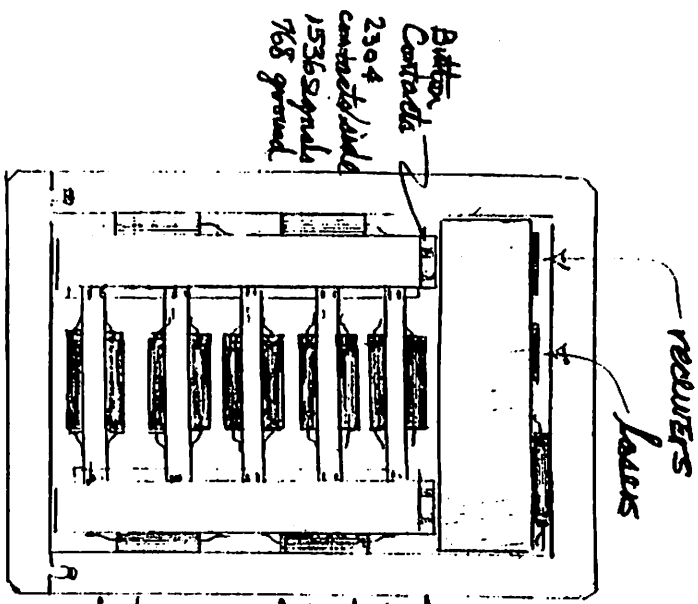




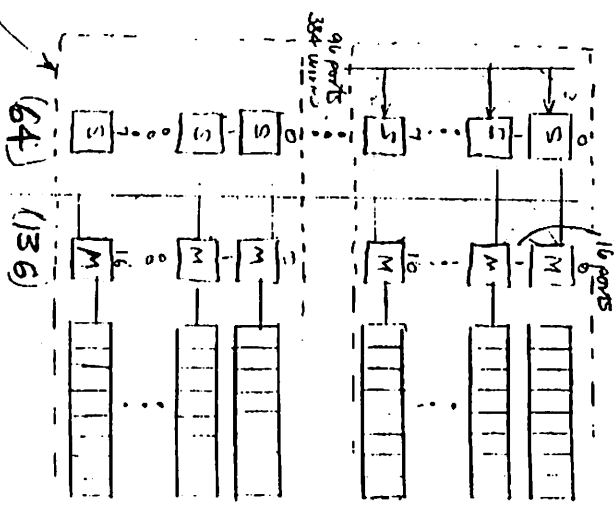
64K processors
5,240,000 Mip
34,300 Mbytes
 1.3×10^{14} bits/sec
mem bandwidth



SSM MODULE PACKAGING



SSM Module
 192 S
 136 M
 1224 SRAM
 16 ports

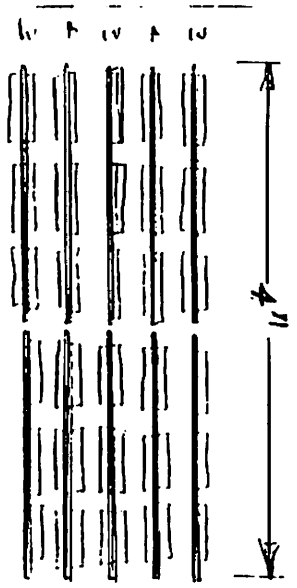
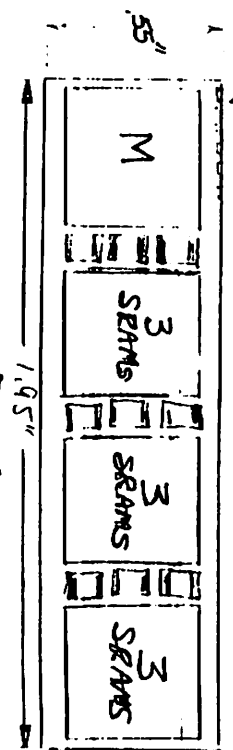


Allow 4" per sub module containing
 16 SMD chips
 17 memory controller
 153 static RAM's
 stacked 3 high



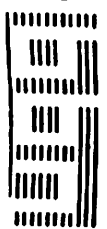
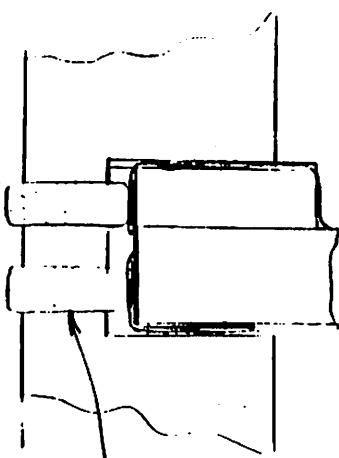
Note submodules may be removed from SSM module and disassembled down to dual M+9 SRAMs

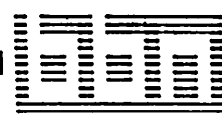
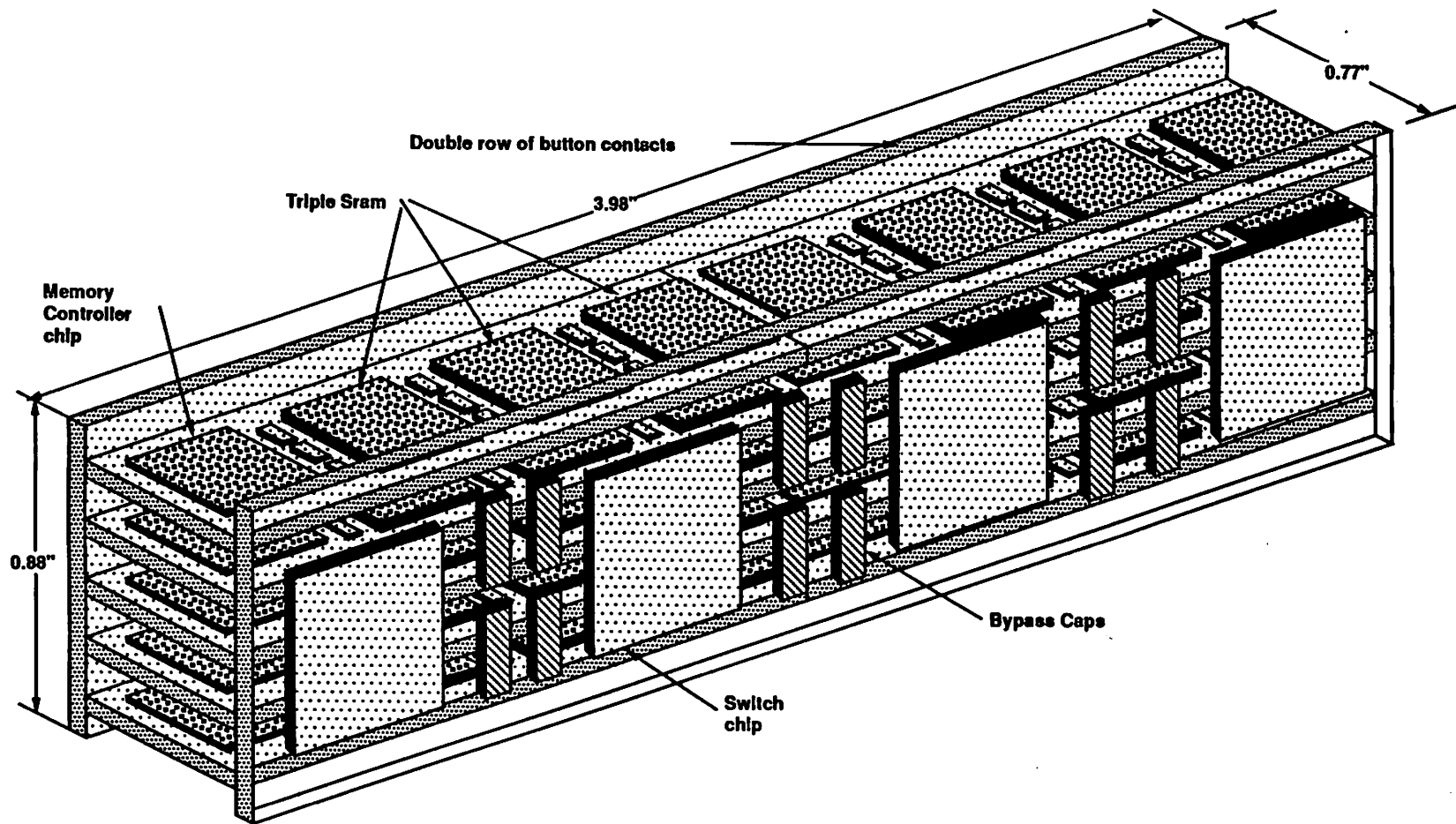
2 rows button contacts on each side
 @ 25 contacts in total 200 #/in



Sub Module side view

Integrate DC/DC converters (very small) as in PC module



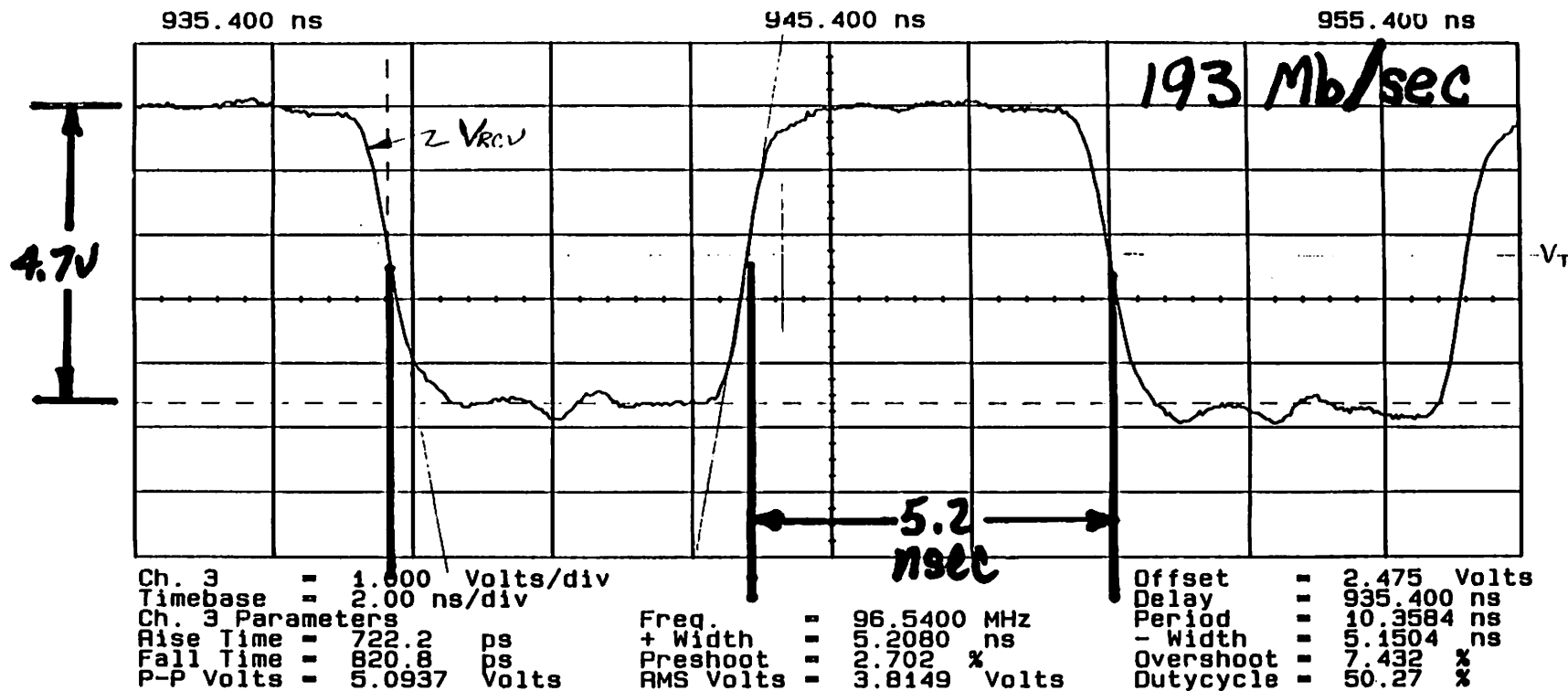


Monarch Enabling Technology

- * High speed signaling**
- * High density packaging**
- * High density DC/DC converters**
- * Integrated free space optical interconnect**
- * Liquid cooling**



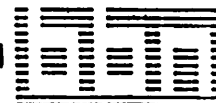
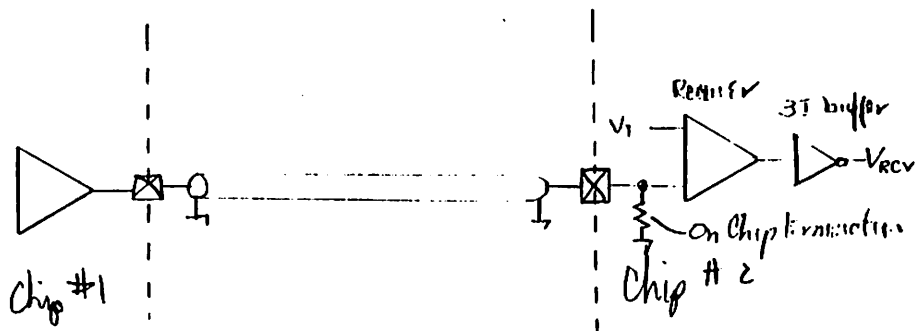
First Generation High Speed Signaling



10-12-79

SCM Signalling
 Data on chip (after ECL → CMOS level conv.)
 RXD (after 3-inverter buffer)

1.6 μ CMOS



Second Generation High Speed Signaling

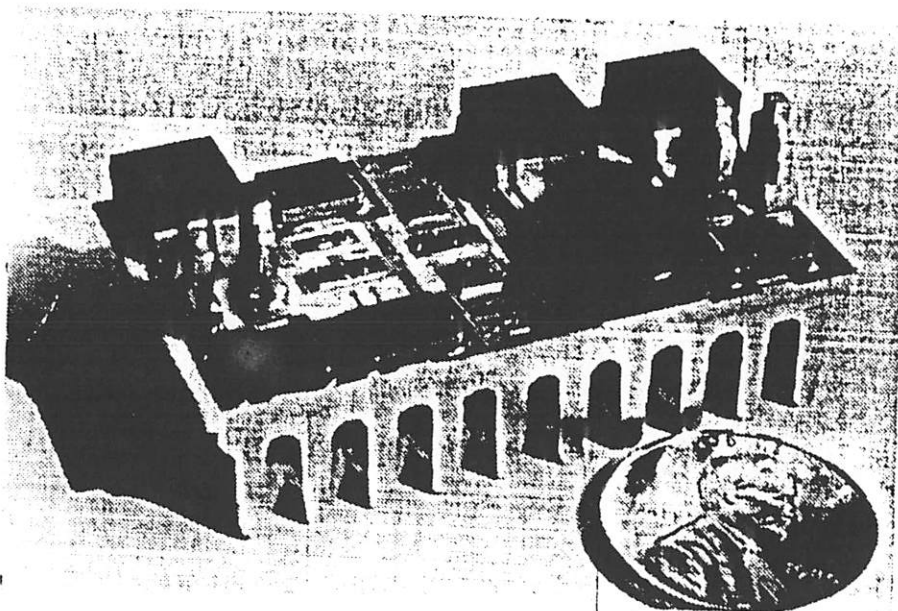
*** Next generation chips will incorporate enhancements to improve noise sensitivity, power discipation, process variation tolerance, and latency**

- PAD latchup protection**
- duty cycle correction of clock and data receivers**
- new chip synchronization scheme**
- gated clock circuits to minimize power discipation**
- process compensated line driver and receiver**
- digital controlled delay line for bit level data recovery**
- better testibility**
- static logic for low speed circuitry**
- low cost, high speed socketable chip package**
- new scheme for transferring data inside chip interior to minimize latency**



High Density Power Supplies

(a)



(b)

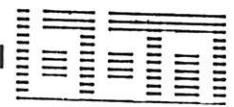
Fig. 1. (a) Block diagram illustrating the concept of a distributed power supply for a mainframe computer. (b) Physical prototype of the 10 MHz, 50 W, converter being developed at MIT for use in distributed power supplies.

2 sample 2 to 1

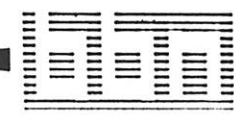
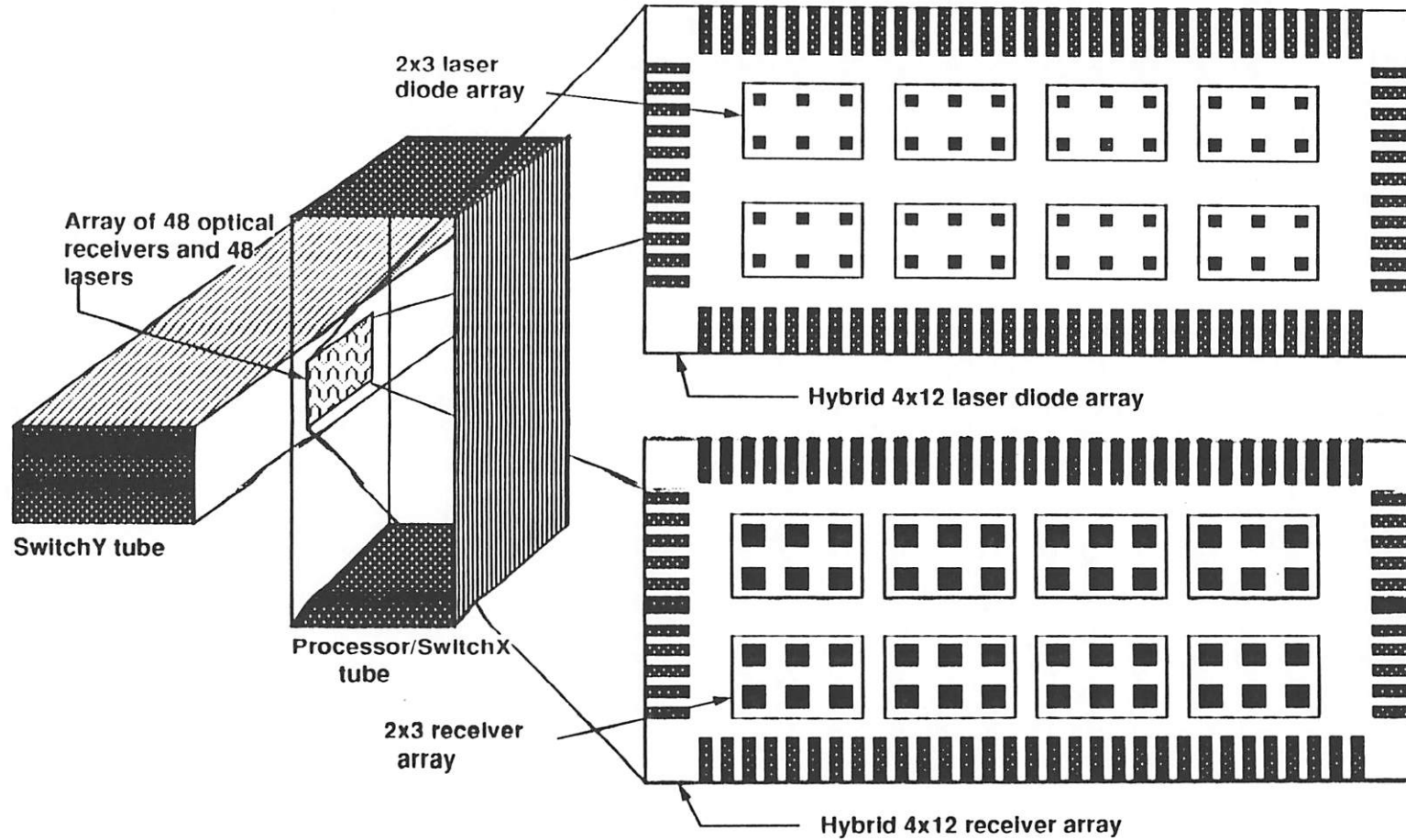
$4.5 \text{ cm} \times 1.75 \text{ cm} \times 1.25 \text{ cm}$
 $1.77" \times .689" \times 0.492"$
 0.603 in^3

penny

1000



Integrated Free Space Optical Interconnect



Monarch

Software

- Programming Model
- Operating System
- Synchronization
- Parallel Debuggers and Tools
- Latency
- Machine Model



Monarch

Programming Model

- **MIMD Multiprocessor**

Virtual Multiprocessor - a partition of the machine.

- **P** Identical Processors - Space Shared

- **Uniform Shared Memory Paradigm**

- Large uniformly accessible shared memory
- Weak Coherence

- Message Passing Paradigm

500 nano-second processor-processor interrupt

- Dataflow Paradigm

I-structure references (tagged memory)



Monarch

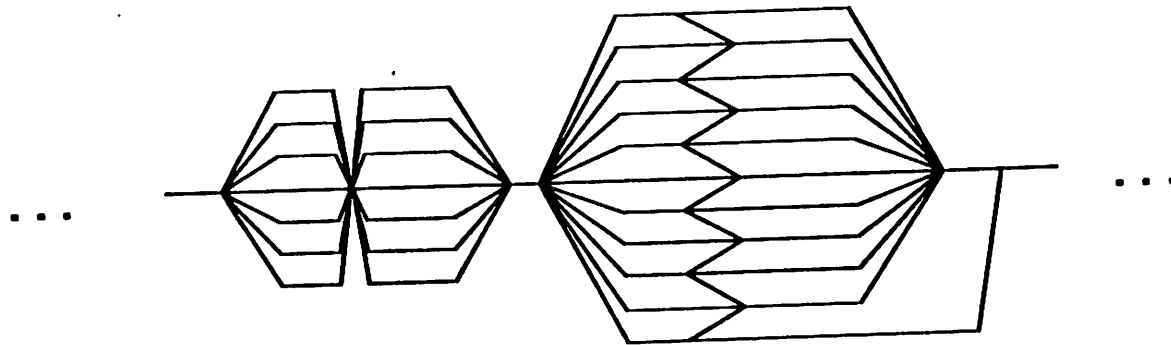
Granularity

- **Fine Grain Parallelism.**
 - **Single Program/Multiple Thread**
 - Procedure Level
 - Explicit Specification of Parallelism
 - Explicit Synchronization
 - **Coarser Grain Parallelism Readily achieved**
- **System Support - Fine Grain Shared Memory**
 - Operating System
 - Runtime Library
 - Programming Language(s)
 - Hardware



Monarch

Program Execution



Process

Sync

Process

Combine

Process

Locally Combine

Process



Monarch

Operating System - MACH

- **BBN ACI enhanced MACH**

Multiprocessor Unix variant - 4.3 BSD compatible

- DARPA supported research (via CMU)
- ACI versions on 2 different multiprocessors
- ACI parallelized I/O and other bottlenecks

- **Task/Thread**

Fits Single Program/Multiple Thread

Large Scale Parallelism

- **MACH Virtual Memory**

Copy-on-Write

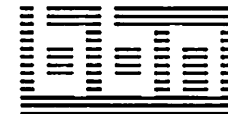


Monarch

MACH TASK

- **Execution Environment**
 - Basic Unit of Resource Allocation
 - Generally High Overhead
- **Virtual Address Space**
 - Ordered collection of mappings to memory objects
- **Access to System Resources**
 - Processors ● Virtual Memory ● Port Capabilities

Unix Process = 1 Task + 1 Thread



Monarch

MACH THREAD

- **FLOW OF CONTROL**

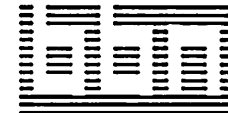
- Independent Program Counter
- Stack Pointer
- All THREADS of TASK share the TASK resources

- **Basic Unit of CPU utilization**

- Generally Low Overhead
- "Lightweight Process"
- "Parallel Coroutine"

- **Key to MACH support of parallelism**

Unix Process = 1 Task + 1 Thread



Monarch

MACH VIRTUAL MEMORY

- **COPY-on-WRITE Sharing**

Copying is delayed until a thread performs a write on COW shared page of memory. Only necessary copying.

- **READ/WRITE Sharing**

Memory object is mapped into multiple address spaces

- **NO Sharing**

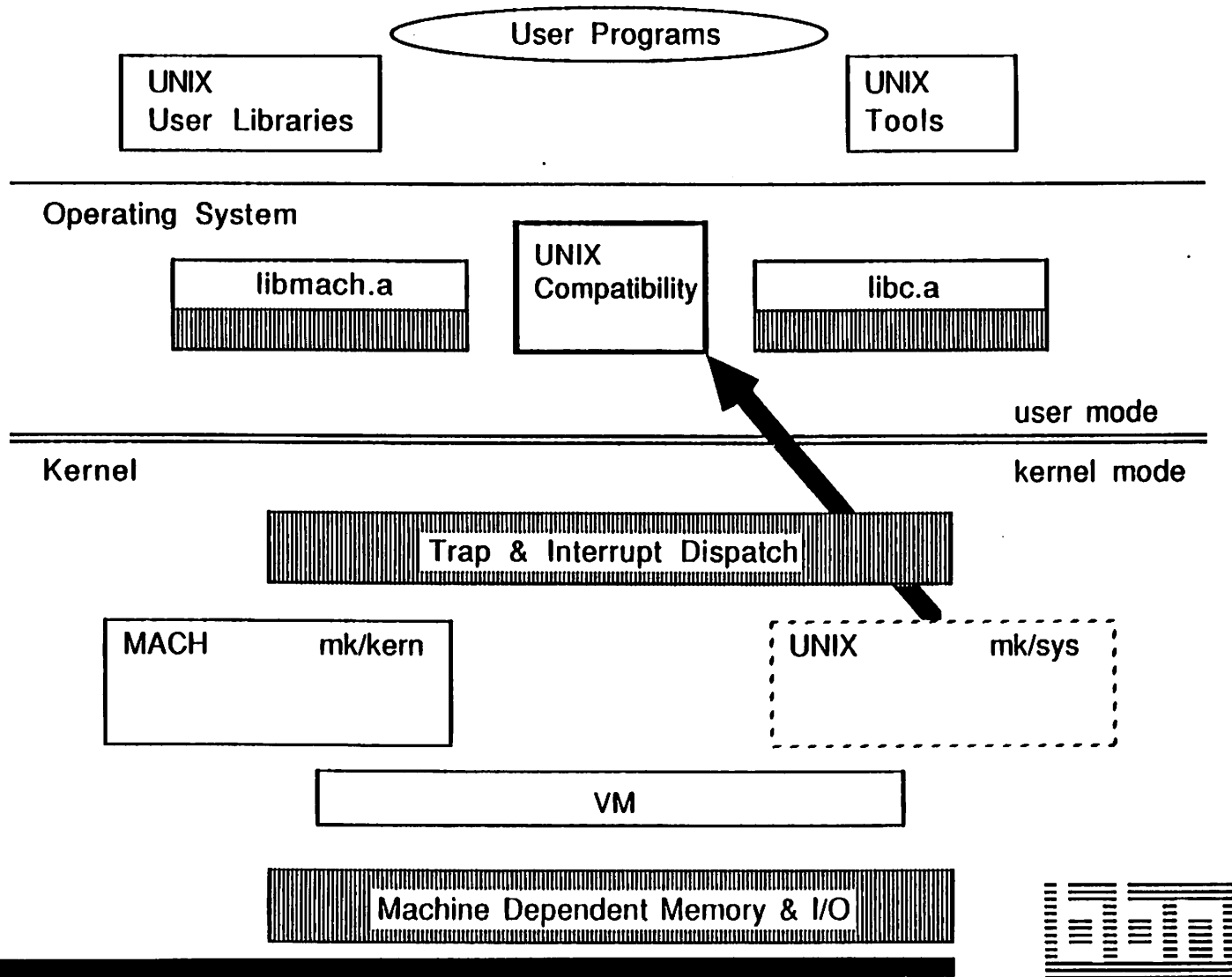
Memory object is private to a task.

INHERITANCE and PROTECTION
on a per page basis



Monarch

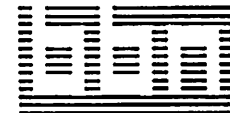
MACH Structure



Monarch

STEAL

- Every Memory Location is **tagged**
- **STEAL** operations sets "stolen" tag bit
- Store clears "stolen" tag bit.
- Load or Steal of stolen location hangs till tag bit is clear
 - Prefetch quietly fails
 - Override Steal No-retry
 - Other memory access variants
- **STEAL** used to construct synchronization mechanisms
 - LOCKS (spin-wait)
 - **LOG P** Barrier
 - Write-once storage
 - **LOG P** Combining and Prefix



Monarch

Simple Locks

A location which is stolen and then written implements a lock directly.

Code to guarantee exclusive access to the variable SharedCounter using the lock variable Lock

Steal Lock

Update(SharedCounter)

Lock = 0 ;Storing into lock "unsteals" it



Monarch

Barrier Synchronization

Use a divide and conquer approach

Synchronize $P / 2$ pairs of processors using simple locks implemented with Steal

Synchronize $P / 4$ pairs of processors, 1 from each pair in the first stage

Repeat recursively

As processors drop out they can poll a location. The last processor can set this location, releasing all other processors simultaneously.



Monarch

Data Synchronization

Use Steal Primitive Directly

1. Presteal Shared Data to be Computed
2. Compute Data
3. Write Data Once Computed
4. Use Data whenever needed relying on Steal mechanism to guarantee integrity of value

Technique useful to avoid task synchronization and pipeline tasks together when a piece of data is computed by a small number of processor and then read by a large number or random set of processors



Monarch

Write-Once Storage

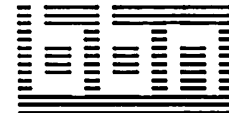
Write once storage can be pre-stolen and then used without synchronization. Honoring the stolen bit ensures that data is not used before it is produced.

for i from 1 to N Steal Y(i, myJ)

synchronize

for i from 1 to N

$Y(i, \text{myJ}) = f(Y(i-1))$



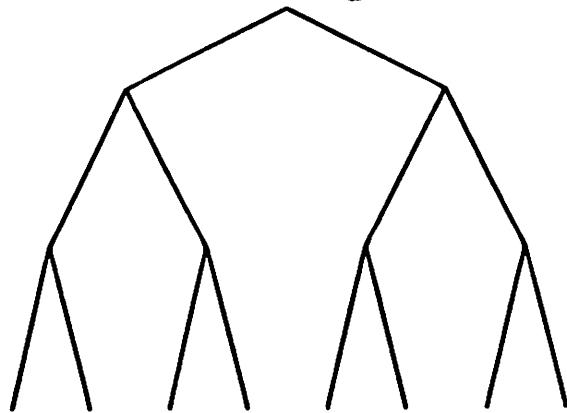
Monarch

Parallel Prefix and Result Combination

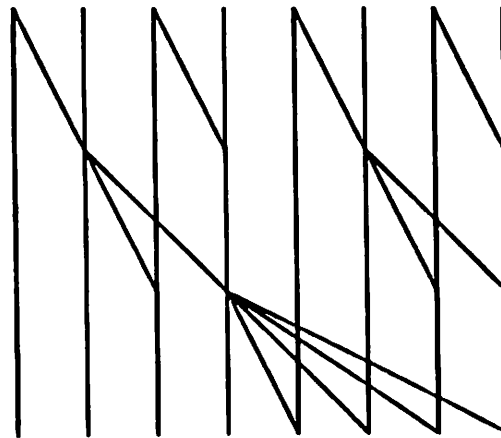
Use a divide and conquer approach as in barrier synchronization.

Combining read allows for simple algorithm

Combining Tree



Prefix Tree



Monarch

Parallel Debuggers and Tools

- **BBN ACI Experience in Parallel Programming**

- Window Oriented Source Level Debug
- **GIST** - Event Logger

Visualizing dynamic behavior of Parallel Programs

- **Parallel Unix Utilities**

- Make
- Grep
- Shell

Wildcards - `cp *.c /usr/darpa/newdir`

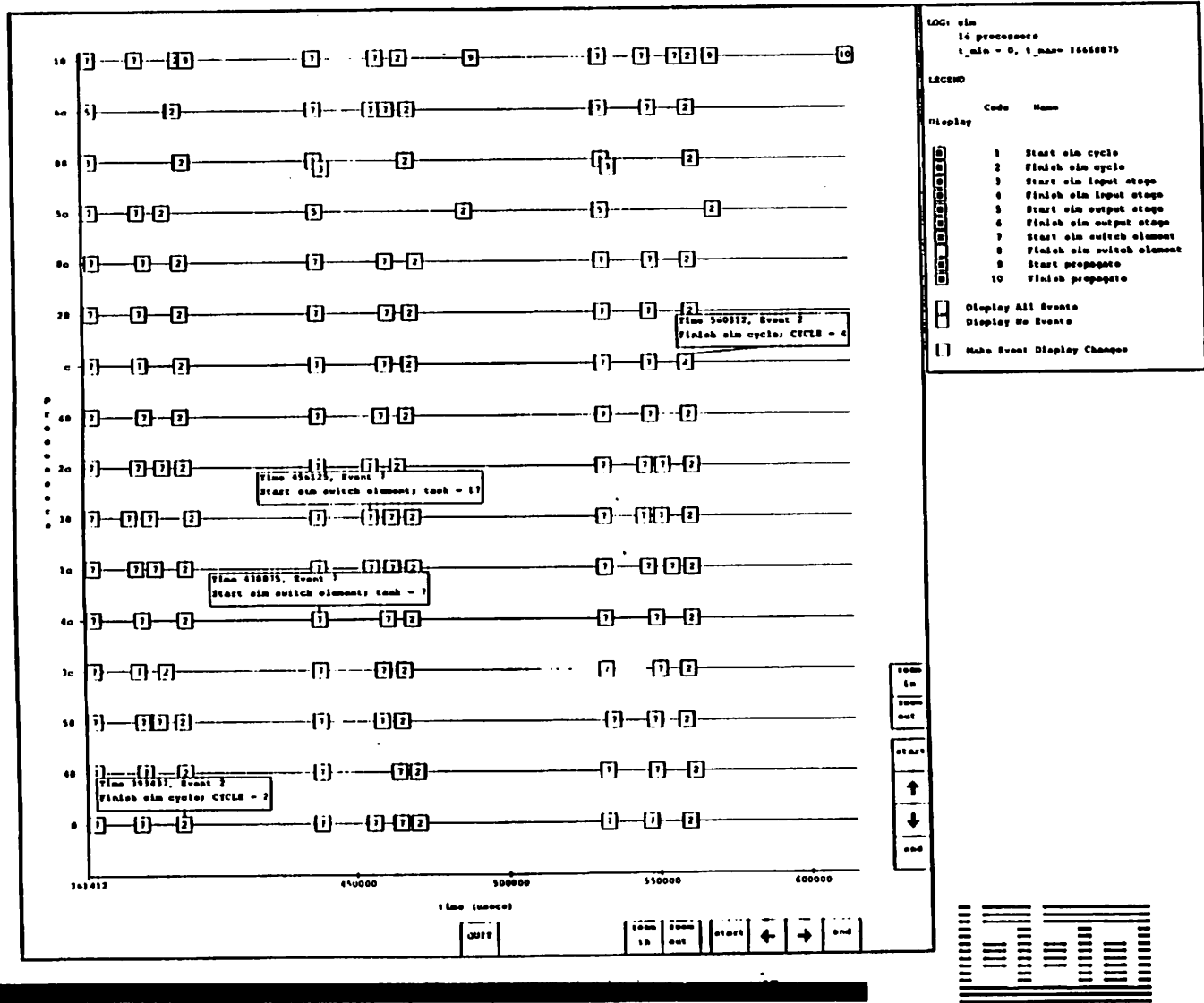
Pipelines - `tbl | eqn | troff | psprint`

Loop - `foreach f (*.c)`



Monarch

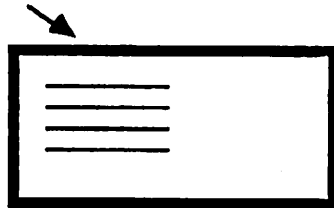
Gist



Monarch

Window Oriented Debugger

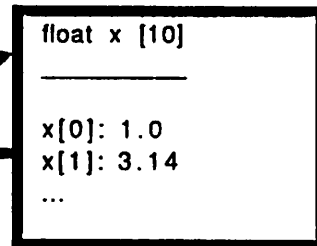
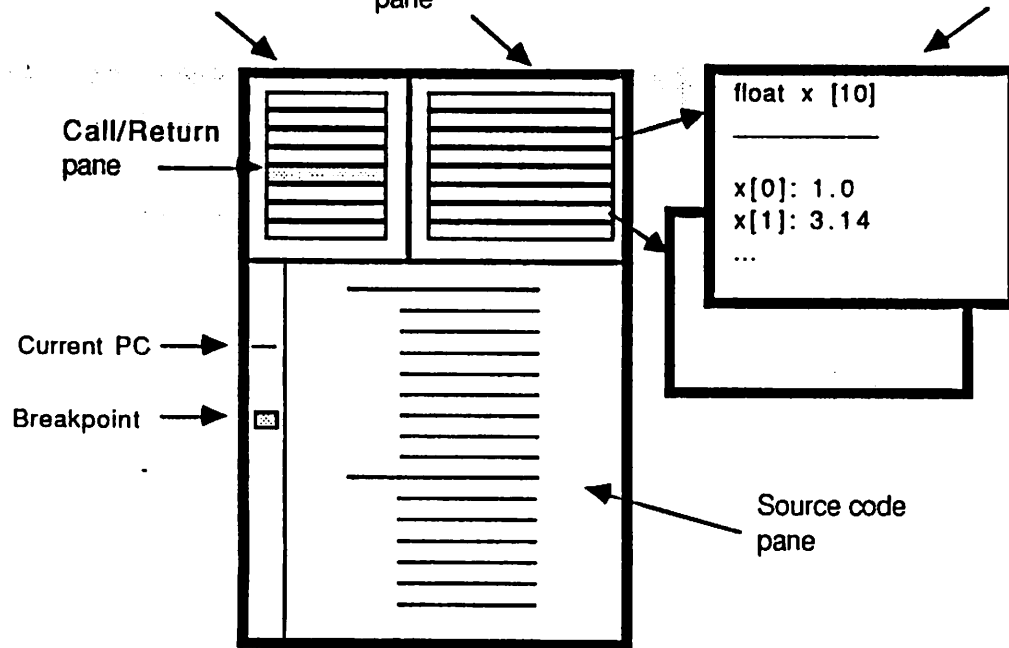
Main debugger window



Process window

Local variable pane

Variable windows



Monarch

Latency

The roundtrip time to memory and back

Performance is limited by operand fetches

- **The Switch Network is Very Good.**
C-MOS ~350ns GAs ~ 200ns
- **Commercial Processor**
Designed for **ultra-low latency** memory systems
Multiprocessor MIPS < Uniprocessor MIPS
- **Latency Masking**
Mechanisms for tolerating long
roundtrip times to memory and back



Monarch

Masking Latency

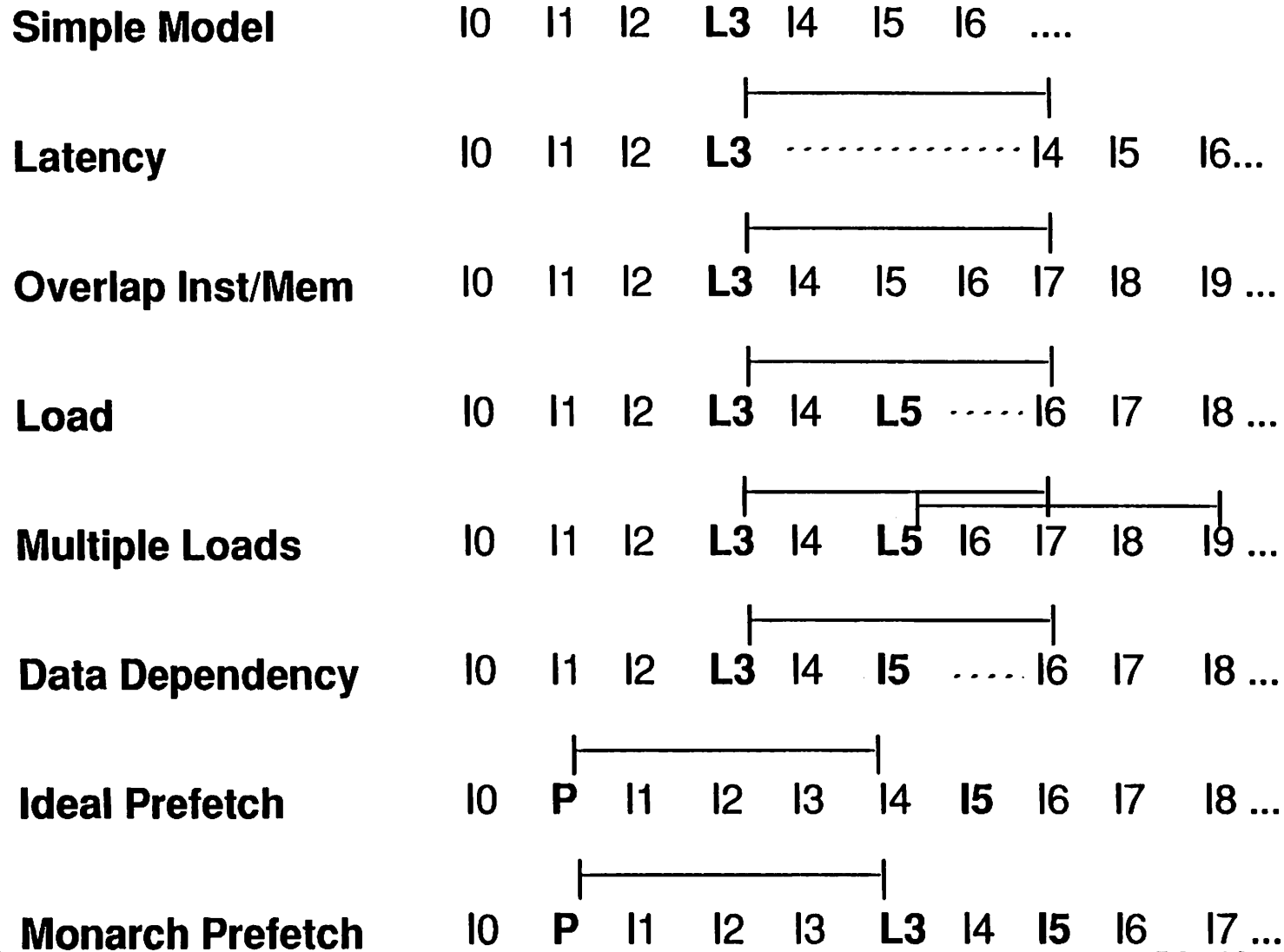
Use extra Bandwidth to regain processor performance.

- **Parallel Memory Access and Computation**
Insufficient natural instruction-level parallelism
 - **Instruction stream blocked by memory access**
Allow multiple outstanding memory accesses.
 - **Instruction stream blocked by data dependency**
Increase instruction level parallelism by generating earlier memory requests - **PREFETCH**



Monarch

Latency



Monarch

Benchmark Simulation

Problem	Problem Name	Cray T3E		Monarch versus Cray		Cray runtime	
		Runtime (microseconds)	processors	processors	processors	processors	processors
1	Small Matrix	2,270,000	14	220	764		
2	Very Sparse Large Matrix	6,960,000,000	43	682	2531		
3	Sparse Matrix	9,900,000	27	425	639		
5	Bit Permutation	94,300,000	323	5153	20500		
6	Modulo 2 Inverse	322,000	7	24	31		
7.1	Sort (Small)	19,500	25	115	130		
7.2	Sort (Small)	390,000,000	1191	18987	74713		
8	Binary Dynamic	36,600,000	55	107	220		
10	Dynamic Programming	1,000,000	7	103	345		
11	Pattern Recognition	966,000	55	743	2123		
12	Pattern Search w/free bits	1,280,000	27	405	1384		
13	Select and summation	466,000	24	370	1347		
14	Every Eleventh Bit	28,000	113	903	1217		
15	Decimate Binary Streams	65,000	31	419	1016		
16	Right Justify	62,000	21	328	1216		
17	Modulo 246 Inverse	291,000	49	766	2798		
18	Frequency Histogram	229,000	36	398	954		
19	Find 100 Zeros	10,000	10	80	122		
20	Transpose Blocks of Words	2,637,000	32	480	1884		
21	Table Lookup and Summation	12,500	56	278	313		
22	Fast Fourier Transform	6,500,000	18	274	1066		
23	Manhattan Distance	780,000,000	19	305	1209		
24	Bit Manipulation	140	6	11	13		
4.1	Frequency Search	11,400,000	735	8906	22800		
4.2		13,230,000	306	4147	8969		
4.3		500,000	407	1667	2000		
4.4		1,320,000	329	2260	2694		
9.1	Basket	363,000	40	81	79		
9.2		34,800,000	26	258	909		
9.3		1,590,000,000	143	2298	9191		
9.4		19,800,000,000	218	3474	13944		

Monarch

Monarch Development Plan

- **Three year R & D Plan:**
 - integrate commercial microprocessor
 - develop and test VLSI components
 - port Mach O/S and develop programming software
 - construct 128 processor demonstration machine
- **Additional Optional Design Study on High Density Packaging**



Monarch

Year 1 Hardware Tasks

Develop Pad Test Chip

- verify latchup protection pad layout
- verify programmable data delay circuit
- verify several process compensation circuits
- verify token recovery circuit
- verify power-up logic

Low cost, high speed package for Switch and Concentrator

Develop Test Board to exercise Pad Test Chip

Finish Development of Switch and Concentrator chips

Develop chip/commercial tester interface and test program

Develop Memory, Utility and Backplane PC boards

Develop Specs for Processor & Memory Controller Interfaces

Develop circuit level simulation



Monarch

Year 2 Hardware Tasks

Develop Processor Interface chip

Develop Memory Controller chip

Produce Memory Controller and Processor Carriers

Chassis, power supply, cooling for 128 Processor Monarch

Develop VME bus I/O adaptor board



Monarch

Year 3 Hardware Tasks

Complete Build of two 8 Processor Monarchs

Complete Build of one 128 Processor Monarch

Run the machines

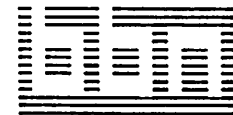
Complete Documentation

- Features

- Hardware Design

- Construction

- Operation & Diagnostics



Monarch

Year 1 Software Tasks

Sun MDS interface

Switch/Concentrator Diagnostic-Debug

Board Level Testing Software

Architectural Level Simulator & Benchmarks

Assembler

Linker

Prefetching C Compiler

Implement Multiprocessing & Synchronization Primitives

Design Machine Dependent Part of Operating System



Monarch

Year 2 Software Tasks

Operating System MDS

Use Simulator to Debug Software

Downloader/Boot Software

Logic Analyzer Based Software Development Station

Machine Dependent Layer of Operating System

Operating System Kernel Port

Essential Operating System Libraries

Operating System I/O



Monarch

Year 3 Software Tasks

Parallel Debugger(s)

Port other Unix Libraries

Port Essential Unix Tools

Identify Operating System Bottlenecks - Redesign and Implement for Massively Parallel Machine

Parallel File System

Documentation

Optional Additional Tasks

Machine Resource Allocation Strategies

Checkpointing

Port X-Windows

Port NFS File System

Parallelize and Improve selected tools



Monarch

Project Costs

Year	Hardware	Software	Mgmt	Total
1	\$1,075K	\$500K	\$150K	\$1,725K
2	\$2,300K	\$1,000K	\$200K	\$3,500K
3	\$1,800K	\$1,000K	\$200K	\$3,000K
Total	\$5,175K	\$2,500K	\$550K	\$8,225K

