WORKSHOP NOTES

لم بر

10th Annual Workshop on Interconnections Within High Speed Digital Systems

9-12 May 1999

Hilton of Santa Fe Santa Fe, New Mexico

Sponsored by the IEEE Lasers & Electro-Optics Society and in cooperation with the IEEE Computer Society and the IEEE Communications Society

DISTRIBUTION STATEMENT A Approved for Public Release Distribution Unlimited

19991220 009

DTIC QUALITY INSPECTED 4

REPORT	DOCUMENTATION	PÅGE	- 00
Public reporting burden for this collection of gathering and maintaining the data needed, collection of information, including succession Davis Highway, Suite 1204, Arington, VA 22	Information is estimated to average 1 and completing and reviewing the colle inst for reducing this burden, to Washin 22-3302, and to the Office of Manageme	hour per response, including the t A locion of information. Send comma giton Headquarters Services, Direc ant and Budget, Paperwork Reductio	FRL-SR-BL-TR-99- 0 299
1. AGENCY USE ONLY (Leave Blank)	2. REPORT DATE 12 May 1999	3. REPORT TYPE AND . Technical	
TITLE AND SUBTITLE 10th Annual Workshop on Inter	connection within High Speed D	Ngital Systems	5. FUNDING NUMBERS G (2712E H832/01
AUTHORS Multiple			
PERFORMING ORGANIZATION IEEE-LEOS 445 Hoes Lane Piscataway, NJ 08855-1331	NAME(S) AND ADDRESS(ES)		8. PERFORMING ORGANIZATION REPORT NUMBER
SPONSORING / MONITORING / AFOSR 801 North Randolph Street, Roo Arlington, VA 22203-1977	NGENCY NAME(S) AND ADDRES	SS(ES)	10. SPONSORING / MONITORING AGENCY REPORT NUMBER F49620-99-1-0300
a distribution / availabilit	ystatement blic release, o	listributin walimited	12b. DISTRIBUTION CODE
a ABSTRACT (Maximum 200 wor 10th Annual Works	ds) hop on Interconnect	ions within High Sp	eed Digital Systems.
SUBJECT TERMS Interconnections, sy optelectronic, and o	stems architectures, ptical interconnecti	, electronic, ions	15. NUMBER OF PAGES 116 16. PRICE CODE
7. SECURITY CLASSIFICATION OF REPORT unclassified	18. SECURITY CLASSIFICAT OF THIS PAGE unclassified	ION 19. SECURITY CLASS OF ABSTRACT unclassified	SIFICATION 20. LIMITATION OF ABSTRACT UL
SN 7540-01-280-5500			Standard Form 298 (Rev. 2-89) Prescribed by ANSI Std. Z39-1 296-102

Program Committee

Workshop Chair Howard Davidson Sun Microsystems Sunnyvale, CA

Program Chair Philippe Marchand UCSD La Jolla, CA

Tutorials Chair George Papen University of Illinois Urbana, IL

Working Group Chair Ashley Saulsbury Sun Microsystems Sunnyvale, CA

International Liasons Peter DeDobbelaere Akzo Nobel Sunnyvale, CA

Henk Neefs University of Gent Gent, BELGIUM

Osamu Wada Fujitsu Laboratories Atsugi, JAPAN Technical Program Committee

Marc Christensen George Mason University, Fairfax, VA Kirk Giboney Hewlett Packard Laboratories, Palo Alto, CA Anthony Lentine Lucent Technologies Bell Laboratories, Holmdel, NJ John Levy Cisco, San Jose, CA Tulin Mangir TM Associates, Santa Monica, CA John Poulton University of North Carolina, Chappell Hill, NC Harold Stone

NEC Research Center, Princeton, IL

Working Group Committee

Lew Aronson Hewlett Packard Laboratories, Palo Alto, CA Giorgio Giaretta Lucent Technologies Bell Laboratories, Holmdel, NJ Charles Kuznia University of Southern California, Los Angeles, CA Rick Lytel Sun Microsystems, Sunnyvale, CA Henk Neefs University of Gent, Gent, BELGIUM Steve Tam Cisco, San Jose, CA

Workshop Scope

The continuing rapid increase in the performance of high speed electronics and communications technologies has led to dramatic improvements in advanced computing and communications systems. The rapid growth of computer internetworking and the rise of new applications such as multimedia and virtual reality are driving the requirements for still higher levels of computing and communications. Interconnections within digital computing and switching systems today are often perceived as a performance bottleneck. The purpose of this Workshop is to determine the interconnection requirements of emerging and future computer and communications systems and to disseminate information about state-of-the-art optical and electrical interconnection technologies at the component, packaging, and systems level.

Because of the multi-disciplinary nature of these problems, this Workshop brings together researchers and engineers with expertise in a variety of fields including electronic, optoelectronic, and optical interconnection technologies, advanced systems architectures as well as the systems level perspective of algorithms and applications. The Workshop is comprised of tutorials and invited talks of the highest caliber as well as a few contributed papers. In addition, all attendees participate in smaller working groups to discuss and address a central focus design problem. Working groups are diverse and multi-disciplinary. In the past, problems ranging from highperformance workstation design to tele-medicine applications have been considered. Historically, this workshop has provided a stimulating, highly interactive environment conducive to thoughtprovoking discussions. Take advantage of this opportunity to contribute to a great Santa Fe experience! More information can be found on the web at http://soliton.ucsd.edu/ihsds/santafe99

SUNDAY, 9 MAY 1999

3:00pm - 6:15pm Tutorial Session

Session Chair: George Papen, University of Illinois, Urbana, IL

3:00pm - 3:45pm Tutorial - 1

Device and Interconnect Technologies for ~100 GHz Mixed-Signal ICs, Mark Rodwell, UC Santa Barbara, Santa Barbara, CA 160 Gb/s TDM optical links will require ICs with > 150 GHz analog bandwidth and a 80 or 160 GHz clock. Mixed-signal ICs (DACs/DDSs/ADCs) for digital processing of 2-20 GHz radar signals will have 2000-transistor complexity and ~100 GHz clock rates. To permit clock rates exceeding 100 GHz, transistor current-gain (ft) and power-gain (fmax) cutoff frequencies must be several hundred GHz. The interconnects must have small capacitance per unit length, and wire lengths, hence transistor spacings, must be small. Given that fast transistors operate at high current densities, effective heatsinking is essential. To prevent circuit-circuit interaction through common-lead inductance ("ground bounce"), low wiring groundreturn inductance is required within the IC and between IC and package. We report a transferred substrate heterojunction bipolar transistor (HBT) IC technology providing scalable submicron HBTs with record 250 GHz ft and 820 GHz fmax. The interconnects, microstrip on a low-epsilon dielectric, have low capacitance and high velocity and a ground plane for low ground-return inductance. An electroplated Au/Ni/Cu metal substrate with Au thermal vias provides effective heatsinking. Demonstrated ICs include 85 GHz amplifiers and 60 GHz M/S latches. To manage power-delay products in larger circuits, low-voltage-swing (nkT/q) circuits are being investigated.

3:45pm - 4:30pm Tutorial - 2

Overview of Nonlinear Optics for High Speed Communication, Bahaa Saleh, Boston University, Boston, MA

4:30pm - 4:45pm Break

4:45pm - 5:30pm Tutorial - 3

Advances in Chip Level Packaging, John Carson, Irvine Sensors Corporation, Costa Mesa, CA

Two major directions in chip level packaging will be observed during the next decade: thinner packages (and therefore thinner chips) and more direct chip attach techniques. Package thickness will be pushed to as low as 0.5 mm for various applications enabled by agressive chip thinning techniques. In direct chip attach, peripheral leads in a footprint smaller than the IC carrier will appear in mainstream applications limited only by printed circuit board constraints. Combined, these two trends will drive toward increased use of three dimensional stacking techniques. Examples of thinned chips on flexible substrates and three dimensional assemblies of multi-chip packages are shown to portend these coming events.

5:30pm - 6:15pm Tutorial - 4

Modeling, Analysis and Simulation of Data Networks, Yusuf Ozturk, San Diego State University, San Diego, CA

The first topic will be more network modeling and analysis oriented. I can demonstrate some traffic collection tools and later incorporating the data collected into commercial simulation tools. We can work around practical problems for capacity planning and projections to the feature. This fits very good into a workshop program. This talk will reflect my experiences and common mistakes network managers and analysis specialists are doing during the data collection process, analysis and simulation of their network. This tutorial will be mostly a demonstration of network design process starting from data source characteristics, network topology selection, modeling and analysis.

6:30pm - 7:30pm Welcome Reception

8:00pm - 8:30pm Kickoff Meeting for Working Group Leaders Session Chair: Ashley Saulsbury, SUN Microsystems, Mountain View, CA

MONDAY, 10 MAY 1998

 8:00am - 8:15am
 Workshop Welcome

 Workshop Chair:
 Howard Davidson, Sun Microsystems, Mountain View, CA

 Session:
 Short Haul Interconnects

 Session Chair:
 Kirk Giboney, Hewlett Packard Laboratories

8:15am - 8:45am

1.1 Overview of 10Gbit Ethernet, Peter Wang, 3COM Technology Development Ctr. Santa Clara, CA

Internet traffic is exploding. Intranet, extranet, E-commerce and Voice-over-IP are all contributing to the growth of data networks. Gigabit Ethernet deployment is ramping up, as are broadband access networks. Carriers are planning for multi-gigabit backbone deployment. Dense Wavelength Division Multiplexing is the talk of the town. Are there alternatives? Is the world ready for 10 Gb/s networking?

This talk will explore the key enabling technologies and the various interconnect options for constructing 10 Gb/s links for the next generation backbone. We will also touch on the challenges of building switching infrastructure for the 10 Gb/s data networks.

8:45am - 9:15am

1.2 PAROLI a Synchronous Optical Interconnection Link with a Through Put of 13 Gbit/s, Karsten Droegemueller,

Siemens AG, GERMANY

Data communication and telecom switching systems require interconnections with high density, high data throughput and low power consumption. The design, realization, and characterization of a multichannel parallel optical interconnection with a 12 fiber ribbon and with an optical data rate of 1,25 Gbit/s per channel is reported. Two versions will be presented. First, a bit synchronous link with an electrical interface consisting of 22 differential data channels operating at 500 Mbit/s each plus one clock channel. Second, an asynchronous link with 12 electrical differential data channels at 1,25 Gbit/s each. On the transmitter side a vertical-cavity surface-emitting laser (VCSEL) array is employed as light source. Results of reliability test of the VCSEL's are given in the presentation

9:15am - 11:30am

 Session:
 Intra-System Interconnects

 Session Chair:
 Rick Lytel, Sun Microsystems, Mountain View, CA

9:15am - 9:45am

1.3 Tb/s Chip I/O – How Close are we to Practical Reality?, *Rick Walker, Hewlett Packard Laboratories, Palo Alto, CA* Computer and Router designers are counting on Tb/s chip-to-chip data transmission capability to continue expanding their system performances to meet the global demand.

Several prototype serial links, with clock and data recovery, have been published at 2-10 Gb/s data rates per pin. Much work is focussed on lowering the power and size of these links to allow hundreds of links to be integrated onto a single chip.

Even with these advances, some scary system issues still remain. Power supply noise can have disastrous effects on PLL and DLL performance. Signal crosstalk can close up an otherwise open eye. Each advance in CMOS scaling reduces the analog circuit options available to the link designer.

The copper signal-transmission infrastructure is not improving at anywhere near a "Moore's law" rate. FR4 has been the standard dielectric for highdensity PCBs for over 20 years, and coax cables are an extremely mature art. Dielectric and skin loss limit data rates to approximately 10Gb/s, and further advances may be slow coming.

This talk will explore these trends and attempt to forecast the future of high-speed serial interconnects.

9:45am - 10:00am Coffee Break

10:00am - 10:30am

1.4 Interconnect Requirements for Digital Cross-connect Systems, Roger Holmstrom and Robert Ward, Tellabs Operations Inc., Lisle, IL The requirements for high-speed and high-density board-to-board interconnects in digital cross-connect systems are such that new and emerging technologies are sought. These requirements are discussed in terms of physical, performance, reliability, and cost metrics. Some alternatives are evaluated. For long reach interconnects, parallel optics are favored. For short reach interconnects, electrical interconnects are chosen.

10:30am - 11:00am

1.5 Moore's Law: The Intra-system I/O Challenge, Craig Theorin, W. L. Gore & Associates, Lompoc, CA

The modularity of recent high speed digital system designs have created the need for intra-system I/O bandwidth in excess of 10 Gbit/sec. Most system architects anticipate this bandwidth requirement to scale with Mooreis law for the foreseeable future, creating a substantial signal integrity challenge for future data links. We will describe the chip-to-chip signal integrity concerns and likely solutions for intra-system I/O in the early part of the next millennium as aggregate bandwidths scale beyond 100 Gbit/sec.

11:00am - 11:30am

1.6 DDR and RAMBUS (High Speed Bus) DRAMD, Mian Quddus, Samsung, KOREA

11:00am - 11:30am Workshop Problem Statement for the 1998 Workshop:

Ashley Saulsbury, Sun Microsystems, Mountain View, CA

12:00pm - 1:30pm Luncheon & Working Group Session I

1:30pm - 3:30pm Working Group Session II

3:30pm - 6:30pm Free Afternoon

6:30pm - 7:30pm Reception

8:30pm - 9:30pm Special Event Speaker: 10 Years of Santa Fe Experience Harold Stone, NEC, Princeton, NJ

TUESDAY, 11 MAY 1998

8:15am - 9:45am

Session: Optical Interconnects for High-performance Computing Systems Session Chair: Harold Stone, NEC, Princeton, NJ

8:15am - 8:45am

2.1 Interconnects in Scalable, Distributed Mulitprocessor Systems, *Jeffrey Kuskin, Silicon Graphics, Inc, Mountain View, CA* Communication among processing nodes (that is, CPUs, memories, and I/O devices) is perhaps the key component in the design of a multiprocessor system. Traditionally, multiprocessors have been constructed by connecting a small number of processing nodes to a common, shared bus. The shared bus provides not only a mechanism for the processing nodes to communicate, but also allows all communication to be broadcast to all nodes on the bus.

The broadcast capability of a shared bus greatly simplifies the overall system design. Unfortunately, electrical and mechanical constraints severely limit the number of nodes that a single shared bus can support. For this reason, multiprocessors that scale to large numbers of processing nodes do away with the single shared bus and instead employ a distributed system design in which processing nodes are interconnected via a high-bandwidth, low-latency, switched routing fabric.

The use of a routing fabric overcomes the scalability limitations of a shared bus, but introduces a number of complications of its own. This talk will explore the use of high-performance interconnects in a distributed multiprocessor system. We begin with a short discussion of the basic distributed multiprocessor node architecture and interconnection fabric, and the difficulties that such an architecture creates. We then describe how these problems are solved in practice, with an emphasis on the role of the interconnection network. We conclude with some thoughts on the increasing importance of communication in multiprocessor system designs and the demands that will be placed on future multiprocessor interconnection networks.

8:45am - 9:15am

2.2 The Role of Optics in Balanced Computer System Design, *Mike Chastain, Hewlett-Packard Company, Richardson, TX* Computer system architects have been waiting and watching the development of parallel optics for five years or more, hoping that breakthroughs in the producibility, pack-aging, and resultant costs would finally make optical links cost competitive with their cop-per counterpoints. The inherent advantages of optical interconnects are well known. The inherent reduction in physical size of both connectors and cables, increased usable communication distance, and reduced susceptibility to EMI and EMC have always been appealing; but the costs have always forced designers to do it in copper just one more time.

In the last few years we have all witnessed the almost exponential climb in CPU clock rates, soon to break the gigahertz barrier, and we have seen virtually every performance feature ever implemented in the fastest supercomputers migrate to single chip CPUs. Along with these improvements has come an equally impressive increase in the data band-widths required to keep these CPUs in execution. Soon we may see single chip CPUs capable of consuming 10 Giga-bytes per second or more.

These enormous bandwidths are forcing CPU and ASIC designers alike to push every integrated circuit connection to the maximum frequency in order to maintain pin counts at manufacturable levels. Intelis recent switch to Rambus DRAM is a clear indication that all vendors, even PC vendors, are faced with this problem.

As we increase the frequency of these products we will test the limits of printed circuit technology. Skin effect losses are already a problem, and within a few years these loses will be replaced in priority by dielectric losses, perhaps limiting the usable connection dis-tance to a single backplane or PC planar. Cables with very good dielectrics will, of course, allow longer distance; but there will always be copper trace in the path. As frequencies increase, every 4 to 5 inches of copper trace will reduce the usable cable length by about three feet in addition to decreases due to the cable dielectric losses. We may do well to connect adjacent racks with copper cables, let alone cross machine rooms.

So it appears, within a few years, that copper interconnects of more than a few meters could easily become very expensive, while VCSEL technology and creative packaging may finally yield cost effective parallel optic interfaces.

The inevitable shift to parallel optical technology may occur because of this juncture of over stressed copper and mass produced optics, but there is still a major disconnect between the future needs of computer systems and the roadmap of optical components. Today the optical roadmap is driven by the telecommunications industry, which seems to be increasing communication frequencies in a fixed 14xî pattern from 622Mhz, to 2.5Ghz, to 10.0Ghz while the computer industry tends to take smaller 12xî increments. The pre-ferred frequency pattern for the computer industry must include 2.5Ghz and 5Ghz. These frequencies are especially important for computer / optic integration because many technologist now believe the useful limit of printed circuit technology is around 5Ghz. Beyond this frequency, i.e. at 10Ghz, it may be impossible to build a reasonable size backplane.

9:15am - 9:45am

2.3 In Pursuit of a Petaflop: Overcoming the Latency/Bandwidth Wall, Peter Kogge, Notre Dame University

The fastest machines on the planet today peak at around a teraflop $(10^{12}, floating point operations per second)$, with plans over the next few years to approach 10-30 TF. This performance, however, is still insufficient for several important classes of applications. Performance levels of a Petaflop $(10^{15}, flops)$ thus become a valuable target to aim for. Unfortunately, achieving this with current conventional technology and architecture seems to be difficult, and destined to wait for the 2010-2015 timeframe.

The twin demons in this wall appear to be latency and bandwidth: getting enough data to the right processing logic in a timely enough fashion that the logic can be kept profitably busy, and doing so in a fashion that the the amount of parallelism that must be found in an application is acceptable.

This talk will address one approach to breaking this wall: the HTMT project (Hybrid Technology MultiThreaded acrhitecture). This multi-institution collaboration is in the middle design phase of a long-term effort started in 1994 to find alternatives to conventional architectures and relevant technologies, and if successful, will result in a petaflops level machine by around 2006.

The solutions used by HTMT encompass both technology and architecture. In technology, superconducting logic, very fast WDM all optical networks, Processing-In-Memory (PIM), and 3D holographic storage form the basic underpinnings for a radically different machine. In architecture, multithreading, active memory, and automatic percolation of data throughout a very deep memory hierarchy all are central players.

This talk will overview the inherent problems associated with achieving a petaflops, and discuss the architecture of the current HTMT design. Although all aspects of the machine will be discussed, emphasis will be placed on the active memories, where PIM technology coupled with the concepts of percolation, allow massive parallelism in the memory system to execute large portions of an application in ways that defeat the bandwidth/latency barriers formed by conventional approaches.

9:45am - 10:15am

2.4 Ultra-High Speed Optical Interconnection Network for Supercomputing, Keren Bergman, Princeton University, Princeton, NJ In an attempt to effectively utilize the immense bandwidth of optical fiber interconnects, we designed a completely novel network architecture specifically for optical implementation. This work is part of an aggressive multidisciplinary architecture study of the next generation high performance computing based on hybrid technologies and multi-threaded (HTMT) latency management. The optical network employs multiple node levels with a routing topology that is based on a minimum logic at the node scheme. Our architecture features radically new traffic control logic, having the property that all routing decisions for the self-routing data packets, are based on a single logic operation at each node. The optical network, named the Data Vortex, can scale to interconnect an ultra-high performance computing system in a massively parallel form. Within the framework of the Data Vortex network we are investigating enabling fiber optic technologies and an implementation that consists of fiber interconnects with wavelength division multiplexed and time division multiplexed (WDM/TDM) payload and header. The development and incorporation into the network of fiber optic modules including high speed fiber lasers, amplifiers, and switching nodes will be discussed in this talk.

10:15am - 10:30am Coffee Break

10:30am - 12:00noon Session: Optical Networks Session Chair: Tulin Mangir, TM Associates, Santa Monica, CA

10:30am - 11:00am

2.5 Ultrafast Optical Interconnect Based on Routing by "Clockwork" in Regular Mesh Networks, David Cotter, British Telecom, UK, F. Chevalier, and D. Harle University of Strathclyde, UK

The effectiveness of multi-processor systems (such as future massive-capacity routers and servers) is critically dependent on the speed and efficiency of interconnection. Full connectivity is required with large message throughput and minimal delay. An option under consideration is an ultra-high speed multi-stage packet-switched network, using fixed-length packets at serial bit rates of 0.1-1 Tbit/s. The packets are routed through the network on optical pipes, with >>routing and digital header processing (such as destination address recognition) performed 'on the fly' in the optical domain.

A key requirement for high performance is that the routing mechanisms and processing at network nodes should be as simple as possible. Here we describe a new strategy for routing in regular mesh interconnection networks, based on a method of automatic global scheduled ('clockwork') switching in the optical domain. Using this strategy, the intermediate routing nodes are merely needed to perform an extraordinarily simple function ('for-me-or-not-for-me' header-address recognition), otherwise traffic is routed onwards automatically in the optical domain with absolutely no further intelligent action performed by the node. The throughput is comparable with conventional store-and-forward packet switching, yet the simplicity of this strategy makes it suitable for implementation in digital optical logic. The clockwork approach enables some special capabilities-such as ultra-low latency signalling, bandwidth reservation, ultra-low response delay, and process scheduling.

11:00am - 11:30am

2.6 Large-scale photonic packet switch using wavelength routing techniques, Koji Sasayama, NTT Network Innovation Laboratories, Kanagawa, JAPAN

This talk describes the large-scale photonic packet switching system being developed in NTT Laboratories. It uses wavelength-division-multiplexing (WDM) techniques to attack Tbit/s-class throughput. The architecture is a simple star with modular structure and effectively combines optical WDM techniques and electronic control circuits. Recent achievements in important key technologies leading to the realization of large-scale photonic packet switches based on the architecture are described. It is confirmed that a 320-Gbit/s system can tolerate the polarization and wavelength dependencies of optical devices. Experiments using rack-mounted prototypes demonstrate the feasibility of the architecture. The experiments showed stable system operation and high-speed WDM switching capability up to the total optical bandwidth of 12.8 nm, as well as successful 10-Gbit/s 4 x 4 broadcast-and-select and 2.5-Gbit/s 16 x 16 wavelength-routing switch operations.

11:30am - 12:00noon

2.7 Latency and Scaling Issues in High-Speed Optical TDM Networks, *Paul R. Prucnal, Princeton University, Princeton, NJ* An overview of optical TDM devices and techniques for ultra-high bit rate data communications is given as well as a discussion of the latency and scaling issues present in these systems. 12:00pm - 1:30pm Luncheon and Working Group Session III

1:30pm - 3:30pm Working Group Session IV

3:30pm - 7:00pm Free Afternoon

 7:00pm - 8:00pm

 Session:
 Optoelectronic and Optical Technologies

 Session Chair:
 Marc Christensen, George Mason University, Fairfax, VA

7:00pm - 7:30pm

2.8 The Commercial Applications of Optoelectronics, A View from the Optoelectronics Industry Development Association (OIDA), Arpad Bergh, OIDA, Washington, DC

The Optoelectronics Industry Development Association (OIDA) was formed in 1991 to advance the worldwide competitiveness of the North American optoelectonics industry and to promote the application of optoelectronics technology. The OE industry is a collection of six ore more distinct industries that all depend on OE technology. This fragmentation represents major challenges and opportunities.

It is difficult to draw a technology roadmap that serves all applications. On the other hand, there are great opportunities to share a common infrastructure that can advance a number of non-competing industries. Over the past eight years OIDA had carried out over thirty market survey and technology roadmap activities to identify emerging markets and shortcomings in domestic technology. Industry wide consensus was developed through informal interactions and recommendations were presented to industry and government for action.

The most prevalent impediments identified in these studies are the exploration of new markets for OE enabled applications and the ability to conduct high volume, low cost manufacturing. This talk will describe some of the initiatives that OIDA has undertaken to overcome these deficiencies.

7:30pm - 8:00pm

2.9 Board and Back-plane Level Optical Circuits Using Integrated Thin-cladding Polymer Fibers, Yao Li, Jan Popelek, and Jun Ai, NEC Research Institute, Princeton, NJ

This talk summarizes recent research activities at NEC Research Institute on optical interconnections using integrated polymer fibers. The objective of the research is to study large-bandwidth, short-distance, packageable optical solutions to address future interconnection needs at circuit board and back-plane levels. We have studied possibility of using embedded polymer fibers to form a 10 GHz board-level optical clock distribution circuit and demonstrated the feasibility of highly efficient and uniform delivery scheme for up to 128 optical termination's. Specialty thin-cladding polymer fiber bundles were integrated into convention PCB's. Various performance data will be presented. We also extended this embedding concept to include polymer fiber image guides (PFIG's), a cost-effective 2D image transmission components. We have fabricated some packaged and connectorized board-level optical circuits to perform point-to-point 2D parallel optical interconnects for future 2D vertical-cavity surface-emitting laser (VCSEL) and optical detector array based optical interconnects. Among demonstrated are some 16 node (6x6 bits/node) optical shuffle and butterfly interconnect circuits using three-layers of FIG embedding. Low insertion loss (< 2 dB) ad moderate resolution (11 lp/mm) were obtained. To further extend the capability of these 2D parallel optical circuits, we are experimenting a hybrid integration of these PFIG's and free-space micro-optic components so that branching and add/drop capabilities at different optical nodes can be incorporated.

8:00pm - 8:30pm

2.10 Development of Monolithically Integrated Transceivers for Single-and Multi-Channel Fiber-Based Optical Interconnects,

Clifton G. Fonstad, Jr., and Joseph F. Ahadian, Massachuestts Institute of Technology, Cambridge, MA

The Epitaxy-on-Electronics (EoE) optoelectronic integration technology, in which optoelectronic device heterostructures are grown epitaxially on fully processed GaAs MESFET electronic circuits, has produced uniquely complex monolithic OEICs combining optical emitters and detectors with high-speed VLSI electronics. This paper describes the development of transceivers for fiber-based optical interconnects using the EoE technology. Recent progress toward implementing the EoE technology with silicon CMOS electronics and more advanced GaAs technologies will also be reviewed.

WEDNESDAY, 12 MAY 1998

8:00am - 9:30am

 Session:
 Working Group Solution Presentations

 Session Chair:
 Ashley Saulsbury, Sun Microsystems, Mountain View. CA

9:30am - 9:45am Coffee Break

9:45am - 11:15am Session: Plenary Session Chair: Philippe Marchand, UCSD, La Jolla, CA

9:45am - 10:30am

3.1 Java, Jini and High Speed Systems of the Future, *Bill Joy, SUN Microsystems, Aspen, CO* Until roughly 1980, high performance systems were built of multiple boards, and organized around a disk operating system. Sun's Solaris and SPARC based products have this form, and applications like Oracle and SAP are focused around management of the information on the disk.

With the emergence of the internet in the last 20 years, systems are more and more often built on networking as the interconnect, often now with TCP/IP playing the role that the disk operating system did, providing the basic interconnect primitives. Sun's work with AOL and Netscape is an example of a major activity for us in this area, defining e-commerce as services relative to the interconnect.

In the future we believe that there will be a third organization for computing systems, those organized around objects and agents. We have built the Java and Jini technologies to support these new kinds of systems.

This talk will discuss these three organizing principles for computer systems (disks, internetworks and objects) and the implications for systems design.

11:15am - 11:30am Workshop Summary Matthew Goodman, Bellcore, Red Bank, NJ

Sunday, 9 May 1999

1999 IEEE Workshop Interconnections within High-Speed Digital Systems

- 3:00pm - 3:45pm Sun, 9 May - Tutorial 1

Device and Interconnect Technologies for ~100 GHz mixed-signal ICs

Mark Rodwell University of California, Santa Barbara

rodwell@ece.ucsb.edu 805-893-3244, 805-893-3262 fax

Device and Interconnect Technologies for ~100 GHz mixed-signal ICs

Two topics:

ICs ***for*** high-frequency interconnects RF/wireless, optical fiber ICs ***needing*** high-frequency interconnects 100 GHz digital logic, GHz ADCs/DACs

The organization:

what are the future applications ? what are the requirements ? what is the state of the art ? challenges for future high speed ICs ...and how my group is attacking them

Applications ICs for Non-nequency interconnects ICs needing* Inch frequency interconnects







÷ 1















































The wiring environment for 100 GHz ICs























Power-delay product in interconnect limit $P_{gate} T_{prop} = (1/2)C_{wire} V_{cc} \Delta V_{logic}$ bipolar logic (static power) $P_{gate} / f_{clock} = (1/2)C_{wire} V_{cc} \Delta V_{logic}$ CMOS logic (dynamic power) $(T_{prop} f_{clock})^{-1} \sim$ number gates between latches For fast, low-power logic: reduce $V_{cc} \Delta V_{logic}$







circuit results: transferredsubstrate technology















Fast ICs for fast interconnects Fast ICs needing fast interconnects ICs for GHz communications: Optical fiber transmission to, beyond 40 Gb/s with electronic data switching millimeter-wave (60/90/180 GHz) wireless networks at mm-wave, bandwidth is cheap & plentiful ...but the hardware must become cheap ADCs, DACs for digital processing of RF signals Challenges for fast ICs Fast transistors: scaling is key Wiring environment: signal, ground and power integrity Interconnect-limited power-delay products Managing high dissipated power densities

4:45pm - 5:30pm Sun, 9 May - Tutorial 3

ADVANCES IN CHIP LEVEL PACKAGING

IRVINE SENSORS CO

work - All rights

John. C. Carson

Irvine Sensors Corporation 3001 Redhill Avenue Bldg 3 Costa Mesa CA 92626

Irvine Sensors Corporation



Unpublished work - All rights reserved irvine Sensors Corporation

Excerpts from SIA I	Roadmap for Co	st/Performance (Category
	1995-2000	2000-2005	2005-2010
feature size (µm)	0.35-0.25	0.18-0.13	0.10-0.07
transistors/cm ²	4M-7M	13M- 25M	50M-120M
pin count	300-1000	1200-2000	2400-3600
package thickness (mm)	1.0-2.0	1.0	0.5 - 1.0
package cost (cents/pin)	1.4 - 4	1-2	0.6-1.3
package size (mm)	23-45	29-50	35-50
lead pitch- peripheral (mm)	0.3-1	0.3-0.65	0.3-0.5
lead pitch - array (mm)	1.0 - 1.5	1.0	0.5- 0.65
power (W)	2-18	2-28	2-55
Performance (MHz)	100-200	200-400	400-1000

Unpublished work - All rights reserved

Irvine Sensors Corporation

TRENDS OBSERVED

Packages are getting thinner

- Number of leads is getting larger
- Package footprint decreasing to approach chip size
- Direct Chip attach techniques are emerging
- Package pin pitch is decreasing
- Package pins are being distributed in array format

HOWEVER, SUPPORTING SUBSTRATES ARE BECOMING MORE AND MORE LIMITING

SYSTEM DESIGNERS ARE SEEKING SOLUTIONS IN MULTI-CHIP AND 3D PACKAGING TECHNIQUES ALONG WITH SYSTEM-IN-A-CHIP APPROACHES

Uspublished work - All rights reserved

Irvine Sensors Corporation

First Year of IC Production	1997 1999 2001 2003 2006 2009 2012
GA SOLDER BALL PITCH	
1.0 mm	
D.B mith	La state and the second s
0.65 mm	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
0.50 mm	and the second
INE PITCH BGA/CSP	
0.50 mm	
0.40 mm	
0.50 mm	
0.25 mm	
Annearth Required	Development Underway
This legend indicates the time during which	ch research, development, and qualification/pre-producton should be taking place for the solution.

Line all index much - All such a summariant

Irvine Sensors Corporation



. -

Year of First Product Shipment Technology Caneration	1997 250 mm	1999 360 run	2001 150 mm	2003 130 nm	2000 100 nm	2009 70 nm	2012 50 mm
Chip Interconnect Pitch (um)							
Wire bond-bell	70	50	50	50	50	50	60
Wire bandwedge	80	45	45	45	45	45	45
TAB*	50	50	50	50	50	50	50
Filp chip (area area)	250	180	150	190	100	70	50
TAB-rape automated bonding						~	
Pirst Year of IC Produ		1987 1999	2001 200		5 2009	2012	
WIRE BOND							1
70 ym							
C pm	1	- T					
50 yan	1	E					
45 will (wedge)		P					
TAB							
50 şuh							
FLIP CHIP	1						
250 =m	1		•				
140 µm		E I	-				
150 ym		<u> </u>					
130 µm		1	· · · · ·				
100 µm	1		100 C. C. S.				
70 em	1			<u></u>	سالىنىيى دەر		
50 pm							
THE REAL PROPERTY NAME		Ormiter	et Underwar			·····	

	Surface Mount	Chip on Flex	Flip-chip on Flex	
Feetprint	Large	Medium	Smail	
IC pad pitch (µ=)	150	150	100	
Cost factor	1.0	1.1	0.8-2.8	
Pretesting of chips	Yes	No	No	
Added wafer process	No	No	Yes	
Therm al Resistance	60 C / W	45 C/W	75 C/W	
Inductance (# H)	1.0-3.0	1.0 • 2.0	0.1 - 0.2	
Capacitance (pF)	0.2	0.2	0.03	
	 reauty available high reliability easy rework established in frastructure 	 good electrical performance compatible with most ICs excellent thermal resistance 	 good control of the second cont	
D isadvaz tagos	 Jow electrical performance high processing tem peratures requires flux cleaning largest outline 	 lower yields no pretesting no rework after glob-top 	poor therm al perform ance or flex wafer bum ping required no pretesting	

Unpublished work - All rights reserved

Irvine Sensors Corporation



Unpublished work - All rights reserved

Irvine Sensors Corporation

Lead count	Conventional TSOP	Fine Pitch QFP	CSF
28	152	56	11
32	168	63	19
52	169	81	31
64	169	99	36



Invine Sensors Corporation



ted work - All rights seenved Irvine Sensors Corporation









abod work - All rights reserved Irvine Sensors Corporation





Unpublished work - All rights reserve

Irvine Sensors Corporation

P	ACKAGING
Traditional IC Packaging	Wafer Level Packaging
Wafer is probed, diced and sorted	Wafer moved directly to packaging
ICs packaged away from fab	ICs packaged in fab
ICs are packaged one at a time	ICs are packaged en masse
Burn in performed in sockets	Burn in performed on wafer
Power and ground taken from PCB	Power and ground distributed in assembled structur
Device tested 2-3 times	Device tested once
High pin counts required	Lower external VO possible
Higher power required	Reduced power requirements
All function in the chip	Function shared between package and chip
More complex substrate required	Simpler substrates possible (lower VO)
Lead inductance concerns	Lead inductance nearly eliminated

Unpublished work - All rights reserved

Irvine Sensors Corporation

SUPPORTING TECHNOLOGIES FOR ADVANCED PACKAGING

Advanced Packaging requires the utilization of the following techniques extensively :

- thinning of silicon wafers containing circuits
- bump bonding for high I/O density interface
- handling of KGD in die form
- handling of die of different sizes and origins, non-electronic chips (e.g. MEMS, Lasers, Detectors, Fluidic Devices)

Therefore, advances in these techniques will help to increase the density and the functionality of advanced packages

Irvine Sensors Corporation

apabluhod work - All rights reserved

Intra Breach Convention ULTRA-THIN SILICON CIRCUITS A Kapton (50 micron thick) based flexible test vehicle has been used to test ultra-thin flash die 25 micron thin 16 Mb Flash die has been successfully tested after mounting on the test vehicle 25 Micron thin memory die mounted on the flexible substrate is bendable with the substrate. A bending radius of 1 mm can be obtained for each micron of silicon thickness

I have not a start of the start

Irvine Sensors Corporation



Irvine Sensors Corporation








NEO STACKING APPROACH

- Starting with KGD, construct a new, or neo-wafer with many dice in a molding compound matrix
- Use a standard neo-die size, just slightly larger than the largest die in the stack
- Add blank silicon to open areas on layers where smaller die are used to enhance thermal conduction between layers if needed
- Perform metalization and thinning in neo-wafer form
- Dice into individual layers
- Laminate into a stack

E SENSORS CORPORATI

- Neo-stacking is a breakthrough in high density packaging technology

 It allows complete systems in a cube
- It allows the combination of massive electronic functions with extreme miniaturization and integral logic and control functions
- dense layer-to-layer interconnects through the epoxy molding layer
- The process is highly manufacturable through industry standard

Irvine Sensors Corporation

automated tooling and batch processing

Unpublished work - All rights reserved

<page-header><section-header><section-header><section-header><complex-block><image><image>

	NEO STACK FA	BRICATIO	N EXAMPL	.E
0.4 ⁹⁸				
0,530	S.			
Aunther of	Layer Type	Routing Layers	Total Chips	Свир Турия
Humber of	Layer Type Cap Substrate	Rotting Layers	Total Chips	Chip Types
Standber of Layers 1 4	Luyer Type Cap Substrate Capacitor	Rowing Layers	Total Chips	Снир Турее
Humber of Layers 1 4 32	Layer Type Cap Subsettrate Cepacitor Flash	Rotting Layers	Total Chipe	Chip Types
6:50 Hundber of Layers 1 4 32 1	Low Type Cap Substrate Cepacitor Flash Plash Driver	Rotting Layers	Total Chips 	Chip Types
6:50 Humber of Layers 1 32 1	Layer Type Cap Subattime Capacitor Plash Plash Driver Microprocessor	Rotting Leven 2-eided 1 2 1	Total Chips 	Chip Types
	Leyer Type Cap Substrate Capacitor Plash Plash Driver Microprocessor FPCA	Rotiting Leyers 2-elded 1 2 1 2 1 2 1 2	Total Chips 	Chep Types
	Liver Type Cap Substrate Capacitor Flash Driver Microprocessor FPGA Bus Driver	Roming Layers	Total Chapa 	Chip Types
Aumber of Layers 4 32 1 1 1 1	Light Type Cap Substrate Capacitor Plash Plash Driver Microprocessor FRGA Bus Driver Boor Flash	Rotting Layers	Total Chips	Сыр Туры - - 1 - 1 - - - - - - - - - - - - -
Humber of Loyen 1 32 1 1 1 1 1 1 1 1 1	Low Type Cap Substrate Cap Substrate Capacitor Flash Driver Microprocessor FPGA Bus Driver Bus Driver Bus Driver Bus Driver Bus Driver Bus Driver	Rotting Layers 2-elded 1 2 1 2 2 2 2 2 2 2 2 2 2	Total Chipa 	Chip Types
Author of 1 4 32 1 1 1 1 1 1 1 1 1 4	Lover Type Cap Substrate Capacitor Fleah Fleah Driver Moroprocessor FPGA Bus Driver Boor Fleah HEE 1354 Interface DRAM	Rotting Layers 2-sided 1 2 1 2 2 2 1 2 1 1 1 1 1	Yotal Chips 	Chip Typen
Member of Layers 1 4 32 1 1 1 1 1 1 1 1	Low Type Cap Substrate Capacitor Flash Flash Driver Microprocessor FPSA Bus Driver Boor Flash EEE 194 Interface DRAM Bottom Ceremic	Rotting Layers 2-elded 1 2 2 2 2 2 2 2 1 2 1 2 1	Total Chips 	Chip Types
Munther of 1 4 32 1	Liver Type Cap Substrate Cap Substrate Capacitor Flash Driver Microprocessor FPGA Bab Driver Boof Flash HEE 1354 Interface DRAM Bottom Ceramic	Rotting Layers 2eided 1 2 2 2 2 1 -	Total Chips 1 1 4 1 4 1 1 2 3 1 52 Total	Chip Types



Unpublished work - All rights reserved

Invine Sensors Corporation







ed work - All rights wed

	SD-PACI	KAGING R	OADMAP	
YEARS	1998	2000	2002	2004
In-plane line density (lines/cm)	500	1000	1500	2000
In-plane total number of metalization layers	2	3	4	5
Side-face line density (lines/cm)	200	400	800	1000
Side face total number of layers	1	2	2	3
Areal line density (new technology) (lines/cm ²)	900	1600	2500	5000
Maximum Operating Frequency -coplanar lines (GHz)	1	2	4	8
Maximum Operating Frequency – Microstrip lines (GHz)	10	20	30	50

d work - All rights

PACKAGE INTERFACES							
100 mm² chip	Bare die with	Flip-chip	Quad Flat Pack	Micro Ball			
	75 μm Wire	0.5 mm bump	with 75 μ m	Grid array			
	bond		wire bond	1 mm bump			
pitch (mm)	0.15	0.25	0.30	0.50			
Footprint (mm ²)	125	125	785	150			
Package/chip area	1.25	1.25	7.85	1.5			
Height (mm)	0.4-0.6	0.5-0.7	1.4	0.84			
Inductance (nH)	1-2	0.05-0.2	1-7	0.5-2.1			
Capacitance (pF)	0.2	0.05- 0.1	0.5-1	0.05-0.2			

Uspublished work -	All rights reserved

Irvine Sensors Corporation

TECHNOLOGY	ADVANTAGES	DISADVANTAGES	APPLICATIONS
Flip-Chip	 Covers least area Has excellent electrical performance 	Lacks die availability Hard to assemble due to planarity Die shrinks results in board redesign	 Low lead count used (watches, vehicle modules, displays High reliability systems Vertically integrated companies
Chip-Scale Package	 System size reduction with standard technology 	 New technology lacks reliability and production infrastructure 	 Memories Portable computing and communications Under 100 leads
Multi Chip	 Early system integration Best electrical performance 	 Needs KGD Lacks die level availability Difficult test 	 Large systems (e.g. avionics) Some automotive and communication systems

lashed work - All rights meanwed

Irvine Sensors Corporation

From	То	Impact
Al pads and metallurgy	Copper pads and metallurgy	bond wires
		bump materials
		passivation
Wire bond	flip-chip	low cost wafer bumping
		underfill materials
leaded packages	area array packages	low cost dense substrates
		encapsulants
single chip packages	direct chip attach	low cost wafer bumping
	Chip scale package	low cost dense substrates
		low cost known good die
200 mm wafers	300 mm wafers	chip thinning
	very thin chips	material handling
		equipment configuration

Unpublished work - Ali rights received

Irvine Sensors Corporation

5:30pm - 6:15pm Sun, 9 May - Tutorial 4







USEFUL SNMP VARIABLES

- IfinOctets
 The total number of octets received on the interface, including framing characters.
- If a UcastPicts The number of subnetwork-unicast packets delivered to a higher-layer protocol.
- IfInNUessiPless
 The number of non-unicast (i.e., subnetwork-broadcast or subnetwork-multicast)
 packets delivered to a higher-layer protocol.
- packets delivered to a higher-layer protocol.

 ifOutOctets
- The total number of octets transmitted out of the interface, including framing characters. iOutUcanPicts
- The total number of packets that higher-level protocols requested be transmitted to a subnetwork-unicast address, including those that were discarded or not sent.
- ifOutNUcsstPkts The total number of packets that higher-level protocols requested be transmitted to a non-unicast (i.e., a subnetwork-broadcast or subnetwork-multicast) address, including those that were discarded or not not.
- iDescr
 Describes the interface. It should include identification information for the physical line and a description of the network.

TRAFFIC GENERATOR PARAMETERS for SIMULATIONS

The data collected must then be related, culled and summarized if they are to be useful.

Ultimately the traffic data must be consolidated to a single estimate of the source traffic destination for each node.

PacketIAT = SampleInterval/(ifInUcastPkts + ifInNUcastPkts)

MeanPacket Size = ifInOctets /(ifInUcas tPkts + ifInNUcast Pkts)

AverageDat aRate = ifInOctets × 8/SampleIn terval

SAMPLING FREQUENCY

- The sampling frequency determines the resolution of the data and
- its storage requirements. • Larger sampling intervals result in smoother data summaries and hide the variation between samples.
- Network performance statistics may have a periodic component. If so, the data sampling period should be less than half of that. The data period should not be divisible by the sampling period, otherwise the samples will consistently be taken at peaks, midpoints or low points of the data.
- When using SNMP to measure Internet/Intranet performance, do not let network management traffic swamp regular traffic. (Heisenberg Uncertainty Principle : We disturb the object we measure -- the more precisely we try to measure it the more we disturb the object)



TRAFFIC CHARACTERISTICS

- LAN Analyzers ,Expert Sniffer , RMON MIB , TCPDump will provide information about packet size distributions. When in doubt about the probability distribution function assume exponential.
- Packet inter-arrival times can also be obtained using LAN Analyzers or sniffers. You may use exponential distribution when you do not have a better model for the inter-arrival. Exponential distribution will not fit to all data. For example : RIP (routing information protocol) sends its routing tables on all interfaces every 30 seconds. While the client requests may be exponentially distributed the responses may be fixed sized closely spaced packets.
- Business cycle defines how the average packet rate fluctuates and the peak value should be used for performance analysis.
- The average packet arrival rate and the average packet size should be multiplied for link bandwidth selection.





BENCHMARKING

- Before a new application is deployed on a network benchmark data should be collected to predict the impact of the new application on the network resources.
- Questions :
- Is enough data being collected to provide a statistically valid sample?
- Will the real network actually experience the type of the traffic being measured?
- While collecting benchmark traffic data for an application the application should run on an isolated segment. Otherwise contamination will occur by the data not participating in the benchmarking study.
- Measurements taken on one type of network (Ethernet for example) does not apply to other networks (Token ring).

BENCHMARKING

Sources of contamination

- Target and measurement media are different. (Differences such as MTU, framing characters etc. should be compensated)
- Network operating systems are dissimilar. • Applications are similar but not identical.
- · Queues existed at the servers and network links at the time of measurement
- Disk I/O delays are lumped in with CPU instruction delays.
- The LAN carried other traffic not related to benchmarking.
- Measurements were taken remote from the the server system. • The server is doing other work in addition to our application.
- The benchmark network and systems are already highly utilized.
- · Even if the application will be deployed over a wide area network with a number of remote clients, benchmarking data still should be collected over a LAN system with ONE client and server on the same LAN.

NETWORK DESIGN PROCESS



DESIGN PARAMETERS

- What type of traffic will be carried ?
- Is the data to be carried time sensitive? If yes, what are the delay and delay jitter requirements?
- Is the data bursty in nature ? Will smoothing affect application ?
- · What are the acceptable bit error rates and packet loss rates?
- Is the traffic symmetric (such as in videoconferencing) ?
- · Are there any applications that will benefit from multicasting or broadcasting (such as digital video broadcasting) ?
- What are the reliability requirements ? What are the tolerable times for recovery in case of failures?
- What are the security requirements ?
- What protocols will be supported ?
- Total network budget and the percentage reserved for network management, analysis and data collection tools.
- Self similarity of network sites for reducing operating costs versus initial deployment costs.

WAN DESIGN

• Type of service decisions

- Frame Relay
- ATM
- DDS Lines
- . SMDS
- · Private versus public leased lines
- · etc
- Topology decisions
 - Star
 - Backbone
 - Tree
 - Mesh
- Link and Node selection and sizing

LAN DESIGN

- Topology decision (Bus, Ring, Tree , Star etc.)
- Access Protocol (CSMA/CD, Token Ring, etc.)
- Frame/ Packet Size
- Transmission Capacity
- Signal propagation delay
- Buffer size
- Processing delays
- Throughput
- User traffic profile
- Data collision and retransmission
- Bridging/Routing decisions
- Security
- Availability

Video Network Design

- Design a network to support video conferencing between three remote locations (Santa Fe, San Diego and Mexico). Each of these locations one or more active participants introducing compressed video at a rate of 384 kbps into the network (Using h.261).
- There may be more than one passive participant (Listening only -destination to the data) at each site.
- Employ multicasting to reduce the duplicate traffic on LAN and WAN links where possible.
- Simulate the proposed model and provide performance measures on the LAN and WAN Links.
- Based on the simulations refine the design (Topology, link and nodes, WAN service type etc).



PERFORMANCE METRICS

- Link Utilization
 - Utilization by application
- Utilization by protocol
 Collision Statistics (CSMA/CD)
- Token Ring Statistics
- Message Delay (maximum , minimum, min)
- Delay Variation
- Node utilization (Router Bridge Switch etc.)
- Processing delays
- Throughput

- Inrougnput
 Buffer size (by node and by port)
 Resistance to node and link failures (Availability)
 Scalability (How easy it will be to increase number of sites or number of users at each site)









NOTES	

Monday, 10 May 1999

8:15am - 8:45am Mon, 10 May - 1.1







Switched Ports: Cumulative 2002: 30 M GbE ports 1000,000 100,0



- Digital multimedia production & distribution - graphic/image rich business presentations
 - digital post-production
 - medical/scientific imagery
- Data mining & database Synchronization
- The video factor - Tele-presence

 - --- When will HDTV move onto the data network?







Distance Assumptions

- No building will move because of 10 GE
- Customers will expect to run 10 GE over the same links as GE
- Customers will expect to use 10 GE with new applications that are emerging or on the horizon
 - Access to RAN/MAN
 - cluster computing interconnects

Con Desired Requirements

- Distance - data center, cluster interconnect: 50 m
 - building risers: 500 m
- campus backbone: 2-10 km
- MAN access: 10-30 km
- MAN/RAN backbone: 30+ km
- Media
- single mode fiber except for shortest reach Topology
- pt-pt, full-duplex, w/ congestion control











- Vertical Cavity Lasers (VCSELs)



Distributed-feedback Laser (DFB)

- Distributed resonators suppresses multi-modes
- Buried heterostructure:low threshold current (~10mA) and high output power
- Thermal cooling » threshold current is temperature dependent » control mode hopping
- isolator
- » eliminate reflection noise
- External modulator
 reduce frequency chilping
 Eliminate cooler, isolator and/or external modulator
 for low-cost shorter reach solutions?







Component	Technology	Availability	hate Rate Capabil	Distant
Loos	EML.	Now	10 Gbps	80 km
	CWILMLINEOS	Now	10 Gbps	600 kr
	Uncooled FP	Jun-99	12.5 Gbps	1 km
	Uncooled DFB	Jun-99	12.5 Gbps	10 km
Photodelector	PIN	Now	12.5 Gbps	NA
	APO	Now	12.5 Gbps	NA
Laser Driver	GeAs	Now	10 Gbps	NKA
	SIGe	Dec-99	12.5 Gbps	NA
TIA	GeAs	Now	12.5 Gbps	NA
	SiGe	Dec-99	12.5 Gbps	NKA
Limiting Amp	GeAs	Now	12.5 Gbps	N/A
	SIGe	Dec-99	12.5 Gbps	NKA
COR	Gele	Now	12.5 Gbps	NA
	SiGe	Dec-99	12.5 Gbps	NKA.
Mus/Demax	GeAs or Bipoler	Now	10 Gbps	NA
	SGe	Dec-99	12.5 Gbps	N/A













8:45am - 9:15am Mon, 10 May - 1.2









	P				
Why Optical Interconnect ?					
Copper: BW x distance produ	ct limited				
BW demand drastically increa	asing				
Cable size is 1/50th of copper	- T				
Optics is the solution to escal	ating EMI problems				
Cost is not much higher than	high performance copper				
PAROLI	K. Drögernüller				
araliel Octical Link	Ber insernation				



(Here)				P				
Appli	Applications of Optical Interconnect							
Applications	Distance	Through Put	DC Couple	Status				
Box-Box	5 - 100 m	10 - 30 Gbit/s	Yes and No	Products available				
Backplane Extend	1 - 5 m	50 Gbit/s	Yes	R&D				
Backplane	25-100 cm	50 - 150 Gbit/s	Yes	R&D				
On Board	5 - 25 cm	50 - 100 Gbit/s	Yes	R&D				
PAROLI Parallel Optical L	ink		K. C)rōgemiller ann				























9:15am - 9:45am Mon, 10 May - 1.3

Tb/s Chip I/O - how close are we to practical reality?

Rick Walker Hewlett-Packard Company Palo Alto, California walker@opus.hpl.hp.com

Agenda

- Applications and Key Specifications
- General Architecture for inter-chip communication
- Limitations
 - Skin-Loss
 - Delay Matching for Multi-phase sampling
 - CMOS Scaling
- Industry Trends
- Conclusions

Current Practice

- Current high-performance systems are skew limited using parallel data clocked at 250-500Mb/s.
- Using clock and data recovery on Gb/s links eliminates the skew problem and improves system BW by factor of 8-16X.
 - What are the limits for advanced systems?

CPU-CPU/Memory Application





Key Specifications

- Speed: As high as possible at least 1Tb/s I/O per chip
- Latency: critical less than 10ns plus time of flight
- BW/link: limited to 4-5 Gb/s by PCB loss
- Power: for a 100W chip, all 250 links should dissipate less than 40W -> 160 mW per link
- Size: a typical processor may be 9cm², if links use 20% of the total area, then each 4Gb/s link cell should be less than 720000um² in size.



Skin Loss and Dielectric Loss

Nearly all cables are well modeled by a product of Skin Loss

 $S(f) = e^{-k_s(1+j)l\sqrt{f}}$, and Dielectric Loss $D(f) = e^{-k_d lf}$ with appropriate k_s,k_d factors. Dielectric Loss dominates in the multi-GHz range. Both plot as straight lines on log(dB) vs log(f) graph.







6dB Equalized Data











Example Multiphase RX Block Diagram



Measurement of a Multi-phase System



Reported Jitter: 8ps rms, 44ps pk-pk at 3.5Gb/s. Measurement of photo shows 26ps difference between widest and narrowest eye, so true eye margin for endend system is $44ps\sqrt{2} + 2 \cdot 26ps = 118ps$, or a total eye closure of 41%.

Attention to delay matching is critical!

Techniques to Improve Delay Matching and Power Supply Noise Immunity



CMOS Scaling Issues

· Gate delay no longer scales with process



See: Chenming Hu, "Low-Voltage CMOS Device Scaling" 1994 ISSCC Digest, pp 86-87.

CMOS Scaling Issues (continued)

- V_t doesn't track with power supply so we gradually lose ability to make ECL-like differential circuits.
- Full-swing circuits show worse delay matching than ECL-like topologies.
- Full-swing circuits show worse power-supply delay modulation than differential circuits.
- V_t matching gets worse due to statistical dopant variations in channel.
 - All of these trends make power supply noise rejection and multi-phase alignment more difficult with each process scaling.



Industry Trends

- 50% of all U.S. Families now have home computers
- Computer performance has surpassed needs of most users: witness the drop of P.C. prices in the last 3 years from a stable \$2K down to \$500 levels.
- Internet host count was doubling every 6 months in 1988, is now doubling every 24 months - we are clearly past the 50% adoption point in the growth curve.
- What applications will continue to drive expensive and exotic improvements in interconnect technology?
 - Without a new "killer app" to drive development, we may by stuck with the limitations of FR4/CMOS for quite some time.

Viability of "exotic" technologies

- Yielded CMOS parts come in at \$10/cm²
- Tb/s chip-chip links are probably feasible in the next few years.
- This performance can be achieved with existing BGA packages across commodity FR-4 PC Backplanes.
- The incremental cost of a Tb/s link in these applications will be about \$18 + connector cost.
 - For optical solutions to take hold in these applications, they must provide either significantly higher performance (>10Tb/s) or cheaper system cost (not likely).

Conclusions

- Still much work to be done, but 1 Tb/s chip I/O seems an attainable target.
- 5Gb/s on 1 meter PCB is the fastest that can be feasibly supported for the foreseeable future with *low latency*.
- Fiber seems to be progressing along either a 1-10-100-1000-10,000MHz or a 622-2488-10,000MHz evolutionary path. There may be an economically important need for 5Gb/s links.
- 10 Tb/s chip I/O is probably out of the question for current high-volume technologies (CMOS, FR-4 PCB). Computer designs and programs may have to give up cache coherency, and move towards cooperative computing architectures to break out of this limitation.

10:00am - 10:30am Mon, 10 May - 1.4





Digital	Rate	Equiv. Voice Channels	Optical	Rate
DS0	64 kbps	1	00-1	51.84 Mbps
DS1 (T1)	1.544 Mbps	24	OC-3	155.52 Mbp
RTD (Digital Radio)	1.544 Mbps		OC-12	622.08 Mbp
DS1C	3.152 Mbps	48	OC-24	1.244 Gbps
DS3 (T3)	44.763 Mbps	672	OC-48	2.488 Gbps
DS3C	89.472 Mbps	1344	OC-192	9.953 Gbps
DS4	274.176 Mbps	4032	OC-768	39.813 Gbp
DS5	470 Mbps			























10:30am - 11:00am Mon, 10 May - 1.5

Moore's Law: The Intra-System I/O Challenge

Craig Theorin May 10, 1999

© W.L. Gore & Associates Inc., 1999

Overview

- Introduction: The implications of Moore
 + Amdahl.
- Link Architecture Options
- Copper Media Scalability
- Fiber Optics Scalability
- Conclusion



Moore + Amdahl = Bandwidth Growth



- Moore Observes Exponential MIPs Growth.
- Amdahl Necessitates Proportional BW Growth to leverage Moore.





Amdahl's Law: Processor, I/O, Memory Balance

System performance is optimized when MIPs=Mbit/sec=Mbytes If processors scale with Moore's Law, so must I/O and Memory



Intra-System Data Bandwidth Trends



Keeping Pace with Moore

	Bandwidth	Bit Width	Rise Time	Spectrum	Zo Discont	Ch-Ch Sk
Year	(Gb/s)	(psec)	(psec)	(GHz)	A.U.	(psec)
1999	1.0	1000	250	1.4	1.00	250
2000	1.6	630	157	2.2	0.63	157
2001	2.5	397	99	3.5	0.40	99
2002	4.0	250	63	5.6	0.25	63
2003	6.3	157	39	8.9	0.16	39
2004	10.1	99	25	14.1	0.10	25
2005	16.0	63	16	22.4	0.06	16
2006	25.4	39	10	35.6	0.04	10
2007	40.3	25	6	56.4	0.02	6
2008	64.0	16	- 4	89.6	0.02	4
2009	101.6	10	2	142.2	0.01	2


Serial vs. Parallel Data Streams

- Serialization converts media cost to launch cost.
 - long reach applications.
- Parallel = N*Serial.
 Maximum I/O BW (ie. 10-100 X serial) & I/O BW*Density.
 BW Reduction for Serial
 - BW Reduction for Serial Stream Processing.

Parallel will be continually obsoleted by increasing SerDes performance/cost ratio.

A decline in serialization performance/cost growth may require parallel.



SerDes Cost/Performance



Ganged Serial vs. Clock Forwarding

- Jitter budget has limited 1 Gbps <u>optical</u> clock forwarding.
 - Clock and Data Jitter accumulate in budget.
 - HiPPI Budget TBD.
 - FO centric designs will end clock forwarding.
- Future High BW Links will look like hybrid of parallel and serial or "ganged serial".
 - Allow deskew of parallel data streams.
 - Scalable for future systems.



© W.L. Gore & Associates Inc., 1995

Data Coding

- Typical Code Functions:
 - Limit low frequency content (ie. run length) for AC coupling.
 - "DC-Balance" the signal to keep duty cycle close to 50%.
- Scrambling: Muxing a PRBS w/ Data
 Statistical max run length and DC balance.
- Multi-Level Coding
 - Lowers max frequency by increasing # of bits/symbol.
- Forward Error Correction (FEC)
 - Using Error Correction Bits to improve BER



Copper Scalability

SI Concern	Cause	Solutions
Loss, Eye Patter	n Skin Effect & Loss Tan	Larger Cable, EQ, EOP, Peaking
EMI	Poor Shielding vs Tr, Imbalance	e:More Shielding, RF Chokes,
Return Loss	Zo Discontinuity	Control, Signal Shielding
Next/Fext	Poor Shielding vs. Tr	More Shielding between Signals
Skew (pair2pair)	Signal Routing, Er Variation	Control, Deskew Circuits.
Imbalance	EMI, Jitter,	Control, More Shielding?

- A common solution is to increase Tr
- "Control", Larger Cable, More Shielding = Cost



Standard Cable Standard Cable Sto Mbps 20 Meters 100 Obm 26 AWG Store Internet Internet

Gb Ethernet Copper Modems

- 1 Gbps data transmitted over 4 pairs
 - 5 Level Code (PAM5) used to send 2 bits/symbol
 - Extra bits for forward error correction (FEC)
- Hybrids are used to bi-directionally couple data into cable.
- Waveform "shaping" for EMI compliance
- Noise Reduction Through DSP
 - NEXT and FEXT
 - Digital Echo (reflection)



© W.L. Gore & Associates Inc., 1999

Optics Scalability

SI Concerns	Cause	Solutions
Optical Eye	Optical Noise from Source	Improved VCSELs and Launch
BER vs. BW	Higher BW Challenges Budget	Longer Wavelengths
Jitter Budget	Source Jitter & Rx BW & Noise	Improved VCSELs and Receiver
Clock/Data Jitter	Sum of clock+data Jitter	Control, Ganged Serial X-mission

- 10 GbE is currently addressing this issue for Serial Optics to 10 Gbps.
- For greater BW Parallel FO is required.
- At a cost.



VCSEL Scaling





Increasing Beam Divergence



IEC Class 1 Allowable Output Power



Summary

- Moore & Amdahl Require 10 Gbps in 5 years and 100 Gbps in 10 years.
- Leverage IC Functionality to solve analog transmission problems.
 - DSP?
 - Encoding, Error Recovery, Peaking, Deskew, etc.
 - Copper EMI challenges will be extreme.
- Parallel Optic data links will address BW needs.
 - Ganged Serial Likely for future Intra-System I/O



© W.L. Gore & Associates Inc., 1999

	NO	125	
		······································	
••••			
	<u>.</u>		
			
		<u> </u>	· · · · · · · · · · · · · · · · · · ·
			·
<u></u>			

•

Tuesday, 11 May 1999

8:15am - 8:45am Tues, 11 May - 2.1





- Cache coherence in action
- Origin 2000 network details
- · Multiprocessor interconnects in the future

sgi









- Defines the convention used to communicate among nodes
- Message passing
 - · Each node has direct access only to its local memory
 - · Communication between nodes is requested explicitly
 - Examples: Intel Paragon, Thinking Machines CM-5, IBM SP2
- Shared memory
 - Physically separate memories appear as a single, unified memory
 - · Each node may access any memory location using normal loads/stores
 - Examples: HP/Convex Exemplar, SGI Origin 2000, Stanford DASH





sgi













Origin 2000 Network Detail 1

· Router characteristics

- 6 ports connected via a crossbar
- · Each port bidirectional at 800 MB/s in each direction
- Provides 4 virtual channels
- · Input-buffered with pipelined crossbar arbitration
- · Best-case (fall-through) input-to-output latency of 50 ns
- · Links (per direction)
 - 20 data bits, 2 clock bits (differential), 1 data framing bit
 - · Clock rate of 200 MHz, sampled on both edges (400 MHz data rate)
 - · Credit-based flow control
 - · Sliding-window, CRC-based error detection/retransmission

sgi

Origin 2000 Network Detail 2

- Implementation
- 850K-gate ASIC
 - IBM CMOS 5L (0.5µ drawn), 5 metal layers
 - 160 mm² die area
 - · Core operates at 3.3V, 100 MHz
 - · 29 W worst-case power dissipation

Cables

- · Shielded, electrically-matched wires, 1-5 m
- Expensive :--(

sgi



Trend: Merging of Network Interface and CPU

- Desire to move network interface "closer" to the CPU
 - · Architecturally
 - · User-level, protected access ("OS bypass")
 - Tied more closely to memory system (address translation, etc.)
 - · Physically
 - Place on same die as CPU
 - Direct datapaths between CPU internals and network interface
- Challenges
 - Development of reasonable interface to user jobs
 - Electrical, mechanical, physical integration of CPU logic and network interface

Future Trends

- · Communication ever-more critical to overall system performance
- · Bandwidth demands growing
 - · CPU bandwidth growing, both of system bus and functional units
 - Memory system bandwidth growing: SDR, DDR, DRDRAM
- · Network latency becoming more of a problem
 - Decreasing in absolute time
 - · But increasing when measured in CPU instruction issue slots
 - Latency impact on overall performance is non-linear
- Will interconnection network become primary limit on overall system performance?



Trend: Active Networks

- Current multiprocessor networks are "passive"
 - · Message unchanged as it flows through network
 - Network does not interpret message contents
 - · Result: network acts mainly as a delay element (though a useful one!)
- Idea: perform computation in the network as well as on CPUs
- Benefits
 - · Moves computation closer to the data on which it operates
 - Offloads CPUs
- Challenges
 - Programming model, compiler and OS support, protection, etc.
 - Details of computational resources, integration into network fabric, etc.

sgi

sgi

Conclusions

· Interconnect is a key component of multiprocessor system performance

· Interconnect latency and bandwidth are both important

Low latency especially critical for cache coherence

Bandwidth for message passing, clustering, traffic bursts

Future interconnects must continue to improve latency and bandwidth

19

By coupling the network more closely to the CPU

By (eventually) making the networks "active"

sgi

8:45am - 9:15am Tues, 11 May - 2.2

The Role of Optics in

Balanced Computer System Design

Mike Chastain Hewlett-Packard chastain@rsn.hp.com

Mike Chastain

Workshop on Interconnections Within High Speed Digital Systems

April 29, 1999

Para	allel F	iber Optic I	Developmen	at	PACKARD
The advan • Physica • Com • Greater • Reduce	ntages o al size i nectors er comm ed susc	f parallel fiber v eduction at high and cables nunication dista eptibility to EM	versus copper i h frequency nce I and EMC	nterconne	ct are well known
Computer • Costs h Computer • Waits f	r indust have alv r indust for tech	ry has watched ways prevented ry is also slow to mology cost cro	parallel fiber d wide spread sy adopt new inf ssover; or some	evelopmen stem inser erconnect e "external	nt for five years tion technologies I" forcing function
Industry in • Real pr • Costs a • Break • Optin • Costs a	investm products are star kthroug mistic p are still	ent is making pa now appearing ting to come do hs in manufacta rojections of hig high relative to	arallel fiber mo from multiple wn uring and pack th volume inser copper for sho	re viable vendors aging tion rt (<10m)	links
Are there	other f	orces, outside op	ptical developm	ent, that r	nay hasten insertion?
ike Chastain	Wa	rkshop on Interconnec	tions Within High Spe	ed Digital Syst	ems April 29, 1995
like Chastein	mside	rkshop on Interconnec	tions Within High Spe Drmance	ed Digital Syst	April 29, 1995
ike Chasteln Col The indu • Single • Single Increasin • Soon Increasin • Desig • Desig	nside Istry is le chip (le chip (o CPUs) ng band gners st 20 mg band	r CPU Perfe now increasing CPUs breaking t CPUs incorpora performance d may require 8 G lwidths forcing i ruggling to mai	ormance CPU performa the Gigahertz b ting "super-cou riving correspo (B/sec, or more maximum freq ntain reasonab	nce at an o parrier, an nputer arc onding inci , to sustain uency at a le pin cour	ems April 29, 1995 EXP HEWLETT EXPONENTIAL TALE d beyond chitecture tricks" rease in bandwidths a performance dl CPU and ASIC pins nos for manufacturing nervasive nrohlem
ike Chastein Col The indu • Single • Single Increasin • Soon Increasin • Desig • Intel 10	mside nside stry is le chip (ccPUs) occus ng banc gners st 's endor	r CPU Perfo now increasing CPUs breaking to CPUs breaking to CPUs incorpora performance d may require 8 G lwidths forcing ruggling to main resement of Ram Server System 1	ormance CPU performa the Gigahertz b ting "super-cor riving correspo (B/sec, or more maximum freq ntain reasonab) bus is an indica Bandwidth	nce at an o aarrier, an nputer aro onding inco , to sustain uency at a le pin com ation of a p	ems April 29, 1995
ike Chasteln Col The indu • Single • Single Increasin • Soon Increasin • Desig • Intel? 10	mside nstry is le chip (le chip (ng CPU) CPUs 1 ng banc gners st 's endor 000 100- 10- 1-	r CPU Perfo now increasing CPUs breaking t CPUs incorpora performance d may require 8 G widths forcing ruggling to main resement of Ram Server System 1	ormance CPU performa the Gigahertz b ting "super-cou riving correspo (B/sec, or more maximum freq ntain reasonab) bus is an indica Bandwidth	nce at an o barrier, an nputer arc onding inci , to sustain uency at a le pin cour ation of a p	ems April 29, 1995

Cons	
	ider Copper Interconnect Limits
At today's i	nterconnect frequencies (up to ~1 Ghz)
• Primary	frequency dependent loss mechanism is skin effect
• Propo	tional to 4
At intercon	nect frequencies beyond 1 Ghz
• Dielectr	ic loss starts to dominate
• Impac At intercon	pert frequencies approaching 2.5 Gbz
• Intercon	neet distance may be limited to a single backplane or PC planer
• New lo	w loss PCB materials will be required
At intercon	nect frequencies approaching 5.0 Ghz
 PCB int 	erconnects may no longer practical
Copper cab	les are still an option; for now
 Designe 	rs will trade copper trace for cable to increase interconnect distance
• 4-5" of	PCB trace is roughly equivalent in loss to 3 feet of copper cable
• Parallel	copper cables will still be limited to adjacent racks
• Six to t	en meters at 622 Mnz, dropping linearly with frequency
• Machine ro	om level interconnects are already in jeopardy without parallel fiber!
ke Chastain	Workshop on Interconnections Within High Speed Digital Systems April 29, 1999
Cons	ider Server Packaging Density
Increased i	nterconnect frequencies coupled with greater interconnect losses
Increased i • Driving	nterconnect frequencies coupled with greater interconnect losses system designers to reduce interconnect distances
Increased i • Driving • Drivin	nterconnect frequencies coupled with greater interconnect losses system designers to reduce interconnect distances g system designers to increase system packaging density
Increased i • Driving • Drivin To achieve	nterconnect frequencies coupled with greater interconnect losses system designers to reduce interconnect distances g system designers to increase system packaging density increased packaging density
Increased i • Driving • Drivin To achieve • More A	nterconnect frequencies coupled with greater interconnect losses system designers to reduce interconnect distances g system designers to increase system packaging density increased packaging density SIC integration to minimize component count
Increased i • Driving • Drivin To achieve • More A • More C	nterconnect frequencies coupled with greater interconnect losses system designers to reduce interconnect distances g system designers to increase system packaging density increased packaging density SIC integration to minimize component count PUs, ASICs, and RAM per PCB area
Increased i • Driving • Drivin To achieve • More A • More C (Frequency	nterconnect frequencies coupled with greater interconnect losses system designers to reduce interconnect distances g system designers to increase system packaging density increased packaging density SIC integration to minimize component count PUs, ASICs, and RAM per PCB area * * Density) is driving power density to the limits
Increased i • Driving • Drivin To achieve • More A • More C (Frequency • More g	nterconnect frequencies coupled with greater interconnect losses system designers to reduce interconnect distances g system designers to increase system packaging density increased packaging density SIC integration to minimize component count PUs, ASICs, and RAM per PCB area ^{, *} Density) is driving power density to the limits ates at greater frequency ==> more power density!
Increased i • Driving • Drivin To achieve • More A • More C (Frequency • More g • More h	nterconnect frequencies coupled with greater interconnect losses system designers to reduce interconnect distances g system designers to increase system packaging density increased packaging density SIC integration to minimize component count PUs, ASICs, and RAM per PCB area * * Density) is driving power density to the limits ates at greater frequency ==> more power density! igh speed I/O ==> more power density!
Increased i • Driving • Drivin To achieve • More A • More C (Frequency • More g • More h • More C	nterconnect frequencies coupled with greater interconnect losses a system designers to reduce interconnect distances g system designers to increase system packaging density increased packaging density SIC integration to minimize component count PUs, ASICs, and RAM per PCB area * Density) is driving power density to the limits ates at greater frequency ==> more power density! igh speed I/O ==> more power density! PUs, ASICs, and RAM per PCB area ==> more power density!
Increased i • Driving • Drivin To achieve • More A • More C (Frequency • More g • More h • More C System poo	nterconnect frequencies coupled with greater interconnect losses a system designers to reduce interconnect distances g system designers to increase system packaging density increased packaging density SIC integration to minimize component count PUs, ASICs, and RAM per PCB area * Density) is driving power density to the limits ates at greater frequency ==> more power density! igh speed I/O ==> more power density! PUs, ASICs, and RAM per PCB area ==> more power density! ver density will soon exceed machine room limitations
Increased i • Driving • Drivin To achieve • More A • More C (Frequency • More b • More C System pov • By 2002	nterconnect frequencies coupled with greater interconnect losses system designers to reduce interconnect distances g system designers to increase system packaging density increased packaging density SIC integration to minimize component count PUs, ASICs, and RAM per PCB area * Density) is driving power density to the limits ates at greater frequency => more power density! gh speed I/O => more power density! PUs, ASICs, and RAM per PCB area => more power density! ver density will soon exceed machine room limitations >3, (4 CPUs + ASICs + 16 GB DRAM + Power) => ~850 watts
Increased i • Driving • Drivin To achieve • More A • More C (Frequency • More b • More b • More C System poo • By 2002 • Existin	nterconnect frequencies coupled with greater interconnect losses system designers to reduce interconnect distances g system designers to increase system packaging density increased packaging density SIC integration to minimize component count PUS, ASICs, and RAM per PCB area * Density) is driving power density to the limits ates at greater frequency => more power density! igh speed I/O => more power density! PUS, ASICs, and RAM per PCB area => more power density! ver density will soon exceed machine room limitations >3, (4 CPUs + ASICs + 16 GB DRAM + Power) => ~850 watts g rooms are designed for 40-70 W/sq.ft. with an 18" raised floor
Increased i • Driving • Drivin To achieve • More A • More C (Frequency • More g • More g • More b • More C System pov • By 200; • Existin • Floor	nterconnect frequencies coupled with greater interconnect losses system designers to reduce interconnect distances g system designers to increase system packaging density increased packaging density SIC integration to minimize component count PUs, ASICs, and RAM per PCB area * Density) is driving power density to the limits ates at greater frequency ==> more power density! igh speed I/O ==> more power density! PUs, ASICs, and RAM per PCB area ==> more power density! ver density will soon exceed machine room limitations >3, (4 CPUs + ASICs + 16 GB DRAM + Power) => ~850 watts g rooms are designed for 40-70 W/sq.ft. with an 18" raised floor area + service area => 19" rack occupies ~14 sq.ft. => 980 watts max maked area till independent (126 W/ref area for area to minimation)
Increased i • Driving • Drivin To achieve • More A • More C (Frequency • More b • More C System pov • By 2002 • Existin • Floor • New sta	nterconnect frequencies coupled with greater interconnect losses system designers to reduce interconnect distances g system designers to increase system packaging density increased packaging density SIC integration to minimize component count 'PUs, ASICs, and RAM per PCB area '* Density) is driving power density to the limits ates at greater frequency =>> more power density! igh speed I/O =>> more power density! PUs, ASICs, and RAM per PCB area =>> more power density! ver density will soon exceed machine room limitations 2-3, (4 CPUs + ASICs + 16 GB DRAM + Power) => ~850 watts g rooms are designed for 40-70 W/sq.ft. with an 18" raised floor area + service area => 19" rack occupies ~14 sq.ft. => 980 watts max indards are still inadequate (125 W/sq.ft. 36" raised floor suggested)
Increased i • Driving • Drivin To achieve • More A • More C (Frequency • More b • More b • More C System pov • By 2002 • Existin • Floor • New sta Result: Page	nterconnect frequencies coupled with greater interconnect losses system designers to reduce interconnect distances g system designers to increase system packaging density increased packaging density SIC integration to minimize component count PUs, ASICs, and RAM per PCB area * Density) is driving power density to the limits ates at greater frequency => more power density! igh speed I/O => more power density! PUs, ASICs, and RAM per PCB area => more power density! ver density will soon exceed machine room limitations 2-3, (4 CPUs + ASICs + 16 GB DRAM + Power) => ~850 watts 3 rooms are designed for 40-70 W/sq.ft. with an 18" raised floor area + service area => 19" rack occupies ~14 sq.ft. => 980 watts max indards are still inadequate (125 W/sq.ft. 36" raised floor suggested) kaging density limited by machine room for foreseeable future
Increased i • Driving • Drivin To achieve • More A • More C (Frequency • More b • More b • More b • More C System pov • By 2000 • Existin • Floor • New sta Result: Pac	nterconnect frequencies coupled with greater interconnect losses system designers to reduce interconnect distances g system designers to increase system packaging density increased packaging density SIC integration to minimize component count (PUs, ASICs, and RAM per PCB area * Density) is driving power density to the limits ates at greater frequency => more power density! (PUs, ASICs, and RAM per PCB area => more power density! PUs, ASICs, and RAM per PCB area => more power density! (PUs, ASICs, and area => more power density! (PUs, ASICs, and area => more power density! (PUs, ASICs, and area => more power density! (PUs, ASICs, area
Increased i	nterconnect frequencies coupled with greater interconnect losses system designers to reduce interconnect distances g system designers to increase system packaging density increased packaging density SIC integration to minimize component count (PUS, ASICS, and RAM per PCB area * Density) is driving power density to the limits ates at greater frequency => more power density! (PUS, ASICS, and RAM per PCB area => more power density! (PUS, ASICS, and RAM per PCB area => more power density! ver density will soon exceed machine room limitations 2-3, (4 CPUS + ASICS + 16 GB DRAM + Power) => ~850 watts g rooms are designed for 40-70 W/sq.ft. with an 18" raised floor area + service area => 19" rack occupies ~14 sq.ft. => 980 watts max indards are still inadequate (125 W/sq.ft. 36" raised floor suggested) taging density limited by machine room for foreseeable future Workshop on Interconnections Within High Speed Digital Systems April 29, 1999
Increased i • Driving • Drivin To achieve • More A • More C (Frequency • More g • More f • More f • More C System por • By 2002 • Existin • Floor • New sta Result: Pac	nterconnect frequencies coupled with greater interconnect losses system designers to reduce interconnect distances g system designers to increase system packaging density increased packaging density SIC integration to minimize component count PUS, ASICs, and RAM per PCB area * Density) is driving power density to the limits ates at greater frequency => more power density! PUS, ASICs, and RAM per PCB area => more power density! Wer density will soon exceed machine room limitations 2.3, (4 CPUs + ASICs + 16 GB DRAM + Power) => ~850 watts g rooms are designed for 40-70 W/sq.ft. with an 18" raised floor area + service area => 19" rack occupies ~14 sq.ft. => 980 watts max undards are still inadequate (125 W/sq.ft. 36" raised floor suggested) *kaging density limited by machine room for foreseeable future Workshop on Interconnections Within High Speed Digital Systems April 29, 1999
Increased i • Driving • Drivin To achieve • More A • More C (Frequency • More g • More f • More f • More C System pov • By 2002 • Existin • Floor • New sta Result: Pac	nterconnect frequencies coupled with greater interconnect losses system designers to reduce interconnect distances g system designers to increase system packaging density increased packaging density SIC integration to minimize component count PUS, ASICs, and RAM per PCB area * Density) is driving power density to the limits ates at greater frequency => more power density! igh speed I/O => more power density! PUS, ASICs, and RAM per PCB area => more power density! ver density will soon exceed machine room limitations 2, (4 CPUs + ASICs + 16 GB DRAM + Power) => ~850 watts g rooms are designed for 40-70 W/sq.ft. with an 18" raised floor area + service area => 19" rack occupies ~14 sq.ft. => 980 watts max undards are still inadequate (125 W/sq.ft. 36" raised floor suggested) *kaging density limited by machine room for foreseeable future Workshop on Interconnections Within High Speed Digital Systems April 28, 1899

Future Server Designs ? May therefore consist of medium (CPU count) SMPs • High frequency signaling on all interconnects High integration (and high power) silicon ·High density packaging, power input, and power dissipation Tightly coupled electrically, but not physically

Tightly integrated coherent cable interconnects
Utilize copper until frequency-versus-distance becomes prohibitive
Shift to optical as frequency increases and/or costs come down

- Frequency may cause a shift in spite of costs!
 Shift to optical now for machine room level interconnects
 Such as the emerging Future I/O standard

To balance system performance versus machine room constraints

- Spread system across multiple racks to distribute thermal load • Match machine room capabilities (power/sq.ft. electrical and thermal)
- Perhaps integrated with storage and I/O components
- · Good volume utilization without adding significant power/sq.ft.
- Mike Chastain

April 29, 1999

_ The Role of Optics in Server Evolution	HEWLETT
ptical interconnects within servers	
• Evolutionary "copper replacement" strategy as frequ	ency increases
 System architects must work closely with the optical 	link community
 Copper link designs must be compatible with optic 	al link limitations
 Optical components must be compatible with serve 	r manufacturing
 Optical packaging must be consistent with server c 	onnector requirements
 Evolutionary "EMI/EMC management" strategy as f 	requency increases
ntical interconnects between servers and/or I/O within	a machine room
• Addressed by the emerging Future I/O standard	
• Parallel optical links at 1/2/4 GB/sec data delivery; 1	1p to 300m (@1GB)
Network like protocols optimized for both cluster an	d I/O communication
Designed for highly reliable, fault tolerant communi	cation

Mike Chastain Workshop on Interconnections Within High Speed Digital Systems April 29, 1999

		_
Server	Architects must Design for Optics	
Server archit	ects are starting to design within limits of availa	ble optics
• Accepting • Leveragi	12 bit wide links as a cost effective limit ng telecom volumes for connectors and cables	
• Accepting • Designin	per-bit encoding and self-clocking for AC coup g in clock recovery and link training sequences	led links
• Accepting • Performi	multiple bit time skews between end points ng parallel word re-assembly in end points	
• Accepting • Designing	a "non-zero BER" at high frequencies g in transparent link retry and ECC recovery m	echanisms

Optical Vendors must Design for Servers	HEWLETT PACKARD
Optical link frequency has been driven by the telecom indust • Telecom road map is 4x per generation; 622Mhz, 2.5Ghz, • Server road map is 2x per generation; 1.25Ghz, 2.5Ghz, 5	ry , 10Ghz ;Ghz, 10Ghz(?!)
Optical link packaging is not consistent with server environm • Server power is generally noisy; Optical links want clean • Servers (most) rely on forced air convection for thermal r • Optical interfaces cannot assume heat conduction to PC • Server manufacturing relies on robotic assembly and test • Optical interfaces should support standard pick&place 1 • Servers need blind-mate optical connectors; with EMI co • 2nd level assemblies to accomplish blind-mate/containm	eents power nanagement B BGA processes ntainment! ent are expensive
Servers need "transparent" optical links • Server silicon must be re-used; copper links may become • Different "products" must make different distance-cost • Electrical interfaces consistent with (same as) copper cab • Same frequency, encoded self-clocked, low voltage differ • Few (if any) special system considerations beyond equival • Example: special system requirements to handle "eye sa	optic links trade-off le interfaces rential interfaces lent copper cable fety"

Mike Chastain

Workshop on Interconnections Within High Speed Digital Systems

Summary

Parallel optical links are finally close to reality, but costs are still high

Parallel optical links (and Future I/O) will address machine room interconnects

CPU frequency and associated dielectric losses will drive server density upward But, existing machine room capability will limit the power density per sq.ft. Future servers may trade PCB trace for cables and distribute the power density But dielectric loss also reduces copper cable length proportional to frequency Therefore server designers may have (non-cost) reason to use internal optical links

Server and optical link designers must work together to enable a smooth transition

Mike Chastain

Workshop on Interconnections Within High Speed Digital Systems

April 29, 1999

9:15am - 9:45am Tues, 11 May - 2.3



Modern technology: The Memory Wall Latency: cannot access data fast enough Bandwidth: cannot get data to logic fast enough Bandwidth: cannot get data to logic fast enough Next level of supercomputing- Petaflops: Impossible without radical change A direct assault on the problem - HTMT Hybrid Technology, MultiThreaded Mix memory & logic, interconnect optically









1













10th Workshop on Interconnections, Santa Fe SNTAFE99.PPT 11

5/12/99









roject	Funded By	Partners	Goal
IM ast	JPL	LM	Build Space Computer
NVA -	DARPA	USC, CalTech	Build Smart Memory
TMT	NASA, NSA	Cal Tech, Princeton, SUNY, Dela	Design 1 st Petaflop Computer









5/12/99 19th Werkshop on Interconnections, Santa Fe SNTAFE99.PPT 20









	ilicon Bu F <u>MT DR</u>	dget AM	for PII	· VI
 Designed to p memory & su Different Vo In 2004, 16 T 	provide prop port for f ortex configure B = 4096 gr	per b iber l ations	alan band s => c	ce of lwidth different #s
• Each Chip: Interface Logic By Area	HRAM	FtPt ASAP FtPt ASAP	FtPt ASAP FtPt ASAP	Fiber WDM Optical Receiver
5/12/99 10th V	Vorkshop on Interconn	ections, Se	inta Fe	SNTAFE99.PPT 25

ð The View from a SPELL Capacity Capacity 2 Way Read BW Words W/Flop atency VIFIop Local CRAM 128 0.00000 ocal SRAM 32M 0.00012 Remote CRAM 12N 0.00 0.0 emote SRAN 128G 0.5 Οł 16,00 Single DRAM Clust 512N 0.002 RAM from 1 Cluster 32G 0.12 67.00 DRAM Clusters 0.00 16.00 0.0 RAM from 4 Cluster 1280 67,00 0.3 Ũ Ū SPELL can never "miss"; PIM must "guess" first!

10th Workshop on Interconnections, Santa Fe SNTAFE99.PPT 26



SRAM PIM Functions

- Initiate Gather/Scatter to/from DRAM
- Recognize when sufficient operands arrive in SRAM context block
- Enqueue/Dequeue SRAM block addresses
- Initiate DMA transfers to/from CRAM context block
- Signal SPELL re task initiation
- Prefix operations like Flt Pt Sum

5/12/39 10th Workshop on Interconnections, Santa Fe SNTAFE99.PPT 28

DRAM PIM Functions

- Initialize data structures
- Stride thru regular data structures, transferring to/from SRAM
- Pointer chase thru linked data structures
- "Join-like" operations
- Reorderings

Ň

- Prefix operations
- I/O transfer management
 - DMA, compress/decompress, ...
- 5/12/99 10th Workshop on Interconnections, Santa Fe SNTAFE99.PPT 29

1. Indexed

5/12/99

Conclusions

- The Twin Demons: Latency & Bandwidth
- PIM Technology: Solves the local problem
- Petaflops: Global problem still present
- HTMT: Attach global problem by:
 - Making memory smart so many transfers only "one way"
 - Utilizing best of emerging optical technology for bulk of remaining

5/12/99 10th Workshop on Interconnections, Santa Fe SNTAFE99.PPT 30

9:45am - 10:15am Tues, 11 May - 2.4





























































10:30am - 11:00am Tues, 11 May - 2.5





- regular topology

- 'clockwork' routing

2



























Configuration of time-slot selector (TSS)

optical gate

optical fiber

FDBA

EECC for the second sec

6

<u>日日日</u> 月日日 72 IIEQ /N WCS

M+1 combiner

M+1 splitter







Broadcast-and-select type switch

PSUISASAWSAS @ NTT 1999 O HTT

Configuration of 32-ch wavelength channel selector (WCS) module







C NTT 1000 00 TTT

FF 6

iene (?) NTT

PSU/SASAWSAD7 @ NTT 1899 @ BTTT

Cascaded small TSS units (nm = M)

Ξ

ର

m+1 splitter 1 0 P

WCS

optical fiber optical gate

EDFA

Large TSS Ø F

Composite optical/electrical buffer configuration



PSUISASAWS/10 @ NTT 1999 @ MTT

Configuration of burst-packet receiver with clock-tank circuit



Packet loss in composite output buffer



PSUSASAMS/11 @ NTT 1998 @ HTT -----

Packet reception for large power fluctuation



D-10-EDEA сч Р ิสยาธ SEL EDFA4 Buf S/N level **EDFA3** Elec. 320-Gbit/s system configuration **EDFA2** Ř 320-Gbit/s system design SLG1 WCS_5 **FDFA1** ino 32coupler b (ab) N/S 451 <u>4</u> 30 Stres Stres Stres TSS WCS 、 SS ו ב EDFA4 (7+1) comb AWG SLG1 AWG EDFA2 (7+1) split EAG1 (7+1) split EAG2 EAG2 (7+1) split EAG2 TSS 32 coupler Power level RPS TSS EDFA.1 RPS. Manchester coded 32coupler BWG AWG 12 Gbit/s FWX 15 [32 / DFB/Mod. 9 -15 f 32 οŝ 0 Ŷ 20 Power level (dBm) - buffer size / outputs: 12 optical buffers (M=12)+64 electrical buffers architecture: broadcast-and-select type with WDM output buffers Packet reception for phase fluctuation TTE (D) 6001 TTE (D) ATTENSION THE B 6661 TIN @ MISHWARANDE 20 320-Gbit/s system specifications synchronized clock packet format: 64 bytes (including 4-byte guard band) idata payload - highway speed: 10 Gbit/s (12 Gbit/s internal speed) - optical frequency channel span: 100 GHz (0.8 nm) 9 Digital-ring oscillator input signal frame1 2 Optical input signal Time (ns) Regenerated clock verhead – Regenerated data - transmission code: IM / Manchester - number of inputs / outputs: 32 - wavelength: 1500 nm band ଟିତ୍ତିନ୍ତ frame0 (g) (a) Ð <u></u>



Experiments using broadcast-and-select switch



Conclusion

Large-scale photonic packet switches simple star architecture with modular structure combination of broadband WDM techniques and electrical control circuits Key technologies needed for sub-Tbit/s switch hybrid-integrated 32-ch wavelength channel selector 10-Gbit/s burst-packet receiver level and S/N design Rack-mounted photonic packet switch prototypes 10 Gbit/s x 4 broadcast-and-select type 2.5 Gbit/s x 16 wavelength-routing type

Experiments using wavelength-routing switch



7:00pm - 7:30pm Tues, 11 May - 2.8





















1998	
Metrology for OE	Feb 98
Optical Communications Road map	May 98
Technology Roadmap for Image Sensors	Jun 98
Annual Forum	Oct 98
1999	
Broadband Communications & Switching Components Technology	Apr 99
Advanced Imaging - "Electronic Eye"	Jun 99
International Standards	Aug 99
Annual Forum	Oct 99












OE DARPA Centers History				
Time Period	No. of Centers	Funding		
1990 - 1993	3	15 million		
1994 - 1997	4	25 Million		
1997 - 2000	2	12 Million		
2000 - 2003	4?	25 Million?		
OIDA		0002365004		

Number of Students Pl	laced in Industry from
OE DARPA Cen	ters 1990-1996
NCIPT	71
OTC	48
OMC	60
COST	<u>47</u> 226
	Manage and the



Board & Back-plane Level Optical Interconnections Using Integrated Thin-cladding Polymer Fibers

Yao Li

NEC Research Institute, 4 Independence Way, Princeton, NJ 08540. e-mail: yao@research.nj.nec.com

Other Contributors & Collaborators

Jun Ai	NECI,	USA
Jan Popelek	NECI,	USA
K. Kasabara	NEC CRL.	Japan
Y. Takiguchi	Hamamatsu, KK	Japan

@ IEEE-Santa Fe, 05/11/99

Talk Outline

- * Introductions,
- * POF's as Short-distance Optical Channels,
- * POF's for Intra-computer Interconnections, Project I: multi-Gb/s on-board clock distributions, Project II: 2D parallel optical circuits on PCB.
- * Some experiments,
- * Summary and Conclusions.

Introduction

- * Bandwidth bottleneck at PCB level, (>500 MHz on-chip & <200 MHz off-chip)
- * Problems of conventional waveguides, (high cost for glass waveguides, large loss for polymer ones)
- * Space parallelism can be spplied by VCSEL's, (1D & 2D arrays with low fabrication cost, low threshold current)
- * POF offers low-cost, high-rigidity, low-loss, (1/4 of glass fiber cost, breakage-free, < 3 dB/m)
- * Low-cost Polymer fiber-image-guides (PFIG's) are also becoming commecially available

7:30pm - 8:00pm Tues, 11 May - 2.9

Basics of Polymer Optical Fibers

- * 1st. POF in 1970's but progress was slow,
- * Main applications now in display & lighting,
- Low material & production cost, (1/4 of cost of silica fibers)
- * High attenuation, (PMMA:12 db/100 m @ 650 nm),
- * Thin-cladding (90% core) & Multimodes,
- * Low operating temperature (-20 to 80 °C),
- * High flexibility and rigidity against breakage.

Two Interconnect Projects at NECI for Board-level POF Circuits

- * Multi-Gb/s Optical Clock Distribution Circuit, (10 Gb/s, 128 port, connectorized integrated optics)
- * 2D Parallel Optical Circuits for VCSEL Arrays, (both point-to-point and multi point capability)

Board-level Optical Clock Distribution Using End-tapered Fiber Bundles





Main Characteristics of Mitsubishi Thin-cladding PMMA Fibers

Temporal & Frequency Domain Measurements

- 30 fs Ti: Sapphire laser at $\lambda = 850$ nm, 10 ps Synchroscan Streak Camera, 1.4 ps Maximum Readout Accuracy

Pulse width & Skew

Bandwidth

μ single 0. Intensity (A.U.) Normalized Power 0.6 400 45.0 pa (0.4 38.9 ps (slagte) 200 0.2 0.0 L 100 200 10 100 Time (ps) Frequency (GHz) **skew** = $\sqrt{45^2 - 39^2} = 22.5$ ps

Embedded Optical Circuit Board for VCSEL Arrays



Main Characteristics of a Prototype PFIG



Fabrication of Circuit Preforms Using Thermo Bending



Preformed Elements





Bending Response at 90°



Prototype Board & VCSEL Transmissions



Power & Resolution Performance Measures





Summary and Conclusions

- * POF suits better for short-distance applications,
- * POF offers better packaging capabilities,
- * Multi Gb/s bandwidth is sustainable,
- * Free-space optics can add value to POF circuits,
- * Packaging capability determines practicality.

NOTES		
· · · · · · · · · · · · · · · · · · ·		

Wednesday, 12 May 1999

NOTES