INTEL/311849-001

# DTIC COPY Touchstone Project **Quarterly Technical Report**

034 AD-A224



Sig Lillevik

**Milestone Event Q4** 

**Intel Corporation** 

December 28, 1989

**Defense Advanced Research Projects Agency** Information Science and Technology Office

Linua Contractor

Contract No. MDA972-89-C-0034

Contractor: Intel Corporation

FACTUAL ACCURACY OR OPINION. Dimited rights are not subject to an appiration date he us

and dilciosure of technical data marked with this legend are set forth in the definition of "limited clause at \$252.277,7013 of the contractlisted above. and disciosure of feoverning right in paragraph (s)(15) of the

This levend, together with the indications of the portions of this data which are subject to limited rights, shall be included on any re production hereof which includes any part of the portions subject to such imitations. This technical data will remain subject to limit ed rights only so long as it remains "unpublished" as defined in paragraph (a) above.

#### DISTRIBUTION STATEMENT A

Approved for public releases Distribution Unlimited

07 16 249 90

CI EARED

DEN PUBLICATION

JUL 1 1 1990

NIRECTORATE FOR FREEDOM OF INFORMATION AND SECURITY REVIEW (DASB-PA)

DEPARTMENT OF DEFENSE

REVIEW OF THIS MATERIAL DOES NOT IMP

DEPARTMENT OF DEFENSE INDORSEMENT C

33

INTEL/311849-001

## **Touchstone Project Quarterly Technical Report**

**Milestone Event Q4** 

STATEMENT "A" per Karen Schroder DARPA Library, 1400 Wilson Blvd. Arlington, VA 1400 Wilson Blvd. Arlington, VA 22209-2308 TELECON 7/25/90

Prepared By

Sig Lillevik

Intel Corporation Intel Scientific Computers 15201 NW Greenbrier Parkway Beaverton, OR 97006

December 28, 1989

Sponsored by

Defense Advanced Research Projects Agency Information Science and Technology Office Research in Concurrent Computer Systems ARPA Order No. <u>6402. 6402-1</u>; Program Code No. <u>8E20 & 9E20</u> Issued by DARPA/CMO under Contract <u>#MDA972-89-C-0034</u>



#### DISCLAIMER

"The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Defense Advanced Research Projects Agency or the U.S. Government."

Accesion For			
NTIS	CRA&I	d	
DTIC	TAB		
Unann	ounced	0	
Justification			
Distribution / Availability Codes			
Distrib	vailability	Codes	
Distrib A Dist	vailability Avail a Spec	Codes	

VG

### Abstract

The focus of this quarter's technical activities has been demonstrating the DELTA numeric node and porting Mach to the IOTA I/O node. In addition, we determined our MRC strategy for the DELTA prototype, completed development of the GAMMA prototype, and finalized our DELTA system architecture.

Based on the i860 microprocessor, the DELTA numeric node offers preliminary 40 MHz performance of 8.76 double-precision MFLOP's on the 100 x 100 LINPACK test. The Mach port was successful but requires future work in the network message server and in improving message-passing performance. For the DELTA prototype, we will use the Caltech MRC in the standard MOSIS package but limit its channel bandwidth to 40 MB/s. Our GAMMA prototype proved highly successful and features up to 128 compute nodes and peak performance of 7.6 double precision GFLOP's. Finally, our expectations for the DELTA prototype are equally impressive as it will scale to 512 nodes in a mesh interconnection network and provide 30.7 peak double precision GFLOP's, up to 32 GBytes of distributed memory, and over 150 GBytes of off-line disk storage.

## Contents

1.	Summary 1
2.	Introduction
3.	Methods, Assumptions, and Procedures
	3.1 DELTA Numeric Node
	3.2 Mach Port
	3.3 MRC Strategy
	3.4 GAMMA Prototype
	3.5 DELTA Architecture
4.	Results and Discussion
	4.1 DELTA Numeric Node (RX-1)
	4.2 Mach Port
	4.3 MRC Strategy
	4.4 GAMMA Prototype 11
	4.5 DELTA Architecture
5.	Conclusions 13
6.	Recommendations15
7.	Distribution List

The technical data contained on this page is subject to the use and disclosure restrictions identified in the restrictive legend on the front cover of this report.

.

## 1. Summary

This quarter our major activities included development of the DELTA numeric node, porting Mach to the IOTA I/O node, determining our MRC strategy for the DELTA prototype, completing development of the GAMMA prototype, and finalizing our DELTA system architecture. We conducted these activities in support of our Q2, 1990 milestone demonstrating the DELTA prototype, our Q3, 1990 milestone comparing three operating system kernels, and our general development leading to the SIGMA prototype.

Each of these activities may be summarized as follows.

- The DELTA numeric node was designed using board-level simulation techniques for all control circuits and high-speed, 40 MHz printed circuit traces.
- A small IOTA prototype system consisting of two, integrated I/O nodes and two disks were used in porting Mach.
- A special surface-mount, *scrambler* board was developed to investigate the *cavity-down* package for the MRC.
- The GAMMA prototype was developed and integrated by porting NX/2 to the DELTA numeric node and operating it in a hypercube environment.
- Our DELTA system hardware design was determined by beginning with "off-theshelf" packaging components, identifying the internal subassemblies, reviewing alternate proposals for each, and preparing technical specifications.

The above activities produced several significant results. First, the DELTA numeric node, based on the i860 microprocessor, offers *preliminary* 40 MHz performance of 8.76 double-precision MFLOP's on the 100 x 100 LINPACK test. With the Mach port, the development process indicates that the Mach task/thread model requires some shared memory mechanism to extend threads across node boundaries. Further experiments indicate that the Caltech MRC in its standard MOSIS package will operate correctly at 40 MB/s so we will use this device in the DELTA prototype. Scalable to 128 nodes, the GAMMA prototype features the DELTA numeric node in a hypercube interconnect and provides 7.6 GFLOP's double and 10.2 GFLOP's single-precision (peak). Similarly, our expectation for the DELTA prototype include scalability to 512 nodes in a mesh interconnect, distributed memory of up to 32 GBytes, and 30.7 GFLOP's double and 40.9 GFLOP's single-precision (peak).

Finally, a number of enhancements or recommendations have been proposed for DELTA prototypes including: i) increasing the level of integration on the numeric node, ii) work with CMU to examine Mach on distributed memory, message passing parallel computers, iii) changing the MRC processing technology and/or packages, and iv) a number of possible hardware and software enhancements for the GAMMA and/or DELTA prototypes. Not all of these recommendations will prove appropriate but some may represent significant advances to scalable parallel computing. In the following quarters, we will be investigating these recommendations and determining their potential impact on the DELTA prototype.

## 2. Introduction

During the quarter just completed (Q4, December 28, 1989), Touchstone Project activities centered on completing the DELTA Numeric Node and porting the Mach operating system to an 80386-based I/O node in the IOTA prototype. In addition, we determined our MRC strategy for the DELTA prototype, completed development of the DELTA pre-prototype (denoted the GAMMA prototype), and finalized our DELTA system architecture design.

Our prototyping plan was to develop the DELTA Numeric Node, denoted the RX-1, with the interface to the system's interconnection network determined by a modular daughter board. Two daughter boards will be used: one for a hypercube interconnect (which exists for the iPSC/2) and one for a mesh. This allowed us to develop the RX-1 and then test it first in a hypercube interconnect (GAMMA prototype) and then later in a mesh interconnect (DELTA). Moreover, the GAMMA prototype provided a interconnection network to debug the system software and confirm the correct operation of the RX-1 in a system environment.

A similar prototyping strategy has been planned for the operating system kernel research. We will port three separate kernels to the DELTA prototype and then perform a study comparing strengths and weaknesses. For the Q2 Milestone we ported the first kernel, the *Reactive Kernel* from Caltech, and this quarter we have ported the second kernel *Mach* developed at Carnegie Mellon University (CMU). Both of these ports have used NX/2 as a stable base in the development process. Finally, we will port all three (RK, Mach, and NX/2) to the DELTA prototype and complete the study for our Q3, 1990 milestone. Each of these ports has allowed us to understand and learn about each kernel and these results will guide us in future kernel development.

For the past two quarters, we have reported on certain anomalies observed with the Caltech MRC. To resolve this issue in time for development of the DELTA prototype, we investigated an alternate packaging scheme and ran additional characterization tests. Results of these tests provided the information necessary to determine a strategy for use of the MRC in the DELTA prototype.

We completed development of the GAMMA prototype to further evaluate the DELTA Numeric Node's capabilities and performance in a system environment. This prototype has been very successful and two additional 128-node GAMMA systems were shipped this quarter to researchers at Oak Ridge National Laboratories and NASA's Ames Research Center. In addition, we sampled the techniques necessary to fabricate, assemble, build, and test the DELTA numeric node within a hypercube interconnection network. With these activities behind us, we are now ready to complete the design and construction of the DELTA prototype.

The technical data contained on this page is subject to the use and disclosure restrictions identified in the restrictive legend on the front cover of this report.

Our last major accomplishment for the quarter entailed finalizing the DELTA system architecture. Several design dimensions were discussed, evaluated, and selected for an overall optimum system. This information provided a set of specifications for the DELTA prototype and allowed many subsystem designs to begin.

The next section of this report discusses the methods, assumptions, and procedures used in completing the above accomplishments. Following this, a section reviews the results of the research and interprets the data. Finally, the last two sections, Conclusions and Recommendations, present the findings of the report, discuss their implications, and provide insight into further research issues and topics.

### 3. Methods, Assumptions, and Procedures

This section explains the experimental approaches used in researching the topics of this report. In addition, it describes the basic assumptions used as a foundation for this study.

### 3.1 DELTA Numeric Node

Several key decisions directed the design of the DELTA numeric node. For example, the RX-1 board form factor was chosen to be identical to the form factor of the iPSC/2 compute node. Similarly, the numeric node was required to use a daughter board network interface similar to the iPSC/2 Direct-Connect Module<sup>TM</sup>. Both of these decisions allowed us to leverage existing packaging and hardware by connecting a Direct-Connect Module (DCM) from the iPSC/2 compute node onto the RX-1 and debugging it in the GAMMA prototype. This approach allowed us to focus on the operation of the RX-1 node and not the interconnection network.

Our final decision concerned the number of prototype boards required for debugging the numeric node in a system (GAMMA) environment. Since a large number of prototype boards are required for system-level testing (approximately 150-175 for a 128-node system), the RX-1 node was designed just as carefully and completely as if it were developed for production.

First, an architectural design specification was written to provide hardware, diagnostic, and software engineers a common point of reference. The circuit was then designed and entered into a schematic database for review and generation of net and parts lists. All PAL's were simulated at the functional level as well as all critical PCB traces at the circuit level. Finally, the printed circuit board was designed and fabricated and initial boards assembled and debugged.

In parallel with the hardware design, firmware for the numeric node was specified, coded, simulated, and eventually debugged on hardware. Features of the firmware assisted the hardware debug. For example, a set of debugging routines were run that exercise specific sections of the numeric node. These loops allowed observations of the circuit signals with logic analyzers and oscilloscopes.

#### 3.2 Mach Port

To manage the port of Mach, we concluded (after discussions with CMU developers) that it was safer to perform the port on a stable release and not one that was still under development. Conversely, it was desirable to port a release that contained as many new and advanced features as possible. For these reasons, we elected to port the stable but advanced CMU-internal Release X95 (July, 1989) which is in between external Release 2.0 and the current external Release 2.5.

Mach was ported to an IOTA I/O node with an eight-channel Direct-Connect Module and not the normal single-channel DCM. This approach enables one I/O node to communicate with the other and the System Resource Manager (SRM) used to download programs into the nodes. In addition, these I/O nodes contained a SCSI interface to Winchester disks because the Mach port requires a disk drive on each node for a local file system and virtual memory swap space. No support was provided for tape drives or any mechanisms for diskless nodes (NFS, etc.).

The entire prototype system consisted of two IOTA I/O nodes connected to two 760 MB Winchester disks all driven by an Intel SYP-301 system resource manager (386-based workstation). A second SYP-301 was used to hold the Mach source code and to perform software builds. Mach binaries were then written onto cartridge tape and transferred from one SYP-301 to the other. Finally, a *light pen* and terminal were used to decode the optical messages generated by the firmware and emitted by the LED's on the I/O node board.

### 3.3 MRC Strategy

Following our MRC characterization activity, we postulated that the cause of the MRC failures at high speed and full load is that the standard MOSIS packaging of the die is not sufficient to meet the chip's transient current demands. Professor Seitz at Caltech then directed the Hewlett-Packard Company in Corvallis, OR to repackage the MRC die i) by double-bonding the die to the pins, and ii) by mounting the die in a *cavity-down* package which places the die upside down, directly onto a ceramic substrate that connects the pads to the pins. Last quarter (September, 1989) we investigated both of these new packages. Results of the double-bonding experiment indicated that this alternative did not solve the problem. In fact, the mis-routing problem increased. Because the pinout of the cavity-down package changed, we could not reliably operate the new package with the existing MRC test fixture.

Rather than design a new test fixture, we initially built a wirewrapped scrambler board that corrected for the change in the MRC pinout. Unfortunately, high-speed signals do not propagate well on wirewrapped boards and the cavity-down package performed poorly. At this time, we designed a small printed circuit board that contained surface mount pins on one side and surface mount sockets on the other. This new scrambler board used controlled-impedance traces and corrected for the pinout change. Finally, using the new scrambler board and existing MRC test fixture the cavity-down package was evaluated.

Armed with the performance of two alternate packages for the MRC, we evaluated the possible options and selected the most promising plan for our DELTA prototype MRC strategy.

#### 3.4 GAMMA Prototype

The GAMMA prototype allows the RX-1 node to be evaluated in a system environment. Although the GAMMA prototype is based on the hypercube interconnection network of the iPSC/2 system, we are confident that from a system point of view the numeric node will function equally well in the mesh interconnection network of the DELTA prototype.

Several activities were involved in exercising the numeric node in the GAMMA hypercube environment—some hardware and some software. First, the DCM used to implement the hypercube interconnection network was attached to the numeric node. Then, the diagnostic software was ported to the i860 architecture and the system hardware verified by correctly running the diagnostics. In parallel, NX/2 was ported to the i860 numeric node to provide message-passing support (buffering, memory management, etc.). This was followed by porting the Concurrent File System and some simple evaluation applications. Then, various system configurations were tested to determine how well the RX-1 node runs in a system environment.

#### 3.5 DELTA Architecture

One of the key decisions affecting the DELTA architecture has been the packaging of the system. Our goal has been to use as much existing or "off-the-shelf" packaging and, otherwise, industry-standard form factors whenever possible, in order to focus on the impact of multi-chassis connections on communication network performance and scalability. To accomplish this goal, we began by assuming the use of many of the IOTA packaging components such as cardcages, cabinets, and power supplies.

The technical data contained on this page is subject to the use and disclosure restrictions identified in the restrictive legend on the front cover of this report.

From that point, we were able to start defining what goes inside the cabinet plus how the cabinets and their subassemblies are interconnected, cooled, powered, and assembled. As previously discussed, we decided to use a form factor for the RX-1 node that was compatible with the IOTA node form factor so the basic cardcage parameters remained the same. Still, what remained undefined were how the MRC's would be connected to each other, the nodes, and from cardcage to cardcage.

Meetings were held to discuss these and other issues and to review and evaluate alternate proposals. As a result, several technical specifications were developed that describe the subsystems that comprise the DELTA prototype. When a unique specifications was approved, the design, development, and integration activities then began.

## 4. Results and Discussion

This section presents the results of the research and interprets the significance of each result.

### 4.1 DELTA Numeric Node (RX-1)

The following list describes the basic features of the DELTA numeric node.

- 40 MHz i860 microprocessor.
- 8 MBytes DRAM with parity, expandable to 64 MBytes.
- IOTA node form factor with support for either the DCM or mesh rout interface module.
- Bi-directional, memory-mapped FIFO interface to the routing module.
- 64 KByte EPROM for booting and debugging.
- 52-bit, i860-readable timer with 100 ns granularity.
- Memory-mapped control and status ports.
- Interrupt controller maskable by control register bits.
- User-defined expansion connector with i860 address, data, and control buses.

The DELTA numeric node is designed to operate reliably from 0 to 50 °C over a voltage range from 4.75 to 5.25 VDC. In addition, the Mean Time Between Failures (MTBF) has been calculated with (1 MB DRAMS and 8 MBytes memory) at 69.4 KHrs for soft errors, 54.7 KHrs for hard errors, and 30.6 KHrs for total errors. For power, the numeric node requires 5.6 A at 5 VDC or 28 Watts.

The i860 component operating at 40 MHz is rated at 60 peak double-precision and 80 peak single-precision MFLOP's and 33 peak MIP's. Preliminary performance of the RX-1 node at 40 MHz has been measured at 8.76 double-precision MFLOP's on the 100 x 100 LINPACK and a harmonic mean of 3.56 double-precision MFLOP's on the Livermore Loops evaluation suite. These performance numbers are only preliminary and will increase with further optimization of the system.

### 4.2 Mach Port

The goals of the Mach port were to implement the Mach operating system on a small IOTA prototype system and to successfully pass messages between two nodes via the Direct-Connect routing modules. Our major results are as follows:

- The tasks necessary to port Mach to the IOTA prototype included: bootstrapping, an interface to the memory management unit, and device drivers.
- Bootstrapping was accomplished using the techniques developed for NX/2 that involve the RS-422 serial line and Direct-Connect routing logic.
- The virtual memory logical page size of 4 KBytes was selected because it corresponds to the 386 physical page size and to the largest disk block that can be transferred in a single operation.
- Three device drivers were written: SCSI disk driver, Com/TTY driver, and DCM driver. These tasks represented a major part of the porting effort.
- Finally, a program was written and executed which successfully passes messages between two nodes,

Mach was successfully ported to the IOTA prototype; this will now serve as a base from which we can compare and contrast Mach with NX/2 and RK. Providing shared virtual memory was deemed necessary if we want to extend Mach tasks across nodes. The design of Mach virtual memory was found to facilitate the addition of this feature.

### 4.3 MRC Strategy

Once the new scrambler board was built, the cavity-down MRC was evaluated in the MRC test fixture at various temperatures, voltage margins, and message loads. Like the double-bonded package, the cavity-down failed to improve the message misroute problem at high speed and high message load.

Faced with these results, additional tests were run on the Caltech MRC (FMRC2.1) in the standard MOSIS 132-pin PGA to determine the minimum speed at which reliable message delivery still occurs. Delays were introduced into the signal handshaking path to reduce channel speeds. It was determined, using the available sample size, that the MRC operated reliably over full temperature and voltage range at a bandwidth of 40 MBytes-per-second.

Given this information we developed the following strategy for the DELTA prototype.

- 1. Use the Caltech MRC in the standard MOSIS package and pinout.
- 2. Limit the speed of the MRC to 40 MBytes-per-second.
- 3. Continue to work with Caltech to refine the MRC and package designs for higher speed operation.

Our earlier characterization of the 1.6 micron MRC indicated that potential performance could reach as high as 65 MBytes-per-second. In summary, our MRC strategy is to operate the device at about 60% of it's potential performance to assure system reliability.

### 4.4 GAMMA Prototype

The following list describes the basic features of the GAMMA system prototype.

- Scalable to 128 compute nodes and 127 I/O nodes in a hypercube interconnection network.
- Distributed memory of up to 1 GByte, expandable to 2 GBytes.
- Peak double-precision performance of 7.6 GFLOP's, and peak single-precision performance of 10.2 GFLOP's.
- Direct-Connect message routing with individual channel bandwidths of 2.8 MB/s.
- On-line disk storage exceeding 100 GBytes.
- Off-line 9-track reel-to-reel tape and 8-mm cartridge tape.
- NX/2 operating system kernel.
- Concurrent File System.
- FORTRAN and C languages.
- DECON debugging system

The GAMMA prototype has proved to be an excellent vehicle to evaluate the DELTA numeric node in a system environment. In fact, the GAMMA prototype represents a significant computer development and a very powerful machine in it own right. For this reason, we will be shipping two 128-node GAMMA prototypes to researchers this quarter. Since the i860-based numeric node has performed so well in the GAMMA prototype, we expect that it will perform even better in the mesh-based DELTA prototype with its high-speed interconnection network.

### 4.5 DELTA Architecture

The DELTA prototype will support a system environment somewhat different than the IOTA or GAMMA prototypes. Instead of accessing the prototype through a frontend computer, a user will access the DELTA prototype through a remote workstation across a local area network. In other words, the DELTA prototype will function as a compute server on a local area network. A separate system console will be used to boot the prototype and run the diagnostic programs.

We will retain much of the basic IOTA and GAMMA packaging but redesign the peripheral module to accommodate the recently introduced 3-1/2 inch disk drives with their excellent capacity of 325 MB formatted and quick 12.5 ms access time. The cabinet will support up to 64 nodes packaged 16 to a cardcage with a maximum of 4 cardcages-per-cabinet. A maximum DELTA system of eight cabinets will hold up to 512 nodes and 512 disks. Physically, the DELTA prototype will measure 21" wide, 37" deep, and 60" high for a total volume of around 27 cubic feet. A maximum 512-node system consisting of eight cabinets would measure 168" (or 14 feet) wide, 37" deep, and 60" tall.

Each cardcage assembly will contain a  $4 \times 4$  mesh slice, a 1000 W switching power supply, and a fan pack consisting of six 125 CFM fans. Attached to the backplane will be a large daughter board called the *routing plane* that holds 16 MRC's and connectors to expand the system from a simple  $4 \times 4$  to the maximum 16 x 32 or 512node prototype. Individual nodes will connect to the mesh network with a *mesh routing interface* that replaces the hypercube Direct-Connect modules used in the GAMMA prototype.

For input/output, the DELTA prototype will incorporate the 80386-based I/O node developed for the IOTA prototype. At least one of the I/O nodes will be attached to a bus interface adapter to create a VME bus and support a TCP/IP Ethernet protocol engine. This gateway node will connect the DELTA prototype to the user's remote workstation and local area network.

The NX/2 operating system, along with the Reactive Kernel and Mach, will be ported to the DELTA prototype to provide message-passing support (buffering, memory management, etc.). In addition, the Concurrent File System developed for the IOTA prototype will also be ported to the DELTA prototype for access to both on-line disk and off-line tape peripheral devices. Both the FORTRAN and C languages will be supported and integrated with the DECON debugger.

## 5. Conclusions

The activities reported in this document have included development of the DELTA numeric node, porting the Mach operating system kernel to an 80386-based I/O node in the IOTA prototype, determining our MRC strategy for the DELTA prototype, development of the GAMMA prototype, and finalizing our DELTA system architecture design. Our major results are as follows:

DELTA Numeric Node	High performance i860-based numeric node supporting 8 Mbytes DRAM (with parity) scalable to 64 Mbytes with 4 MB DRAM, bi-directional FIFO based interface to the interconnection network, and peak performance of 60 MFLOP's single- precision and 80 MFLOP's double-precision.
Mach Port	Porting process required adaptations for bootstrapping, interfacing to the memory management unit, and device drivers. Of these activities, the greatest effort involve writing the device drivers. A shared memory mechanism would enhance Mach message passing support.
MRC Strategy	Use the Caltech MRC in the standard MOSIS package and pinout, reduce the bandwidth to 40 MB/s, and continue to work with Caltech to refine the design for higher-speed bandwidth.
GAMMA Prototype	Scalable to 128 compute nodes and 127 I/O nodes, hypercube interconnection network, distributed memory of up to 2 GBytes, peak double-precision performance of 7.6 GFLOP's and single-precision performance of 10.2 GFLOP's, on-line disk storage exceeding 100 GBytes, NX/2 operating system kernel, Concurrent File System, FORTRAN and C languages, and DECON debugger.

#### **DELTA Architecture**

Scalable to 512 nodes, mesh interconnection network with 40 MB/s channels, distributed memory of up to 32 GBytes, peak double-precision performance of 30.7 GFLOP's and single-precision performance of 40.9 GFLOP's, on-line disk storage exceeding 150 GBytes, NX/2 operating system kernel, concurrent file system, FORTRAN and C languages, DECON debugger, and remote workstation user environment.

The DELTA numeric node represents a state-of-the art microprocessor design with outstanding stand-alone performance. Moreover, it illustrates the advantages of using commodity or stock microprocessors, memories, and interconnection network interface components in a numeric node design. As these prototype nodes transition to commercial products, the price-performance of such numeric nodes may far exceed designs using custom silicon and esoteric semiconductor technologies.

The port of the Mach operating system kernel to the IOTA I/O node indicates that a kernel designed to adapt to a shared memory multiprocessor architecture requires similar mechanisms in a loosely-coupled, distributed memory architecture. This shared memory mechanism facilitates the implementation of message passing in Mach.

Our MRC strategy for the DELTA prototype will allow investigation of mesh-based interconnection networks and the system issues associated with a mesh-based parallel computer. The fact that we will slow down the device by some 60% should not impact our study or the results obtained from the DELTA prototype. If faster MRC's are later developed with identical signal definitions and pinouts, then they may easily be used to replace the slower devices.

The GAMMA prototype provides an excellent example of how a large number of microprocessors can be combined to create a powerful and scalable supercomputer. If we extend this trend even further, then it becomes clear that tomorrow's supercomputers will consist of hundreds or thousands of microcomputer-based processing nodes. Yet these scalable machines will outperform the traditional supercomputers and at a very affordable price.

As we finalized the DELTA architecture, we developed the characteristics and features of a new, mesh-based parallel supercomputer prototype. Many of the attributes of the DELTA architecture are based on results of the IOTA and GAMMA prototypes and other roots such as the Intel iPSC/2. This system will allow us to evaluate the impact of a number of new and important technologies--both hardware and software.

## 6. Recommendations

As we continue our research, several of the experiments reported in this document may benefit from additional study or require us to pursue alternate design approaches. We are considering the following possibilities:

#### Improved DELTA Numeric Node (RX-2)

- Page mode RAM (PMRAM) cache for increased memory performance.
- An ASIC that would enhance the reliability, improve performance, and lower costs. Its major features would include a data pipeline, error detection and correction, integration of miscellaneous functions, PMRAM cache comparators, and faster performance monitoring port.
- A new stepping of the i860 device with clock frequency increased to 50 MHz which would offer a 25% performance improvement.
- Three functional changes to improve the interconnection network interface including: a wider FIFO data path, single-cycle FIFO transfers, and message start/stop assistance.
- A *Flash* electrically erasable EPROM that would allow the hardware to reprogram the boot and diagnostic PROM without physically removing the device from the board.
- Combining several of the node's state machines into one device for increased reliability, improved performance, and less cost.

#### Mach Port

- Implementation of a transport module for the DCM.
- Implementation of shared virtual memory.
- Examining ways to improve Mach message-passing performance.
- Eliminating the requirement of a local disk.
- Working cooperatively with CMU to investigate Mach on distributed memory, message passing, parallel computers.

The technical data contained on this page is subject to the use and disclosure restrictions identified in the restrictive legend on the front cover of this report.

#### MRC Strategy

- Consider alternate package styles with decreased lead inductance and improved power distribution.
- Increasing the speed of the device by changing to more advanced processing technologies with smaller feature size.
- Incorporating temperature-compensated, slew-rated limited input and output buffers for higher reliability and simplified system design.
- A *Slack Chip* to allow burst-like operations and to introduce differential signalling between backplanes.

#### **GAMMA Prototype**

- Port X-Windows client program to the IOTA I/O node.
- Install a VME-based frame buffer on an IOTA I/O node and support improved graphics capabilities.
- Investigate hosting the C and FORTRAN compilers on the concurrent file system.
- Implement redundant arrays of inexpensive disks (RAID) techniques.
- Develop a more powerful and easier to use debugger.

#### **DELTA Architecture**

- All the improvements to the GAMMA prototype described above.
- The use of battery-backed power supplies to increase file system data integrity and streamline file cache protocols.
- Incorporate an FDDI protocol engine with the IOTA I/O node as a gateway to higher-speed local area networks.
- The development of a new I/O node that provides a multi-slot VME interface with both master and slave operations plus hardware for support of RAID parity (exclusive-OR) generation.

As discussed above, many of the experiments reported in this document may be extended or augmented with several of these recommendations. Not all of them may prove appropriate or offer any clear advantage. They do represent a list of potential ideas, projects, and/or significant advances to parallel computing technology.

## 7. Distribution List

Delivery of this report has been made to the following:

DARPA/ISTO Attn: Stephen L. Squires 1400 Wilson Boulevard Arlington, VA 22209-2308 (one copy)

DARPA/RMO/Retrieval Services 1400 Wilson Boulevard Arlington, VA 22209-2308 (one copy)

Defense Technical Information Services Building 5, Cameron Station Attn: Selections Alexandria, VA 22304 (two copies)

DISCEMINATION (S Algher Dod