

---

<b>M.1</b>	Introduction	M-2
<b>M.2</b>	The Early Development of Computers (Chapter 1)	M-2
<b>M.3</b>	The Development of Memory Hierarchy and Protection (Chapter 2 and Appendix B)	M-9
<b>M.4</b>	The Evolution of Instruction Sets (Appendices A, J, and K)	M-17
<b>M.5</b>	The Development of Pipelining and Instruction-Level Parallelism (Chapter 3 and Appendices C and H)	M-27
<b>M.6</b>	The Development of SIMD Supercomputers, Vector Computers, Multimedia SIMD Instruction Extensions, and Graphical Processor Units (Chapter 4)	M-45
<b>M.7</b>	The History of Multiprocessors and Parallel Processing (Chapter 5 and Appendices F, G, and I)	M-55
<b>M.8</b>	The Development of Clusters (Chapter 6)	M-74
<b>M.9</b>	Historical Perspectives and References	M-79
<b>M.10</b>	The History of Magnetic Storage, RAID, and I/O Buses (Appendix D)	M-84

M

---

## Historical Perspectives and References

If ... history ... teaches us anything, it is that man in his quest for knowledge and progress is determined and cannot be deterred.

**John F. Kennedy**

*Address at Rice University (1962)*

Those who cannot remember the past are condemned to repeat it.

**George Santayana**

*The Life of Reason (1905), Vol. 2, Chapter 3*

---

## M.1

### Introduction

This appendix provides historical background on some of the key ideas presented in the chapters. We may trace the development of an idea through a series of machines or describe significant projects. If you are interested in examining the initial development of an idea or machine or are interested in further reading, references are provided at the end of each section.

Section M.2 starts us off with the invention of the digital computer and corresponds to Chapter 1. Section M.3, on memory hierarchy, corresponds to Chapter 2 and Appendix B. Section M.4, on instruction set architecture, covers Appendices A, J, and K. Section M.5, on pipelining and instruction-level parallelism, corresponds to Chapter 3 and Appendices C and H. Section M.6, on data-level parallelism in vector, SIMD, and GPU architectures, corresponds to Chapter 4. Section M.7, on multiprocessors and parallel programming, covers Chapter 5 and Appendices F, G, and I. Section M.8, on the development of clusters, covers Chapter 6. Finally, Section M.9, on I/O, corresponds to Appendix D.

---

## M.2

### The Early Development of Computers (Chapter 1)

In this historical section, we discuss the early development of digital computers and the development of performance measurement methodologies.

#### The First General-Purpose Electronic Computers

J. Presper Eckert and John Mauchly at the Moore School of the University of Pennsylvania built the world's first fully operational electronic general-purpose computer. This machine, called ENIAC (Electronic Numerical Integrator and Calculator), was funded by the U.S. Army and became operational during World War II, but it was not publicly disclosed until 1946. ENIAC was used for computing artillery firing tables. The machine was enormous—100 feet long, 8½ feet high, and several feet wide. Each of the 20 ten-digit registers was 2 feet long. In total, there were 18,000 vacuum tubes.

Although the size was three orders of magnitude bigger than the size of the average machines built today, it was more than five orders of magnitude slower, with an add taking 200 microseconds. The ENIAC provided conditional jumps and was programmable, which clearly distinguished it from earlier calculators. Programming was done manually by plugging up cables and setting switches and required from a half hour to a whole day. Data were provided on punched cards. The ENIAC was limited primarily by a small amount of storage and tedious programming.

In 1944, John von Neumann was attracted to the ENIAC project. The group wanted to improve the way programs were entered and discussed storing programs as numbers; von Neumann helped crystallize the ideas and wrote a memo proposing a stored-program computer called EDVAC (Electronic Discrete Variable

Automatic Computer). Herman Goldstine distributed the memo and put von Neumann's name on it, much to the dismay of Eckert and Mauchly, whose names were omitted. This memo has served as the basis for the commonly used term *von Neumann computer*. Several early inventors in the computer field believe that this term gives too much credit to von Neumann, who conceptualized and wrote up the ideas, and too little to the engineers, Eckert and Mauchly, who worked on the machines. Like most historians, your authors (winners of the 2000 IEEE von Neumann Medal) believe that all three individuals played a key role in developing the stored-program computer. Von Neumann's role in writing up the ideas, in generalizing them, and in thinking about the programming aspects was critical in transferring the ideas to a wider audience.

In 1946, Maurice Wilkes of Cambridge University visited the Moore School to attend the latter part of a series of lectures on developments in electronic computers. When he returned to Cambridge, Wilkes decided to embark on a project to build a stored-program computer named EDSAC (Electronic Delay Storage Automatic Calculator). (The EDSAC used mercury delay lines for its memory; hence, the phrase "delay storage" in its name.) The EDSAC became operational in 1949 and was the world's first full-scale, operational, stored-program computer [Wilkes, Wheeler, and Gill 1951; Wilkes 1985, 1995]. (A small prototype called the Mark I, which was built at the University of Manchester and ran in 1948, might be called the first operational stored-program machine.) The EDSAC was an accumulator-based architecture. This style of instruction set architecture remained popular until the early 1970s. (Appendix A starts with a brief summary of the EDSAC instruction set.)

In 1947, Mauchly took the time to help found the Association for Computing Machinery. He served as the ACM's first vice-president and second president. That same year, Eckert and Mauchly applied for a patent on electronic computers. The dean of the Moore School, by demanding that the patent be turned over to the university, may have helped Eckert and Mauchly conclude that they should leave. Their departure crippled the EDVAC project, which did not become operational until 1952.

Goldstine left to join von Neumann at the Institute for Advanced Study at Princeton in 1946. Together with Arthur Burks, they issued a report based on the 1944 memo [Burks, Goldstine, and von Neumann 1946]. The paper led to the IAS machine built by Julian Bigelow at Princeton's Institute for Advanced Study. It had a total of 1024 40-bit words and was roughly 10 times faster than ENIAC. The group thought about uses for the machine, published a set of reports, and encouraged visitors. These reports and visitors inspired the development of a number of new computers, including the first IBM computer, the 701, which was based on the IAS machine. The paper by Burks, Goldstine, and von Neumann was incredible for the period. Reading it today, you would never guess this landmark paper was written more than 50 years ago, as most of the architectural concepts seen in modern computers are discussed there (e.g., see the quote at the beginning of Chapter 2).

In the same time period as ENIAC, Howard Aiken was designing an electro-mechanical computer called the Mark-I at Harvard. The Mark-I was built by a team

of engineers from IBM. He followed the Mark-I with a relay machine, the Mark-II, and a pair of vacuum tube machines, the Mark-III and Mark-IV. The Mark-III and Mark-IV were built after the first stored-program machines. Because they had separate memories for instructions and data, the machines were regarded as reactionary by the advocates of stored-program computers. The term *Harvard architecture* was coined to describe this type of machine. Though clearly different from the original sense, this term is used today to apply to machines with a single main memory but with separate instruction and data caches.

The Whirlwind project [Redmond and Smith 1980] began at MIT in 1947 and was aimed at applications in real-time radar signal processing. Although it led to several inventions, its overwhelming innovation was the creation of magnetic core memory, the first reliable and inexpensive memory technology. Whirlwind had 2048 16-bit words of magnetic core. Magnetic cores served as the main memory technology for nearly 30 years.

## Important Special-Purpose Machines

During World War II, major computing efforts in both Great Britain and the United States focused on special-purpose code-breaking computers. The work in Great Britain was aimed at decrypting messages encoded with the German Enigma coding machine. This work, which occurred at a location called Bletchley Park, led to two important machines. The first, an electromechanical machine, conceived of by Alan Turing, was called BOMB [see Good in Metropolis, Howlett, and Rota 1980]. The second, much larger and electronic machine, conceived and designed by Newman and Flowers, was called COLOSSUS [see Randall in Metropolis, Howlett, and Rota 1980]. These were highly specialized cryptanalysis machines, which played a vital role in the war by providing the ability to read coded messages, especially those sent to U-boats. The work at Bletchley Park was highly classified (indeed, some of it is still classified), so its direct impact on the development of ENIAC, EDSAC, and other computers is difficult to trace, but it certainly had an indirect effect in advancing the technology and gaining understanding of the issues.

Similar work on special-purpose computers for cryptanalysis went on in the United States. The most direct descendent of this effort was the company Engineering Research Associates (ERA) [see Thomash in Metropolis, Howlett, and Rota 1980], which was founded after the war to attempt to commercialize on the key ideas. ERA built several machines that were sold to secret government agencies, and it was eventually purchased by Sperry-Rand, which had earlier purchased the Eckert Mauchly Computer Corporation.

Another early set of machines that deserves credit was a group of special-purpose machines built by Konrad Zuse in Germany in the late 1930s and early 1940s [see Bauer and Zuse in Metropolis, Howlett, and Rota 1980]. In addition to producing an operating machine, Zuse was the first to implement floating point, which von Neumann claimed was unnecessary! His early machines used a mechanical store that was smaller than other electromechanical solutions of the

time. His last machine was electromechanical but, because of the war, was never completed.

An important early contributor to the development of electronic computers was John Atanasoff, who built a small-scale electronic computer in the early 1940s [Atanasoff 1940]. His machine, designed at Iowa State University, was a special-purpose computer (called the ABC, for Atanasoff Berry Computer) that was never completely operational. Mauchly briefly visited Atanasoff before he built ENIAC, and several of Atanasoff's ideas (e.g., using binary representation) likely influenced Mauchly. The presence of the Atanasoff machine, delays in filing the ENIAC patents (the work was classified, and patents could not be filed until after the war), and the distribution of von Neumann's EDVAC paper were used to break the Eckert–Mauchly patent [Larson 1973]. Though controversy still rages over Atanasoff's role, Eckert and Mauchly are usually given credit for building the first working, general-purpose, electronic computer [Stern 1980]. Atanasoff, however, demonstrated several important innovations included in later computers. Atanasoff deserves much credit for his work, and he might fairly be given credit for the world's first special-purpose electronic computer and for possibly influencing Eckert and Mauchly.

## Commercial Developments

In December 1947, Eckert and Mauchly formed Eckert-Mauchly Computer Corporation. Their first machine, the BINAC, was built for Northrop and was shown in August 1949. After some financial difficulties, the Eckert-Mauchly Computer Corporation was acquired by Remington-Rand, later called Sperry-Rand. Sperry-Rand merged the Eckert-Mauchly acquisition, ERA, and its tabulating business to form a dedicated computer division, called UNIVAC. UNIVAC delivered its first computer, the UNIVAC I, in June 1951. The UNIVAC I sold for \$250,000 and was the first successful commercial computer—48 systems were built! Today, this early machine, along with many other fascinating pieces of computer lore, can be seen at the Computer History Museum in Mountain View, California. Other places where early computing systems can be visited include the Deutsches Museum in Munich and the Smithsonian Institution in Washington, D.C., as well as numerous online virtual museums.

IBM, which earlier had been in the punched card and office automation business, didn't start building computers until 1950. The first IBM computer, the IBM 701 based on von Neumann's IAS machine, shipped in 1952 and eventually sold 19 units [see Hurd in Metropolis, Howlett, and Rota 1980]. In the early 1950s, many people were pessimistic about the future of computers, believing that the market and opportunities for these “highly specialized” machines were quite limited. Nonetheless, IBM quickly became the most successful computer company. Their focus on reliability and customer- and market-driven strategies were key. Although the 701 and 702 were modest successes, IBM's follow-up machines, the 650, 704, and 705 (delivered in 1954 and 1955) were significant successes, each selling from 132 to 1800 computers.

Several books describing the early days of computing have been written by the pioneers [Goldstine 1972; Wilkes 1985, 1995], as well as Metropolis, Howlett, and Rota [1980], which is a collection of recollections by early pioneers. There are numerous independent histories, often built around the people involved [Slater 1987], as well as a journal, *Annals of the History of Computing*, devoted to the history of computing.

### Development of Quantitative Performance Measures: Successes and Failures

In the earliest days of computing, designers set performance goals—ENIAC was to be 1000 times faster than the Harvard Mark-I, and the IBM Stretch (7030) was to be 100 times faster than the fastest machine in existence. What wasn't clear, though, was how this performance was to be measured. In looking back over the years, it is a consistent theme that each generation of computers obsoletes the performance evaluation techniques of the prior generation.

The original measure of performance was time to perform an individual operation, such as addition. Since most instructions took the same execution time, the timing of one gave insight into the others. As the execution times of instructions in a machine became more diverse, however, the time for one operation was no longer useful for comparisons. To take these differences into account, an *instruction mix* was calculated by measuring the relative frequency of instructions in a computer across many programs. The Gibson mix [Gibson 1970] was an early popular instruction mix. Multiplying the time for each instruction times its weight in the mix gave the user the *average instruction execution time*. (If measured in clock cycles, average instruction execution time is the same as average cycles per instruction.) Since instruction sets were similar, this was a more accurate comparison than add times. From average instruction execution time, then, it was only a small step to MIPS (as we have seen, the one is the inverse of the other). MIPS had the virtue of being easy for the layperson to understand.

As CPUs became more sophisticated and relied on memory hierarchies and pipelining, there was no longer a single execution time per instruction; MIPS could not be calculated from the mix and the manual. The next step was benchmarking using kernels and synthetic programs. Curnow and Wichmann [1976] created the Whetstone synthetic program by measuring scientific programs written in Algol 60. This program was converted to FORTRAN and was widely used to characterize scientific program performance. An effort with similar goals to Whetstone, the Livermore FORTRAN Kernels, was made by McMahon [1986] and researchers at Lawrence Livermore Laboratory in an attempt to establish a benchmark for supercomputers. These kernels, however, consisted of loops from real programs.

As it became clear that using MIPS to compare architectures with different instruction sets would not work, a notion of relative MIPS was created. When the VAX-11/780 was ready for announcement in 1977, DEC ran small benchmarks that were also run on an IBM 370/158. IBM marketing referred to the 370/158 as a

1 MIPS computer, and, because the programs ran at the same speed, DEC marketing called the VAX-11/780 a 1 MIPS computer. Relative MIPS for a machine  $M$  was defined based on some reference machine as:

$$\text{MIPS}_M = \frac{\text{Performance}_M}{\text{Performance}_{\text{reference}}} \times \text{MIPS}_{\text{reference}}$$

The popularity of the VAX-11/780 made it a popular reference machine for relative MIPS, especially since relative MIPS for a 1 MIPS computer is easy to calculate: If a machine was five times faster than the VAX-11/780, for that benchmark its rating would be 5 relative MIPS. The 1 MIPS rating was unquestioned for 4 years, until Joel Emer of DEC measured the VAX-11/780 under a time-sharing load. He found that the VAX-11/780 native MIPS rating was 0.5. Subsequent VAXes that ran 3 native MIPS for some benchmarks were therefore called 6 MIPS machines because they ran six times faster than the VAX-11/780. By the early 1980s, the term *MIPS* was almost universally used to mean relative MIPS.

The 1970s and 1980s marked the growth of the supercomputer industry, which was defined by high performance on floating-point-intensive programs. Average instruction time and MIPS were clearly inappropriate metrics for this industry, hence the invention of MFLOPS (millions of floating-point operations per second), which effectively measured the inverse of execution time for a benchmark. Unfortunately, customers quickly forget the program used for the rating, and marketing groups decided to start quoting peak MFLOPS in the supercomputer performance wars.

SPEC (System Performance and Evaluation Cooperative) was founded in the late 1980s to try to improve the state of benchmarking and make a more valid basis for comparison. The group initially focused on workstations and servers in the UNIX marketplace, and these remain the primary focus of these benchmarks today. The first release of SPEC benchmarks, now called SPEC89, was a substantial improvement in the use of more realistic benchmarks. SPEC2006 still dominates processor benchmarks almost two decades later.

## References

- Amdahl, G. M. [1967]. “Validity of the single processor approach to achieving large scale computing capabilities,” *Proc. AFIPS Spring Joint Computer Conf.*, April 18–20, 1967, Atlantic City, N.J., 483–485.
- Atanasoff, J. V. [1940]. “Computing machine for the solution of large systems of linear equations,” Internal Report, Iowa State University, Ames.
- Azizi, O., Mahesri, A., Lee, B. C., Patel, S. J., & Horowitz, M. [2010]. Energy-performance tradeoffs in processor architecture and circuit design: a marginal cost analysis. *Proc. International Symposium on Computer Architecture*, 26–36.
- Bell, C. G. [1984]. “The mini and micro industries,” *IEEE Computer* 17:10 (October), 14–30.
- Bell, C. G., J. C. Mudge, and J. E. McNamara [1978]. *A DEC View of Computer Engineering*, Digital Press, Bedford, Mass.



- Burks, A. W., H. H. Goldstine, and J. von Neumann [1946]. "Preliminary discussion of the logical design of an electronic computing instrument," Report to the U.S. Army Ordnance Department, p. 1; also appears in *Papers of John von Neumann*, W. Aspray and A. Burks, eds., MIT Press, Cambridge, Mass., and Tomash Publishers, Los Angeles, Calif., 1987, 97–146.
- Curnow, H. J., and B. A. Wichmann [1976]. "A synthetic benchmark," *The Computer J.* 19:1, 43–49.
- Dally, William J., "High Performance Hardware for Machine Learning," Cadence Embedded Neural Network Summit, February 9, 2016. [http://ip.cadence.com/uploads/presentations/1000AM\\_Dally\\_Cadence\\_ENN.pdf](http://ip.cadence.com/uploads/presentations/1000AM_Dally_Cadence_ENN.pdf)
- Flemming, P. J., and J. J. Wallace [1986]. "How not to lie with statistics: The correct way to summarize benchmarks results," *Communications of the ACM* 29:3 (March), 218–221.
- Fuller, S. H., and W. E. Burr [1977]. "Measurement and evaluation of alternative computer architectures," *Computer* 10:10 (October), 24–35.
- Gibson, J. C. [1970]. "The Gibson mix," Rep. TR. 00.2043, IBM Systems Development Division, Poughkeepsie, N.Y. (research done in 1959).
- Goldstine, H. H. [1972]. *The Computer: From Pascal to von Neumann*, Princeton University Press, Princeton, N.J.
- Gray, J., and C. van Ingen [2005]. *Empirical Measurements of Disk Failure Rates and Error Rates*, MSR-TR-2005-166, Microsoft Research, Redmond, Wash.
- Jain, R. [1991]. *The Art of Computer Systems Performance Analysis: Techniques for Experimental Design, Measurement, Simulation, and Modeling*, Wiley, New York.
- Kemmel, R. [2000]. "Fibre Channel: A comprehensive introduction," *Internet Week* (April).
- Larson, E. R. [1973]. "Findings of fact, conclusions of law, and order for judgment," File No. 4-67, Civ. 138, *Honeywell v. Sperry-Rand and Illinois Scientific Development*, U.S. District Court for the State of Minnesota, Fourth Division (October 19).
- Lubeck, O., J. Moore, and R. Mendez [1985]. "A benchmark comparison of three supercomputers: Fujitsu VP-200, Hitachi S810/20, and Cray X-MP/2," *Computer* 18:12 (December), 10–24.
- Landstrom, B. [2014]. "The Cost Of Downtime," <http://www.interxion.com/blogs/2014/07/the-cost-of-downtime/>
- McMahon, F. M. [1986]. *The Livermore FORTRAN Kernels: A Computer Test of Numerical Performance Range*, Tech. Rep. UCRL-55745, Lawrence Livermore National Laboratory, University of California, Livermore.
- Metropolis, N., J. Howlett, and G. C. Rota, eds. [1980]. *A History of Computing in the Twentieth Century*, Academic Press, New York.
- Mukherjee S. S., C. Weaver, J. S. Emer, S. K. Reinhardt, and T. M. Austin [2003]. "Measuring architectural vulnerability factors," *IEEE Micro* 23:6, 70–75.
- Oliker, L., A. Canning, J. Carter, J. Shalf, and S. Ethier [2004]. "Scientific computations on modern parallel vector systems," *Proc. ACM/IEEE Conf. on Supercomputing*, November 6–12, 2004, Pittsburgh, Penn., 10.

- Patterson, D. [2004]. “Latency lags bandwidth,” *Communications of the ACM* 47:10 (October), 71–75.
- Redmond, K. C., and T. M. Smith [1980]. *Project Whirlwind—The History of a Pioneer Computer*, Digital Press, Boston.
- Shurkin, J. [1984]. *Engines of the Mind: A History of the Computer*, W. W. Norton, New York.
- Slater, R. [1987]. *Portraits in Silicon*, MIT Press, Cambridge, Mass.
- Smith, J. E. [1988]. “Characterizing computer performance with a single number,” *Communications of the ACM* 31:10 (October), 1202–1206.
- SPEC. [1989]. *SPEC Benchmark Suite Release 1.0* (October 2).
- SPEC. [1994]. *SPEC Newsletter* (June).
- Stern, N. [1980]. “Who invented the first electronic digital computer?” *Annals of the History of Computing* 2:4 (October), 375–376.
- Touma, W. R. [1993]. *The Dynamics of the Computer Industry: Modeling the Supply of Workstations and Their Components*, Kluwer Academic, Boston.
- Weicker, R. P. [1984]. “Dhrystone: A synthetic systems programming benchmark,” *Communications of the ACM* 27:10 (October), 1013–1030.
- Wilkes, M. V. [1985]. *Memoirs of a Computer Pioneer*, MIT Press, Cambridge, Mass.
- Wilkes, M. V. [1995]. *Computing Perspectives*, Morgan Kaufmann, San Francisco.
- Wilkes, M. V., D. J. Wheeler, and S. Gill [1951]. *The Preparation of Programs for an Electronic Digital Computer*, Addison-Wesley, Cambridge, Mass.

---

## M.3

### The Development of Memory Hierarchy and Protection (Chapter 2 and Appendix B)

Although the pioneers of computing knew of the need for a memory hierarchy and coined the term, the automatic management of two levels was first proposed by Kilburn et al. [1962]. It was demonstrated with the Atlas computer at the University of Manchester. This computer appeared the year before the IBM 360 was announced. Although IBM planned for its introduction with the next generation (System/370), the operating system TSS was not up to the challenge in 1970. Virtual memory was announced for the 370 family in 1972, and it was for this computer that the term *translation lookaside buffer* was coined [Case and Padegs 1978]. The only computers today without virtual memory are a few supercomputers, embedded processors, and older personal computers.

Both the Atlas and the IBM 360 provided protection on pages, and the GE 645 was the first system to provide paged segmentation. The earlier Burroughs computers provided virtual memory using segmentation, similar to the segmented address scheme of the Intel 8086. The 80286, the first 80x86 to have the protection mechanisms described in Appendix C, was inspired by the Multics protection software that ran on the GE 645. Over time, computers

evolved more elaborate mechanisms. The most elaborate mechanism was *capabilities*, which attracted the greatest interest in the late 1970s and early 1980s [Fabry 1974; Wulf, Levin, and Harbison 1981]. Wilkes [1982], one of the early workers on capabilities, had this to say:

*Anyone who has been concerned with an implementation of the type just described [capability system], or has tried to explain one to others, is likely to feel that complexity has got out of hand. It is particularly disappointing that the attractive idea of capabilities being tickets that can be freely handed around has become lost ....*

*Compared with a conventional computer system, there will inevitably be a cost to be met in providing a system in which the domains of protection are small and frequently changed. This cost will manifest itself in terms of additional hardware, decreased runtime speed, and increased memory occupancy. It is at present an open question whether, by adoption of the capability approach, the cost can be reduced to reasonable proportions. [p. 112]*

Today there is little interest in capabilities either from the operating systems or the computer architecture communities, despite growing interest in protection and security.

Bell and Strecker [1976] reflected on the PDP-11 and identified a small address space as the only architectural mistake that is difficult to recover from. At the time of the creation of PDP-11, core memories were increasing at a very slow rate. In addition, competition from 100 other minicomputer companies meant that DEC might not have a cost-competitive product if every address had to go through the 16-bit data path twice, hence the architect's decision to add only 4 more address bits than found in the predecessor of the PDP-11.

The architects of the IBM 360 were aware of the importance of address size and planned for the architecture to extend to 32 bits of address. Only 24 bits were used in the IBM 360, however, because the low-end 360 models would have been even slower with the larger addresses in 1964. Unfortunately, the architects didn't reveal their plans to the software people, and programmers who stored extra information in the upper 8 "unused" address bits foiled the expansion effort. (Apple made a similar mistake 20 years later with the 24-bit address in the Motorola 68000, which required a procedure to later determine "32-bit clean" programs for the Macintosh when later 68000s used the full 32-bit virtual address.) Virtually every computer since then will check to make sure the unused bits stay unused and trap if the bits have the wrong value.

As mentioned in the text, system virtual machines were pioneered at IBM as part of its investigation into virtual memory. IBM's first computer with virtual memory was the IBM 360/67, introduced in 1967. IBM researchers wrote the program CP-67 that created the illusion of several independent 360 computers. They then wrote an interactive, single-user operating system called CMS that ran on these virtual machines. CP-67 led to the product VM/370, and today IBM sells z/VM for its mainframe computers [Meyer and Seawright 1970; Van Vleck 2005].

A few years after the Atlas paper, Wilkes published the first paper describing the concept of a cache [1965]:

*The use is discussed of a fast core memory of, say, 32,000 words as slave to a slower core memory of, say, one million words in such a way that in practical cases the effective access time is nearer that of the fast memory than that of the slow memory. [p. 270]*

This two-page paper describes a direct-mapped cache. Although this is the first publication on caches, the first implementation was probably a direct-mapped instruction cache built at the University of Cambridge. It was based on tunnel diode memory, the fastest form of memory available at the time. Wilkes stated that G. Scarrott suggested the idea of a cache memory.

Subsequent to that publication, IBM started a project that led to the first commercial computer with a cache, the IBM 360/85 [Liptay 1968]. Gibson [1967] described how to measure program behavior as memory traffic as well as miss rate and showed how the miss rate varies between programs. Using a sample of 20 programs (each with 3 million references!), Gibson also relied on average memory access time to compare systems with and without caches. This precedent is more than 40 years old, and yet many used miss rates until the early 1990s.

Conti, Gibson, and Pitkowsky [1968] described the resulting performance of the 360/85. The 360/91 outperforms the 360/85 on only 3 of the 11 programs in the paper, even though the 360/85 has a slower clock cycle time (80 ns versus 60 ns), less memory interleaving (4 versus 16), and a slower main memory (1.04 microsecond versus 0.75 microsecond). This paper was also the first to use the term *cache*.

Others soon expanded the cache literature. Strecker [1976] published the first comparative cache design paper examining caches for the PDP-11. Smith [1982] later published a thorough survey paper that used the terms *spatial locality* and *temporal locality*; this paper has served as a reference for many computer designers.

Although most studies relied on simulations, Clark [1983] used a hardware monitor to record cache misses of the VAX-11/780 over several days. Clark and Emer [1985] later compared simulations and hardware measurements for translations.

Hill [1987] proposed the three C's used in Appendix B to explain cache misses. Jouppi [1998] retrospectively said that Hill's three C's model led directly to his invention of the victim cache to take advantage of faster direct-mapped caches and yet avoid most of the cost of conflict misses. Sugumar and Abraham [1993] argued that the baseline cache for the three C's model should use optimal replacement; this would eliminate the anomalies of least recently used (LRU)-based miss classification and allow conflict misses to be broken down into those caused by mapping and those caused by a nonoptimal replacement algorithm.

One of the first papers on nonblocking caches was by Kroft [1981]. Kroft [1998] later explained that he was the first to design a computer with a cache at

Control Data Corporation, and when using old concepts for new mechanisms he hit upon the idea of allowing his two-ported cache to continue to service other accesses on a miss.

Baer and Wang [1988] did one of the first examinations of the multilevel inclusion property. Wang, Baer, and Levy [1989] then produced an early paper on performance evaluation of multilevel caches. Later, Jouppi and Wilton [1994] proposed multilevel exclusion for multilevel caches on chip.

In addition to victim caches, Jouppi [1990] also examined prefetching via streaming buffers. His work was extended by Farkas, Jouppi, and Chow [1995] to streaming buffers that work well with nonblocking loads and speculative execution for in-order processors, and later Farkas et al. [1997] showed that, while out-of-order processors can tolerate unpredictable latency better, they still benefit. They also refined memory bandwidth demands of stream buffers.

Proceedings of the Symposium on Architectural Support for Compilers and Operating Systems (ASPLOS) and the International Computer Architecture Symposium (ISCA) from the 1990s are filled with papers on caches. (In fact, some wags claimed ISCA really stood for the International *Cache* Architecture Symposium.)

Chapter 2 relies on the measurements of SPEC2000 benchmarks collected by Cantin and Hill [2001]. There are several other papers used in Chapter 2 that are cited in the captions of the figures that use the data: Agarwal and Pudar [1993]; Barroso, Gharachorloo, and Bugnion [1998]; Farkas and Jouppi [1994]; Jouppi [1990]; Lam, Rothberg, and Wolf [1991]; Lebeck and Wood [1994]; McCalpin [2005]; Mowry, Lam, and Gupta [1992]; and Torrellas, Gupta, and Hennessy [1992].

## References

- Agarwal, A. [1987]. “Analysis of Cache Performance for Operating Systems and Multiprogramming,” Ph.D. thesis, Tech. Rep. No. CSL-TR-87-332, Stanford University, Palo Alto, Calif.
- Agarwal, A., and S. D. Pudar [1993]. “Column-associative caches: A technique for reducing the miss rate of direct-mapped caches,” *20th Annual Int’l. Symposium on Computer Architecture (ISCA)*, May 16–19, 1993, San Diego, Calif. (*Computer Architecture News* 21:2 (May), 179–190).
- Baer, J.-L., and W.-H. Wang [1988]. “On the inclusion property for multi-level cache hierarchies,” *Proc. 15th Annual Int’l. Symposium on Computer Architecture (ISCA)*, May 30–June 2, 1988, Honolulu, Hawaii, 73–80.
- Barham, P., B. Dragovic, K. Fraser, S. Hand, T. Harris, A. Ho, and R. Neugebauer [2003]. “Xen and the art of virtualization,” *Proc. of the 19th ACM Symposium on Operating Systems Principles*, October 19–22, 2003, Bolton Landing, N.Y.
- Barroso, L. A., K. Gharachorloo, and E. Bugnion [1998]. “Memory system characterization of commercial workloads,” *Proc. 25th Annual Int’l. Symposium on Computer Architecture (ISCA)*, July 3–14, 1998, Barcelona, Spain, 3–14.

- Bell, C. G., and W. D. Strecker [1976]. "Computer structures: What have we learned from the PDP-11?" *Proc. Third Annual Int'l. Symposium on Computer Architecture (ISCA)*, January 19–21, 1976, Tampa, Fla., 1–14.
- Bhandarkar, D. P. [1995]. *Alpha Architecture Implementations*, Digital Press, Newton, Mass.
- Borg, A., R. E. Kessler, and D. W. Wall [1990]. "Generation and analysis of very long address traces," *Proc. 17th Annual Int'l. Symposium on Computer Architecture (ISCA)*, May 28–31, 1990, Seattle, Wash., 270–279.
- Cantin, J. F., and M. D. Hill [2001]. "Cache performance for selected SPEC CPU2000 benchmarks," <http://www.cs.wisc.edu/multifacet/misc/spec2000cache-data/>.
- Cantin, J., and M. Hill [2003]. "Cache performance for SPEC CPU2000 benchmarks, version 3.0," <http://www.cs.wisc.edu/multifacet/misc/spec2000cache-data/index.html>.
- Case, R. P., and A. Padeys [1978]. "The architecture of the IBM System/370," *Communications of the ACM* 21:1, 73–96. Also appears in D. P. Siewiorek, C. G. Bell, and A. Newell, *Computer Structures: Principles and Examples*, McGraw-Hill, New York, 1982, 830–855.
- Clark, B., T. Deshane, E. Dow, S. Evanchik, M. Finlayson, J. Herne, and J. Neefe Matthews [2004]. "Xen and the art of repeated research," *Proc. USENIX Annual Technical Conf.*, June 27–July 2, 2004, Boston, 1135–1144.
- Clark, D. W. [1983]. "Cache performance of the VAX-11/780," *ACM Trans. on Computer Systems* 1:1, 24–37.
- Clark, D. W., and J. S. Emer [1985]. "Performance of the VAX-11/780 translation buffer: Simulation and measurement," *ACM Trans. on Computer Systems* 3:1 (February), 31–62.
- Compaq Computer Corporation. [1999]. *Compiler Writer's Guide for the Alpha 21264*, Order Number EC-RJ66A-TE, June.
- Conti, C., D. H. Gibson, and S. H. Pitkowsky [1968]. "Structural aspects of the System/360 Model 85. Part I. General organization," *IBM Systems J.* 7:1, 2–14.
- Crawford, J., and P. Gelsinger [1988]. *Programming the 80386*, Sybex, Alameda, Calif.
- Cvetanovic, Z., and R. E. Kessler [2000]. "Performance analysis of the Alpha 21264-based Compaq ES40 system," *Proc. 27th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 10–14, 2000, Vancouver, Canada, 192–202.
- Fabry, R. S. [1974]. "Capability based addressing," *Communications of the ACM* 17:7 (July), 403–412.
- Farkas, K. I., P. Chow, N. P. Jouppi, and Z. Vranesic [1997]. "Memory-system design considerations for dynamically-scheduled processors," *Proc. 24th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 2–4, 1997, Denver, Colo., 133–143.
- Farkas, K. I., and N. P. Jouppi [1994]. "Complexity/performance trade-offs with non-blocking loads," *Proc. 21st Annual Int'l. Symposium on Computer Architecture (ISCA)*, April 18–21, 1994, Chicago.



- Farkas, K. I., N. P. Jouppi, and P. Chow [1995]. "How useful are non-blocking loads, stream buffers and speculative execution in multiple issue processors?" *Proc. First IEEE Symposium on High-Performance Computer Architecture*, January 22–25, 1995, Raleigh, N.C., 78–89.
- Gao, Q. S. [1993]. "The Chinese remainder theorem and the prime memory system," *20th Annual Int'l. Symposium on Computer Architecture (ISCA)*, May 16–19, 1993, San Diego, Calif. (*Computer Architecture News* 21:2 (May), 337–340).
- Gee, J. D., M. D. Hill, D. N. Pnevmatikatos, and A. J. Smith [1993]. "Cache performance of the SPEC92 benchmark suite," *IEEE Micro* 13:4 (August), 17–27.
- Gibson, D. H. [1967]. "Considerations in block-oriented systems design," *AFIPS Conf. Proc.* 30, 75–80.
- Handy, J. [1993]. *The Cache Memory Book*, Academic Press, Boston.
- Heald, R., K. Aingaran, C. Amir, M. Ang, M. Boland, A. Das, P. Dixit, G. Gouldsberry, J. Hart, T. Horel, W.-J. Hsu, J. Kaku, C. Kim, S. Kim, F. Klass, H. Kwan, R. Lo, H. McIntyre, A. Mehta, D. Murata, S. Nguyen, Y.-P. Pai, S. Patel, K. Shin, K. Tam, S. Vishwanthaiiah, J. Wu, G. Yee, and H. You [2000]. "Implementation of third-generation SPARC V9 64-b microprocessor," *ISSCC Digest of Technical Papers*, 412–413 and slide supplement.
- Hill, M. D. [1987]. "Aspects of Cache Memory and Instruction Buffer Performance," Ph.D. thesis, Tech. Rep. UCB/CSD 87/381, Computer Science Division, University of California, Berkeley.
- Hill, M. D. [1988]. "A case for direct mapped caches," *Computer* 21:12 (December), 25–40.
- Horel, T., and G. Lauterbach [1999]. "UltraSPARC-III: Designing third-generation 64-bit performance," *IEEE Micro* 19:3 (May–June), 73–85.
- Hughes, C. J., P. Kaul, S. V. Adve, R. Jain, C. Park, and J. Srinivasan [2001]. "Variability in the execution of multimedia applications and implications for architecture," *Proc. 28th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 30–July 4, 2001, Goteborg, Sweden, 254–265.
- IEEE. [2005]. "Intel virtualization technology, computer," *IEEE Computer Society* 38:5 (May), 48–56.
- Jouppi, N. P. [1990]. "Improving direct-mapped cache performance by the addition of a small fully-associative cache and prefetch buffers," *Proc. 17th Annual Int'l. Symposium on Computer Architecture (ISCA)*, May 28–31, 1990, Seattle, Wash., 364–373.
- Jouppi, N. P. [1998]. "Retrospective: Improving direct-mapped cache performance by the addition of a small fully-associative cache and prefetch buffers," in G. S. Sohi, ed., *25 Years of the International Symposia on Computer Architecture (Selected Papers)*, ACM, New York, 71–73.
- Jouppi, N. P., and S. J. E. Wilton [1994]. "Trade-offs in two-level on-chip caching," *Proc. 21st Annual Int'l. Symposium on Computer Architecture (ISCA)*, April 18–21, 1994, Chicago, 34–45.
- Kessler, R. E. [1999]. "The Alpha 21264 microprocessor," *IEEE Micro* 19:2 (March/April), 24–36.

- Kilburn, T., D. B. G. Edwards, M. J. Lanigan, and F. H. Sumner [1962]. "One-level storage system," *IRE Trans. on Electronic Computers* EC-11 (April) 223–235. Also appears in D. P. Siewiorek, C. G. Bell, and A. Newell, *Computer Structures: Principles and Examples*, McGraw-Hill, New York, 1982, 135–148.
- Kroft, D. [1981]. "Lockup-free instruction fetch/prefetch cache organization," *Proc. Eighth Annual Int'l. Symposium on Computer Architecture (ISCA)*, May 12–14, 1981, Minneapolis, Minn., 81–87.
- Kroft, D. [1998]. "Retrospective: Lockup-free instruction fetch/prefetch cache organization," in G. S. Sohi, ed., *25 Years of the International Symposia on Computer Architecture (Selected Papers)*, ACM, New York, 20–21.
- Kunimatsu, A., N. Ide, T. Sato, Y. Endo, H. Murakami, T. Kamei, M. Hirano, F. Ishihara, H. Tago, M. Oka, A. Ohba, T. Yutaka, T. Okada, and M. Suzuoki [2000]. "Vector unit architecture for emotion synthesis," *IEEE Micro* 20:2 (March–April), 40–47.
- Lam, M. S., E. E. Rothberg, and M. E. Wolf [1991]. "The cache performance and optimizations of blocked algorithms," *Proc. Fourth Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, April 8–11, 1991, Santa Clara, Calif. (*SIGPLAN Notices* 26:4 (April), 63–74).
- Lebeck, A. R., and D. A. Wood [1994]. "Cache profiling and the SPEC benchmarks: A case study," *Computer* 27:10 (October), 15–26.
- Liptay, J. S. [1968]. "Structural aspects of the System/360 Model 85. Part II. The cache," *IBM Systems J.* 7:1, 15–21.
- Luk, C.-K., and T. C. Mowry [1999]. "Automatic compiler-inserted prefetching for pointer-based applications," *IEEE Trans. on Computers*, 48:2 (February), 134–141.
- McCalpin, J. D. [2005]. "STREAM: Sustainable Memory Bandwidth in High Performance Computers," [www.cs.virginia.edu/stream/](http://www.cs.virginia.edu/stream/).
- McFarling, S. [1989]. "Program optimization for instruction caches," *Proc. Third Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, April 3–6, 1989, Boston, 183–191.
- Menon, A., J. Renato Santos, Y. Turner, G. Janakiraman, and W. Zwaenepoel [2005]. "Diagnosing performance overheads in the xen virtual machine environment," *Proc. First ACM/USENIX Int'l. Conf. on Virtual Execution Environments*, June 11–12, 2005, Chicago, 13–23.
- Meyer, R. A., and L. H. Seawright [1970]. "A virtual machine time sharing system," *IBM Systems J.* 9:3, 199–218.
- Mowry, T. C., S. Lam, and A. Gupta [1992]. "Design and evaluation of a compiler algorithm for prefetching," *Proc. Fifth Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, October 12–15, 1992, Boston (*SIGPLAN Notices* 27:9 (September), 62–73).
- Oka, M., and M. Suzuoki [1999]. "Designing and programming the emotion engine," *IEEE Micro* 19:6 (November–December), 20–28.
- Pabst, T. [2000]. "Performance Showdown at 133 MHz FSB—The Best Platform for Coppermine," [www6.tomshardware.com/mainboard/00q1/000302/](http://www6.tomshardware.com/mainboard/00q1/000302/).



- Palacharla, S., and R. E. Kessler [1994]. "Evaluating stream buffers as a secondary cache replacement," *Proc. 21st Annual Int'l. Symposium on Computer Architecture (ISCA)*, April 18–21, 1994, Chicago, 24–33.
- Przybylski, S. A. [1990]. *Cache Design: A Performance-Directed Approach*, Morgan Kaufmann, San Francisco.
- Przybylski, S. A., M. Horowitz, and J. L. Hennessy [1988]. "Performance trade-offs in cache design," *Proc. 15th Annual Int'l. Symposium on Computer Architecture (ISCA)*, May 30–June 2, 1988, Honolulu, Hawaii, 290–298.
- Reinman, G., and N. P. Jouppi. [1999]. "Extensions to CACTI."
- Robin, J., and C. Irvine [2000]. "Analysis of the Intel Pentium's ability to support a secure virtual machine monitor," *Proc. USENIX Security Symposium*, August 14–17, 2000, Denver, Colo.
- Saavedra-Barrera, R. H. [1992]. "CPU Performance Evaluation and Execution Time Prediction Using Narrow Spectrum Benchmarking," Ph.D. dissertation, University of California, Berkeley.
- Samples, A. D., and P. N. Hilfinger [1988]. *Code Reorganization for Instruction Caches*, Tech. Rep. UCB/CSD 88/447, University of California, Berkeley.
- Sites, R. L. (ed.) [1992]. *Alpha Architecture Reference Manual*, Digital Press, Burlington, Mass.
- Skadron, K., and D. W. Clark [1997]. "Design issues and tradeoffs for write buffers," *Proc. Third Int'l. Symposium on High-Performance Computer Architecture*, February 1–5, 1997, San Antonio, Tex., 144–155.
- Smith, A. J. [1982]. "Cache memories," *Computing Surveys* 14:3 (September), 473–530.
- Smith, J. E., and J. R. Goodman [1983]. "A study of instruction cache organizations and replacement policies," *Proc. 10th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 5–7, 1982, Stockholm, Sweden, 132–137.
- Stokes, J. [2000]. "Sound and Vision: A Technical Overview of the Emotion Engine," <http://arstechnica.com/hardware/reviews/2000/02/ee.ars>.
- Strecker, W. D. [1976]. "Cache memories for the PDP-11?" *Proc. Third Annual Int'l. Symposium on Computer Architecture (ISCA)*, January 19–21, 1976, Tampa, Fla., 155–158.
- Sugumar, R. A., and S. G. Abraham [1993]. "Efficient simulation of caches under optimal replacement with applications to miss characterization," *Proc. ACM SIGMETRICS Conf. on Measurement and Modeling of Computer Systems*, May 17–21, 1993, Santa Clara, Calif., 24–35.
- Tarjan, D., S. Thoziyoor, and N. Jouppi [2006]. CACTI 4.0. Technical Report HPL-2006-86, HP Laboratories.
- Torrellas, J., A. Gupta, and J. Hennessy [1992]. "Characterizing the caching and synchronization performance of a multiprocessor operating system," *Proc. Fifth Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, October 12–15, 1992, Boston (*SIGPLAN Notices* 27:9 (September), 162–174).
- Van Vleck, T. [2005]. "The IBM 360/67 and CP/CMS," <http://www.multicians.org/thvv/360-67.html>.

- Wang, W.-H., J.-L. Baer, and H. M. Levy [1989]. "Organization and performance of a two-level virtual-real cache hierarchy," *Proc. 16th Annual Int'l. Symposium on Computer Architecture (ISCA)*, May 28–June 1, 1989, Jerusalem, 140–148.
- Wilkes, M. [1965]. "Slave memories and dynamic storage allocation," *IEEE Trans. Electronic Computers* EC-14:2 (April), 270–271.
- Wilkes, M. V. [1982]. "Hardware support for memory protection: Capability implementations," *Proc. Symposium on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, March 1–3, 1982, Palo Alto, Calif., 107–116.
- Wulf, W. A., R. Levin, and S. P. Harbison [1981]. *Hydra/C.mmp: An Experimental Computer System*, McGraw-Hill, New York.

---

## M.4

## The Evolution of Instruction Sets (Appendices A, J, and K)

*One's eyebrows should rise whenever a future architecture is developed with a stack- or register-oriented instruction set.*

**Meyers [1978, p. 20]**

The earliest computers, including the UNIVAC I, the EDSAC, and the IAS computers, were accumulator-based computers. The simplicity of this type of computer made it the natural choice when hardware resources were very constrained. The first general-purpose register computer was the Pegasus, built by Ferranti, Ltd., in 1956. The Pegasus had eight general-purpose registers, with R0 always being zero. Block transfers loaded the eight registers from the drum memory.

### Stack Architectures

In 1963, Burroughs delivered the B5000. The B5000 was perhaps the first computer to seriously consider software and hardware-software trade-offs. Barton and the designers at Burroughs made the B5000 a stack architecture (as described in Barton [1961]). Designed to support high-level languages such as ALGOL, this stack architecture used an operating system (MCP) written in a high-level language. The B5000 was also the first computer from a U.S. manufacturer to support virtual memory. The B6500, introduced in 1968 (and discussed in Hauck and Dent [1968]), added hardware-managed activation records. In both the B5000 and B6500, the top two elements of the stack were kept in the processor and the rest of the stack was kept in memory. The stack architecture yielded good code density, but only provided two high-speed storage locations. The authors of both the original IBM 360 paper [Amdahl, Blaauw, and Brooks 1964] and the original PDP-11 paper [Bell et al. 1970] argued against the stack organization. They cited three major points in their arguments against stacks:

- Performance is derived from fast registers, not the way they are used.
- The stack organization is too limiting and requires many swap and copy operations.
- The stack has a bottom, and when placed in slower memory there is a performance loss.

Stack-based hardware fell out of favor in the late 1970s and, except for the Intel 80x86 floating-point architecture, essentially disappeared; for example, except for the 80x86, none of the computers listed in the SPEC report uses a stack.

In the 1990s, however, stack architectures received a shot in the arm with the success of the Java Virtual Machine (JVM). The JVM is a software interpreter for an intermediate language produced by Java compilers, called *Java bytecodes* [Lindholm and Yellin 1999]. The purpose of the interpreter is to provide software compatibility across many platforms, with the hope of “write once, run everywhere.” Although the slowdown is about a factor of 10 due to interpretation, there are times when compatibility is more important than performance, such as when downloading a Java “applet” into an Internet browser.

Although a few have proposed hardware to directly execute the JVM instructions (see McGhan and O’Connor [1998]), thus far none of these proposals has been significant commercially. The hope instead is that *just-in-time* (JIT) Java compilers—which compile during runtime to the native instruction set of the computer running the Java program—will overcome the performance penalty of interpretation. The popularity of Java has also led to compilers that compile directly into the native hardware instruction sets, bypassing the illusion of the Java bytecodes.

## Computer Architecture Defined

IBM coined the term *computer architecture* in the early 1960s. Amdahl, Blaauw, and Brooks [1964] used the term to refer to the programmer-visible portion of the IBM 360 instruction set. They believed that a *family* of computers of the same architecture should be able to run the same software. Although this idea may seem obvious to us today, it was quite novel at that time. IBM, although it was the leading company in the industry, had five different architectures before the 360; thus, the notion of a company standardizing on a single architecture was a radical one. The 360 designers hoped that defining a common architecture would bring six different divisions of IBM together. Their definition of architecture was

... the structure of a computer that a machine language programmer must understand to write a correct (timing independent) program for that machine.

The term *machine language programmer* meant that compatibility would hold, even in machine language, while *timing independent* allowed different implementations. This architecture blazed the path for binary compatibility, which others have followed.

The IBM 360 was the first computer to sell in large quantities with both byte addressing using 8-bit bytes and general-purpose registers. The 360 also had register-memory and limited memory-memory instructions. Appendix K summarizes this instruction set.

In 1964, Control Data delivered the first supercomputer, the CDC 6600. As Thornton [1964] discussed, he, Cray, and the other 6600 designers were among the first to explore pipelining in depth. The 6600 was the first general-purpose, load-store computer. In the 1960s, the designers of the 6600 realized the need to simplify architecture for the sake of efficient pipelining. Microprocessor and minicomputer designers largely neglected this interaction between architectural simplicity and implementation during the 1970s, but it returned in the 1980s.

## High-Level Language Computer Architecture

In the late 1960s and early 1970s, people realized that software costs were growing faster than hardware costs. McKeeman [1967] argued that compilers and operating systems were getting too big and too complex and taking too long to develop. Because of inferior compilers and the memory limitations of computers, most systems programs at the time were still written in assembly language. Many researchers proposed alleviating the software crisis by creating more powerful, software-oriented architectures. Tanenbaum [1978] studied the properties of high-level languages. Like other researchers, he found that most programs are simple. He argued that architectures should be designed with this in mind and that they should optimize for program size and ease of compilation. Tanenbaum proposed a stack computer with frequency-encoded instruction formats to accomplish these goals; however, as we have observed, program size does not translate directly to cost-performance, and stack computers faded out shortly after this work.

Strecker's article [1978] discusses how he and the other architects at DEC responded to this by designing the VAX architecture. The VAX was designed to simplify compilation of high-level languages. Compiler writers had complained about the lack of complete orthogonality in the PDP-11. The VAX architecture was designed to be highly orthogonal and to allow the mapping of a high-level language statement into a single VAX instruction. Additionally, the VAX designers tried to optimize code size because compiled programs were often too large for available memories. Appendix K summarizes this instruction set.

The VAX-11/780 was the first computer announced in the VAX series. It is one of the most successful—and most heavily studied—computers ever built. The cornerstone of DEC's strategy was a single architecture, VAX, running a single operating system, VMS. This strategy worked well for over 10 years. The large number of papers reporting instruction mixes, implementation measurements, and analysis of the VAX makes it an ideal case study [Clark and Levy 1982; Wiecek 1982]. Bhandarkar and Clark [1991] gave a quantitative analysis of the disadvantages of the VAX versus a RISC computer, essentially a technical explanation for the demise of the VAX.

While the VAX was being designed, a more radical approach, called *high-level language computer architecture* (HLLCA), was being advocated in the research community. This movement aimed to eliminate the gap between high-level languages and computer hardware—what Gagliardi [1973] called the “semantic gap”—by bringing the hardware “up to” the level of the programming language. Meyers [1982] provided a good summary of the arguments and a history of high-level language computer architecture projects. HLLCA never had a significant commercial impact. The increase in memory size on computers eliminated the code size problems arising from high-level languages and enabled operating systems to be written in high-level languages. The combination of simpler architectures together with software offered greater performance and more flexibility at lower cost and lower complexity.

## Reduced Instruction Set Computers

In the early 1980s, the direction of computer architecture began to swing away from providing high-level hardware support for languages. Ditzel and Patterson [1980] analyzed the difficulties encountered by the high-level language architectures and argued that the answer lay in simpler architectures. In another paper [Patterson and Ditzel 1980], these authors first discussed the idea of Reduced Instruction Set Computers (RISCs) and presented the argument for simpler architectures. Clark and Strecker [1980], who were VAX architects, rebutted their proposal.

The simple load-store computers such as MIPS are commonly called RISC architectures. The roots of RISC architectures go back to computers like the 6600, where Thornton, Cray, and others recognized the importance of instruction set simplicity in building a fast computer. Cray continued his tradition of keeping computers simple in the CRAY-1. Commercial RISCs are built primarily on the work of three research projects: the Berkeley RISC processor, the IBM 801, and the Stanford MIPS processor. These architectures have attracted enormous industrial interest because of claims of a performance advantage of anywhere from two to five times over other computers using the same technology.

Begun in 1975, the IBM project was the first to start but was the last to become public. The IBM computer was designed as a 24-bit ECL minicomputer, while the university projects were both MOS-based, 32-bit microprocessors. John Cocke is considered the father of the 801 design. He received both the Eckert–Mauchly and Turing awards in recognition of his contribution. Radin [1982] described the highlights of the 801 architecture. The 801 was an experimental project that was never designed to be a product. In fact, to keep down costs and complexity, the computer was built with only 24-bit registers.

In 1980, Patterson and his colleagues at Berkeley began the project that was to give this architectural approach its name (see Patterson and Ditzel [1980]). They built two computers called RISC-I and RISC-II. Because the IBM project was not widely known or discussed, the role played by the Berkeley group in promoting the RISC approach was critical to acceptance of the technology. They also built one of

the first instruction caches to support hybrid-format RISCs (see Patterson et al. [1983]). It supported 16-bit and 32-bit instructions in memory but 32 bits in the cache. The Berkeley group went on to build RISC computers targeted toward Smalltalk, described by Ungar et al. [1984], and LISP, described by Taylor et al. [1986].

In 1981, Hennessy and his colleagues at Stanford published a description of the Stanford MIPS computer. Efficient pipelining and compiler-assisted scheduling of the pipeline were both important aspects of the original MIPS design. MIPS stood for Microprocessor without Interlocked Pipeline Stages, reflecting the lack of hardware to stall the pipeline, as the compiler would handle dependencies.

These early RISC computers—the 801, RISC-II, and MIPS—had much in common. Both university projects were interested in designing a simple computer that could be built in VLSI within the university environment. All three computers used a simple load-store architecture and fixed-format 32-bit instructions, and emphasized efficient pipelining. Patterson [1985] described the three computers and the basic design principles that have come to characterize what a RISC computer is, and Hennessy [1984] provided another view of the same ideas, as well as other issues in VLSI processor design.

In 1985, Hennessy published an explanation of the RISC performance advantage and traced its roots to a substantially lower CPI—under 2 for a RISC processor and over 10 for a VAX-11/780 (though not with identical workloads). A paper by Emer and Clark [1984] characterizing VAX-11/780 performance was instrumental in helping the RISC researchers understand the source of the performance advantage seen by their computers.

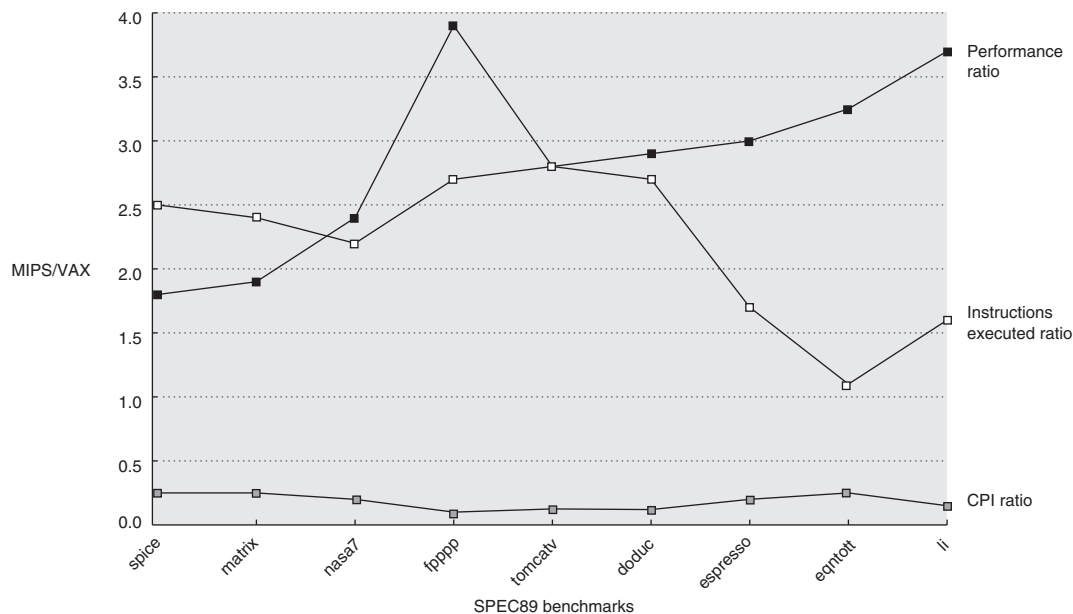
Since the university projects finished up, in the 1983–1984 time frame, the technology has been widely embraced by industry. Many manufacturers of the early computers (those made before 1986) claimed that their products were RISC computers. These claims, however, were often born more of marketing ambition than of engineering reality.

In 1986, the computer industry began to announce processors based on the technology explored by the three RISC research projects. Moussouris et al. [1986] described the MIPS R2000 integer processor, while Kane's book [1986] provides a complete description of the architecture. Hewlett-Packard converted their existing minicomputer line to RISC architectures; Lee [1989] described the HP Precision Architecture. IBM never directly turned the 801 into a product. Instead, the ideas were adopted for a new, low-end architecture that was incorporated in the IBM RT-PC and described in a collection of papers [Waters 1986]. In 1990, IBM announced a new RISC architecture (the RS 6000), which is the first superscalar RISC processor. In 1987, Sun Microsystems began delivering computers based on the SPARC architecture, a derivative of the Berkeley RISC-II processor; SPARC is described in Garner et al. [1988]. The PowerPC joined the forces of Apple, IBM, and Motorola. Appendix K summarizes several RISC architectures.

To help resolve the RISC versus traditional design debate, designers of VAX processors later performed a quantitative comparison of VAX and a RISC processor for implementations with comparable organizations. Their choices were the

VAX 8700 and the MIPS M2000. The differing goals for VAX and MIPS have led to very different architectures. The VAX goals, simple compilers and code density, led to powerful addressing modes, powerful instructions, efficient instruction encoding, and few registers. The MIPS goals were high performance via pipelining, ease of hardware implementation, and compatibility with highly optimizing compilers. These goals led to simple instructions, simple addressing modes, fixed-length instruction formats, and a large number of registers.

Figure M.1 shows the ratio of the number of instructions executed, the ratio of CPIs, and the ratio of performance measured in clock cycles. Since the organizations were similar, clock cycle times were assumed to be the same. MIPS executes about twice as many instructions as the VAX, while the CPI for the VAX is about six times larger than that for the MIPS. Hence, the MIPS M2000 has almost three times the performance of the VAX 8700. Furthermore, much less hardware is needed to build the MIPS processor than the VAX processor. This cost-performance gap is the reason why the company that used to make the VAX introduced a MIPS-based product and then has dropped the VAX completely and switched to Alpha, which is quite similar to MIPS. Bell and Strecker [1998] summarized the debate inside the company. Today, DEC, once the second largest computer company and the major success of the minicomputer industry, exists only as remnants within HP and Intel.



**Figure M.1** Ratio of MIPS M2000 to VAX 8700 in instructions executed and performance in clock cycles using SPEC89 programs. On average, MIPS executes a little over twice as many instructions as the VAX, but the CPI for the VAX is almost six times the MIPS CPI, yielding almost a threefold performance advantage. (Based on data from Bhandarkar and Clark [1991].)



Looking back, only one Complex Instruction Set Computer (CISC) instruction set survived the RISC/CISC debate, and that one had binary compatibility with PC software. The volume of chips is so high in the PC industry that there is a sufficient revenue stream to pay the extra design costs—and sufficient resources due to Moore’s law—to build microprocessors that translate from CISC to RISC internally. Whatever loss in efficiency occurred (due to longer pipeline stages and bigger die size to accommodate translation on the chip) was overcome by the enormous volume and the ability to dedicate IC processing lines specifically to this product.

Interestingly, Intel also concluded that the future of the 80x86 line was doubtful. They created the IA-64 architecture to support 64-bit addressing and to move to a RISC-style instruction set. The embodiment of the IA-64 (see Huck et al. [2000]) architecture in the Itanium-1 and Itanium-2 has been a mixed success. Although high performance has been achieved for floating-point applications, the integer performance was never impressive. In addition, the Itanium implementations have been large in transistor count and die size and power hungry. The complexity of the IA-64 instruction set, standing at least in partial conflict with the RISC philosophy, no doubt contributed to this area and power inefficiency.

AMD decided instead to just stretch the architecture from a 32-bit address to a 64-bit address, much as Intel had done when the 80386 stretched it from a 16-bit address to a 32-bit address. Intel later followed AMD’s example. In the end, the tremendous marketplace advantage of the 80x86 presence was too much even for Intel, the owner of this legacy, to overcome!

## References

- Alexander, W. G., and D. B. Wortman [1975]. “Static and dynamic characteristics of XPL programs,” *IEEE Computer* 8:11 (November), 41–46.
- Amdahl, G. M., G. A. Blaauw, and F. P. Brooks, Jr. [1964]. “Architecture of the IBM System 360,” *IBM J. Research and Development* 8:2 (April), 87–101.
- Barton, R. S. [1961]. “A new approach to the functional design of a computer,” *Proc. Western Joint Computer Conf.*, May 9–11, 1961, Los Angeles, Calif., 393–396.
- Bell, G., R. Cady, H. McFarland, B. DeLagi, J. O’Laughlin, R. Noonan, and W. Wulf [1970]. “A new architecture for mini-computers: The DEC PDP-11,” *Proc. AFIPS SJCC*, May 5–7, 1970, Atlantic City, N.J., 657–675.
- Bell, G., and W. D. Strecker [1998]. “Computer structures: What have we learned from the PDP-11?” in G. S. Sohi, ed., *25 Years of the International Symposium on Computer Architecture (Selected Papers)*, ACM, New York, 138–151.
- Bhandarkar, D. P. [1995]. *Alpha Architecture and Implementations*, Digital Press, Newton, Mass.
- Bhandarkar, D., and D. W. Clark [1991]. “Performance from architecture: Comparing a RISC and a CISC with similar hardware organizations,” *Proc. Fourth*



- Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, April 8–11, 1991, Palo Alto, Calif., 310–319.
- Bier, J. [1997]. “The evolution of DSP processors,” paper presented at University of California, Berkeley, November 14.
- Boddie, J. R. [2000]. “History of DSPs,” [www.lucent.com/micro/dsp/dsphist.html](http://www.lucent.com/micro/dsp/dsphist.html).
- Case, R. P., and A. Padegs [1978]. “The architecture of the IBM System/370,” *Communications of the ACM* 21:1, 73–96.
- Chow, F. C. [1983]. “A Portable Machine-Independent Global Optimizer—Design and Measurements,” Ph.D. thesis, Stanford University, Palo Alto, Calif.
- Clark, D., and H. Levy [1982]. “Measurement and analysis of instruction set use in the VAX-11/780,” *Proc. Ninth Annual Int'l. Symposium on Computer Architecture (ISCA)*, April 26–29, 1982, Austin, Tex., 9–17.
- Clark, D., and W. D. Strecker [1980]. “Comments on ‘the case for the reduced instruction set computer,’” *Computer Architecture News* 8:6 (October), 34–38.
- Crawford, J., and P. Gelsinger [1988]. *Programming the 80386*, Sybex Books, Alameda, Calif.
- Darcy, J. D., and D. Gay [1996]. “FLECKmarks: Measuring floating point performance using a full IEEE compliant arithmetic benchmark,” CS 252 class project, University of California, Berkeley (see <http://www.sonic.net/~jddarcy/Research/fleckmrk.pdf>).
- Digital Semiconductor. [1996]. *Alpha Architecture Handbook, Version 3*, Digital Press, Maynard, Mass.
- Ditzel, D. R., and D. A. Patterson [1980]. “Retrospective on high-level language computer architecture,” *Proc. Seventh Annual Int'l. Symposium on Computer Architecture (ISCA)*, May 6–8, 1980, La Baule, France, 97–104.
- Emer, J. S., and D. W. Clark [1984]. “A characterization of processor performance in the VAX-11/780,” *Proc. 11th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 5–7, 1984, Ann Arbor, Mich., 301–310.
- Furber, S. B. [2000]. *ARM system-on-chip architecture*. Addison-Wesley, Boston, Mass.
- Gagliardi, U. O. [1973]. “Report of workshop 4—software-related advances in computer hardware,” *Proc. Symposium on the High Cost of Software*, September 17–19, 1973, Monterey, Calif., 99–120.
- Game, M., and A. Booker [1999]. “CodePack code compression for PowerPC processors,” *MicroNews*, 5:1.
- Garner, R., A. Agarwal, F. Briggs, E. Brown, D. Hough, B. Joy, S. Kleiman, S. Muchnick, M. Namjoo, D. Patterson, J. Pendleton, and R. Tuck [1988]. “Scalable processor architecture (SPARC),” *Proc. IEEE COMPCON*, February 29–March 4, 1988, San Francisco, 278–283.
- Hauck, E. A., and B. A. Dent [1968]. “Burroughs’ B6500/B7500 stack mechanism,” *Proc. AFIPS SJCC*, April 30–May 2, 1968, Atlantic City, N.J., 245–251.
- Hennessy, J. [1984]. “VLSI processor architecture,” *IEEE Trans. on Computers* C-33:11 (December), 1221–1246.

- Hennessy, J. [1985]. "VLSI RISC processors," *VLSI Systems Design* 6:10 (October), 22–32.
- Hennessy, J., N. Jouppi, F. Baskett, and J. Gill [1981]. "MIPS: A VLSI processor architecture," in *CMU Conference on VLSI Systems and Computations*, Computer Science Press, Rockville, Md.
- Hewlett-Packard. [1994]. *PA-RISC 2.0 Architecture Reference Manual*, 3rd ed., Hewlett-Packard, Palo Alto, Calif.
- Hitachi. [1997]. *SuperH RISC Engine SH7700 Series Programming Manual*, Hitachi, Santa Clara, Calif.
- Huck, J. et al. [2000]. "Introducing the IA-64 Architecture" *IEEE Micro*, 20:5, (September–October), 12–23.
- IBM. [1994]. *The PowerPC Architecture*, Morgan Kaufmann, San Francisco.
- Intel. [2001]. "Using MMX instructions to convert RGB to YUV color conversion," [cedar.intel.com/cgi-bin/ids.dll/content/content.jsp?cntKey=Legacy::irtm\\_AP548\\_9996&cntType=IDS\\_EDITORIAL](http://cedar.intel.com/cgi-bin/ids.dll/content/content.jsp?cntKey=Legacy::irtm_AP548_9996&cntType=IDS_EDITORIAL).
- Kahan, J. [1990]. "On the advantage of the 8087's stack," unpublished course notes, Computer Science Division, University of California, Berkeley.
- Kane, G. [1986]. *MIPS R2000 RISC Architecture*, Prentice Hall, Englewood Cliffs, N.J.
- Kane, G. [1996]. *PA-RISC 2.0 Architecture*, Prentice Hall, Upper Saddle River, N.J.
- Kane, G., and J. Heinrich [1992]. *MIPS RISC Architecture*, Prentice Hall, Englewood Cliffs, N.J.
- Kissell, K. D. [1997]. "MIPS16: High-density for the embedded market," *Proc. Real Time Systems '97*, June 15, 1997, Las Vegas, Nev.
- Kozyrakis, C. [2000]. "Vector IRAM: A media-oriented vector processor with embedded DRAM," paper presented at Hot Chips 12, August 13–15, 2000, Palo Alto, Calif, 13–15.
- Lee, R. [1989]. "Precision architecture," *Computer* 22:1 (January), 78–91.
- Levy, H., and R. Eckhouse [1989]. *Computer Programming and Architecture: The VAX*, Digital Press, Boston.
- Lindholm, T., and F. Yellin [1999]. *The Java Virtual Machine Specification*, 2nd ed., Addison-Wesley, Reading, Mass.
- Lunde, A. [1977]. "Empirical evaluation of some features of instruction set processor architecture," *Communications of the ACM* 20:3 (March), 143–152.
- Magenheimer, D. J., L. Peters, K. W. Pettis, and D. Zuras [1988]. "Integer multiplication and division on the HP precision architecture," *IEEE Trans. on Computers* 37:8, 980–990.
- McGhan, H., and M. O'Connor [1998]. "PicoJava: A direct execution engine for Java bytecode," *Computer* 31:10 (October), 22–30.
- McKeeman, W. M. [1967]. "Language directed computer design," *Proc. AFIPS Fall Joint Computer Conf.*, November 14–16, 1967, Washington, D.C., 413–417.
- Meyers, G. J. [1978]. "The evaluation of expressions in a storage-to-storage architecture," *Computer Architecture News* 7:3 (October), 20–23.

- Meyers, G. J. [1982]. *Advances in Computer Architecture*, 2nd ed., Wiley, New York.
- MIPS. [1997]. *MIPS16 Application Specific Extension Product Description*.
- Mitsubishi. [1996]. *Mitsubishi 32-Bit Single Chip Microcomputer M32R Family Software Manual*, Mitsubishi, Cypress, Calif.
- Morse, S., B. Ravenal, S. Mazor, and W. Pohlman [1980]. "Intel microprocessors—8080 to 8086," *Computer* 13:10 (October).
- Moussouris, J., L. Crudele, D. Freitas, C. Hansen, E. Hudson, S. Przybylski, T. Riordan, and C. Rowen [1986]. "A CMOS RISC processor with integrated system functions," *Proc. IEEE COMPCON*, March 3–6, 1986, San Francisco, 191.
- Muchnick, S. S. [1988]. "Optimizing compilers for SPARC," *Sun Technology* 1:3 (Summer), 64–77.
- Palmer, J., and S. Morse [1984]. *The 8087 Primer*, John Wiley & Sons, New York, 93.
- Patterson, D. [1985]. "Reduced instruction set computers," *Communications of the ACM* 28:1 (January), 8–21.
- Patterson, D. A., and D. R. Ditzel [1980]. "The case for the reduced instruction set computer," *Computer Architecture News* 8:6 (October), 25–33.
- Patterson, D. A., P. Garrison, M. Hill, D. Lioupis, C. Nyberg, T. Sippel, and K. Van Dyke [1983]. "Architecture of a VLSI instruction cache for a RISC," *10th Annual Int'l. Conf. on Computer Architecture Conf. Proc.*, June 13–16, 1983, Stockholm, Sweden, 108–116.
- Radin, G. [1982]. "The 801 minicomputer," *Proc. Symposium Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, March 1–3, 1982, Palo Alto, Calif., 39–47.
- Riemens, A., K. A. Vissers, R. J. Schutten, F. W. Sijstermans, G. J. Hekstra, and G. D. La Hei [1999]. "Trimedia CPU64 application domain and benchmark suite," *Proc. IEEE Int'l. Conf. on Computer Design: VLSI in Computers and Processors (ICCD'99)*, October 10–13, 1999, Austin, Tex., 580–585.
- Ropers, A., H. W. Lollman, and J. Wellhausen [1999]. *DSPstone: Texas Instruments TMS320C54x*, Tech. Rep. Nr. IB 315 1999/9-ISS-Version 0.9, Aachen University of Technology, Aachen, Germany ([www.ert.rwth-aachen.de/Projekte/Tools/coal/dspstone\\_c54x/index.html](http://www.ert.rwth-aachen.de/Projekte/Tools/coal/dspstone_c54x/index.html)).
- Shustek, L. J. [1978]. "Analysis and Performance of Computer Instruction Sets," Ph.D. dissertation, Stanford University, Palo Alto, Calif.
- Silicon Graphics. [1996]. *MIPS V Instruction Set* (see [http://www.sgi.com/MIPS/arch/ISA5/#MIPSV\\_indx](http://www.sgi.com/MIPS/arch/ISA5/#MIPSV_indx)).
- Sites, R. L., and R. Witek, eds. [1995]. *Alpha Architecture Reference Manual*, 2nd ed., Digital Press, Newton, Mass.
- Strauss, W. [1998]. "DSP Strategies 2002," [www.usadata.com/market\\_research/spr\\_05/spr\\_r127-005.htm](http://www.usadata.com/market_research/spr_05/spr_r127-005.htm).
- Strecker, W. D. [1978]. "VAX-11/780: A virtual address extension of the PDP-11 family," *Proc. AFIPS National Computer Conf.*, June 5–8, 1978, Anaheim, Calif., 47, 967–980.

- Sun Microsystems. [1989]. *The SPARC Architectural Manual*, Version 8, Part No. 800-1399-09, Sun Microsystems, Santa Clara, Calif.
- Tanenbaum, A. S. [1978]. "Implications of structured programming for machine architecture," *Communications of the ACM* 21:3 (March), 237–246.
- Taylor, G., P. Hilfinger, J. Larus, D. Patterson, and B. Zorn [1986]. "Evaluation of the SPUR LISP architecture," *Proc. 13th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 2–5, 1986, Tokyo.
- Texas Instruments [2000]. "History of innovation: 1980s," [www.ti.com/corp/docs/company/history/1980s.shtml](http://www.ti.com/corp/docs/company/history/1980s.shtml).
- Thornton, J. E. [1964]. "Parallel operation in Control Data 6600," *Proc. AFIPS Fall Joint Computer Conf., Part II*, October 27–29, 1964, San Francisco, 26, 33–40.
- Ungar, D., R. Blau, P. Foley, D. Samples, and D. Patterson [1984]. "Architecture of SOAR: Smalltalk on a RISC," *Proc. 11th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 5–7, 1984, Ann Arbor, Mich., 188–197.
- van Eijndhoven, J. T. J., F. W. Sijstermans, K. A. Vissers, E. J. D. Pol, M. I. A. Tromp, P. Struik, R. H. J. Bloks, P. van der Wolf, A. D. Pimentel, and H. P. E. Vranken [1999]. "Trimedia CPU64 architecture," *Proc. IEEE Int'l. Conf. on Computer Design: VLSI in Computers and Processors (ICCD'99)*, October 10–13, 1999, Austin, Tex., 586–592.
- Wakerly, J. [1989]. *Microcomputer Architecture and Programming*, Wiley, New York.
- Waters, F. (ed.) [1986]. *IBM RT Personal Computer Technology*, SA 23-1057, IBM, Austin, Tex.
- Weaver, D. L., and T. Germond [1994]. *The SPARC Architectural Manual*, Version 9, Prentice Hall, Englewood Cliffs, N.J.
- Weiss, S., and J. E. Smith [1994]. *Power and PowerPC*, Morgan Kaufmann, San Francisco.
- Wiecek, C. [1982]. "A case study of the VAX 11 instruction set usage for compiler execution," *Proc. Symposium on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, March 1–3, 1982, Palo Alto, Calif., 177–184.
- Wulf, W. [1981]. "Compilers and computer architecture," *Computer* 14:7 (July), 41–47.

---

**M.5**

## **The Development of Pipelining and Instruction-Level Parallelism (Chapter 3 and Appendices C and H)**

### **Early Pipelined CPUs**

The first general-purpose pipelined processor is considered to be Stretch, the IBM 7030. Stretch followed the IBM 704 and had a goal of being 100 times faster than the 704. The goal was a stretch from the state of the art at that time, hence the

nickname. The plan was to obtain a factor of 1.6 from overlapping fetch, decode, and execute, using a four-stage pipeline. Bloch [1959] and Bucholtz [1962] described the design and engineering trade-offs, including the use of ALU bypasses.

A series of general pipelining descriptions that appeared in the late 1970s and early 1980s provided most of the terminology and described most of the basic techniques used in simple pipelines. These surveys include Keller [1975], Ramamoorthy and Li [1977], and Chen [1980], as well as Kogge [1981], whose book is devoted entirely to pipelining. Davidson and his colleagues [1971, 1975] developed the concept of pipeline reservation tables as a design methodology for multicycle pipelines with feedback (also described in Kogge [1981]). Many designers use a variation of these concepts, in either designing pipelines or in creating software to schedule them.

The RISC processors were originally designed with ease of implementation and pipelining in mind. Several of the early RISC papers, published in the early 1980s, attempt to quantify the performance advantages of the simplification in instruction set. The best analysis, however, is a comparison of a VAX and a MIPS implementation published by Bhandarkar and Clark in 1991, 10 years after the first published RISC papers (see Figure M.1). After 10 years of arguments about the implementation benefits of RISC, this paper convinced even the most skeptical designers of the advantages of a RISC instruction set architecture.

J. E. Smith and his colleagues have written a number of papers examining instruction issue, exception handling, and pipeline depth for high-speed scalar CPUs. Kunkel and Smith [1986] evaluated the impact of pipeline overhead and dependences on the choice of optimal pipeline depth; they also provided an excellent discussion of latch design and its impact on pipelining. Smith and Pleszkun [1988] evaluated a variety of techniques for preserving precise exceptions. Weiss and Smith [1984] evaluated a variety of hardware pipeline scheduling and instruction issue techniques.

The MIPS R4000 was one of the first deeply pipelined microprocessors and is described by Killian [1991] and by Heinrich [1993]. The initial Alpha implementation (the 21064) has a similar instruction set and similar integer pipeline structure, with more pipelining in the floating-point unit.

## The Introduction of Dynamic Scheduling

In 1964, CDC delivered the first CDC 6600. The CDC 6600 was unique in many ways. In addition to introducing scoreboarding, the CDC 6600 was the first processor to make extensive use of multiple functional units. It also had peripheral processors that used multithreading. The interaction between pipelining and instruction set design was understood, and a simple, load-store instruction set was used to promote pipelining. The CDC 6600 also used an advanced packaging technology. Thornton [1964] described the pipeline and I/O processor architecture, including the concept of out-of-order instruction execution. Thornton's book [1970] provides an excellent description of the entire processor, from technology

to architecture, and includes a foreword by Cray. (Unfortunately, this book is currently out of print.) The CDC 6600 also has an instruction scheduler for the FORTRAN compilers, described by Thorlin [1967].

### The IBM 360 Model 91: A Landmark Computer

The IBM 360/91 introduced many new concepts, including tagging of data, register renaming, dynamic detection of memory hazards, and generalized forwarding. Tomasulo's algorithm is described in his 1967 paper. Anderson, Sparacio, and Tomasulo [1967] described other aspects of the processor, including the use of branch prediction. Many of the ideas in the 360/91 faded from use for nearly 25 years before being broadly resurrected in the 1990s. Unfortunately, the 360/91 was not successful, and only a handful were sold. The complexity of the design made it late to the market and allowed the Model 85, which was the first IBM processor with a cache, to outperform the 91.

### Branch-Prediction Schemes

The 2-bit dynamic hardware branch-prediction scheme was described by J. E. Smith [1981]. Ditzel and McLellan [1987] described a novel branch-target buffer for CRISP, which implements branch folding. The correlating predictor we examine was described by Pan, So, and Rameh [1992]. Yeh and Patt [1992, 1993] generalized the correlation idea and described multilevel predictors that use branch histories for each branch, similar to the local history predictor used in the 21264. McFarling's tournament prediction scheme, which he refers to as a combined predictor, is described in his 1993 technical report. There are a variety of more recent papers on branch prediction based on variations in the multilevel and correlating predictor ideas. Kaeli and Emma [1991] described return address prediction, and Evers et al. [1998] provided an in-depth analysis of multilevel predictors. The data shown in Chapter 3 are from Skadron et al. [1999]. There are several schemes for prediction that may offer some additional benefit beyond tournament predictors. Eden and Mudge [1998] and Jimenez and Lin [2002] have described such approaches.

### The Development of Multiple-Issue Processors

IBM did pioneering work on multiple issue. In the 1960s, a project called ACS was under way in California. It included multiple-issue concepts, a proposal for dynamic scheduling (although with a simpler mechanism than Tomasulo's scheme, which used backup registers), and fetching down both branch paths. The project originally started as a new architecture to follow Stretch and surpass the CDC 6600/6800. ACS started in New York but was moved to California, later changed to be S/360 compatible, and eventually canceled. John Cocke was one of the intellectual forces behind the team that included a number of IBM veterans and younger contributors,



many of whom went on to other important roles in IBM and elsewhere: Jack Bertram, Ed Sussenguth, Gene Amdahl, Herb Schorr, Fran Allen, Lynn Conway, and Phil Dauber, among others. While the compiler team published many of their ideas and had great influence outside IBM, the architecture ideas were not widely disseminated at that time. The most complete accessible documentation of this important project is at [www.cs.clemson.edu/~mark/acs.html](http://www.cs.clemson.edu/~mark/acs.html), which includes interviews with the ACS veterans and pointers to other sources. Sussenguth [1999] is a good overview of ACS.

Most of the early multiple-issue processors that actually reached the market followed an LIW or VLIW design approach. Charlesworth [1981] reported on the Floating Point Systems AP-120B, one of the first wide-instruction processors containing multiple operations per instruction. Floating Point Systems applied the concept of software pipelining both in a compiler and by handwriting assembly language libraries to use the processor efficiently. Because the processor was an attached processor, many of the difficulties of implementing multiple issue in general-purpose processors (for example, virtual memory and exception handling) could be ignored.

One of the interesting approaches used in early VLIW processors, such as the AP-120B and i860, was the idea of a pipeline organization that requires operations to be “pushed through” a functional unit and the results to be caught at the end of the pipeline. In such processors, operations advance only when another operation pushes them from behind (in sequence). Furthermore, an instruction specifies the destination for an instruction issued earlier that will be pushed out of the pipeline when this new operation is pushed in. Such an approach has the advantage that it does not specify a result destination when an operation first issues but only when the result register is actually written. This separation eliminates the need to detect write after write (WAW) and write after read (WAR) hazards in the hardware. The disadvantage is that it increases code size since no-ops may be needed to push results out when there is a dependence on an operation that is still in the pipeline and no other operations of that type are immediately needed. Instead of the “push-and-catch” approach used in these two processors, almost all designers have chosen to use *self-draining pipelines* that specify the destination in the issuing instruction and in which an issued instruction will complete without further action. The advantages in code density and simplifications in code generation seem to outweigh the advantages of the more unusual structure.

Several research projects introduced some form of multiple issue in the mid-1980s. For example, the Stanford MIPS processor had the ability to place two operations in a single instruction, although this capability was dropped in commercial variants of the architecture, primarily for performance reasons. Along with his colleagues at Yale, Fisher [1983] proposed creating a processor with a very wide instruction (512 bits) and named this type of processor a VLIW. Code was generated for the processor using trace scheduling, which Fisher [1981] had developed originally for generating horizontal microcode. The implementation of trace scheduling for the Yale processor is described by Fisher et al. [1984] and by Ellis [1986].

Although IBM canceled ACS, active research in the area continued in the 1980s. More than 10 years after ACS was canceled, John Cocke made a new proposal for a superscalar processor that dynamically made issue decisions; he and Tilak Agerwala described the key ideas in several talks in the mid-1980s and coined the term *superscalar*. He called the design America; it is described by Agerwala and Cocke [1987]. The IBM Power1 architecture (the RS/6000 line) is based on these ideas (see Bakoglu et al. [1989]).

J. E. Smith [1984] and his colleagues at Wisconsin proposed the decoupled approach that included multiple issue with limited dynamic pipeline scheduling. A key feature of this processor is the use of queues to maintain order among a class of instructions (such as memory references) while allowing it to slip behind or ahead of another class of instructions. The Astronautics ZS-1 described by Smith et al. [1987] embodies this approach with queues to connect the load-store unit and the operation units. The Power2 design uses queues in a similar fashion. J. E. Smith [1989] also described the advantages of dynamic scheduling and compared that approach to static scheduling.

The concept of speculation has its roots in the original 360/91, which performed a very limited form of speculation. The approach used in recent processors combines the dynamic scheduling techniques of the 360/91 with a buffer to allow in-order commit. Smith and Pleszkun [1988] explored the use of buffering to maintain precise interrupts and described the concept of a reorder buffer. Sohi [1990] described adding renaming and dynamic scheduling, making it possible to use the mechanism for speculation. Patt and his colleagues were early proponents of aggressive reordering and speculation. They focused on checkpoint and restart mechanisms and pioneered an approach called HPSm, which is also an extension of Tomasulo's algorithm [Hwu and Patt 1986].

The use of speculation as a technique in multiple-issue processors was evaluated by Smith, Johnson, and Horowitz [1989] using the reorder buffer technique; their goal was to study available ILP in nonscientific code using speculation and multiple issue. In a subsequent book, Johnson [1990] described the design of a speculative superscalar processor. Johnson later led the AMD K-5 design, one of the first speculative superscalars.

In parallel with the superscalar developments, commercial interest in VLIW approaches also increased. The Multiflow processor (see Colwell et al. [1987]) was based on the concepts developed at Yale, although many important refinements were made to increase the practicality of the approach. Among these was a control-lable store buffer that provided support for a form of speculation. Although more than 100 Multiflow processors were sold, a variety of problems, including the difficulties of introducing a new instruction set from a small company and competition from commercial RISC microprocessors that changed the economics in the mini-computer market, led to the failure of Multiflow as a company.

Around the same time as Multiflow, Cydrome was founded to build a VLIW-style processor (see Rau et al. [1989]), which was also unsuccessful commercially. Dehnert, Hsu, and Bratt [1989] explained the architecture and performance of the



Cydrome Cydra 5, a processor with a wide-instruction word that provides dynamic register renaming and additional support for software pipelining. The Cydra 5 is a unique blend of hardware and software, including conditional instructions and register rotation, aimed at extracting ILP. Cydrome relied on more hardware than the Multiflow processor and achieved competitive performance primarily on vector-style codes. In the end, Cydrome suffered from problems similar to those of Multiflow and was not a commercial success. Both Multiflow and Cydrome, although unsuccessful as commercial entities, produced a number of people with extensive experience in exploiting ILP as well as advanced compiler technology; many of those people have gone on to incorporate their experience and the pieces of the technology in newer processors. Fisher and Rau [1993] edited a comprehensive collection of papers covering the hardware and software of these two important processors.

Rau had also developed a scheduling technique called *polycyclic scheduling*, which is a basis for most software-pipelining schemes (see Rau, Glaeser, and Picard [1982]). Rau's work built on earlier work by Davidson and his colleagues on the design of optimal hardware schedulers for pipelined processors. Other historical LIW processors have included the Apollo DN 10000 and the Intel i860, both of which could dual-issue FP and integer operations.

## Compiler Technology and Hardware Support for Scheduling

Loop-level parallelism and dependence analysis were developed primarily by D. Kuck and his colleagues at the University of Illinois in the 1970s. They also coined the commonly used terminology of *antidependence* and *output dependence* and developed several standard dependence tests, including the GCD and Banerjee tests. The latter test was named after Uptal Banerjee and comes in a variety of flavors. Recent work on dependence analysis has focused on using a variety of exact tests ending with a linear programming algorithm called Fourier–Motzkin. D. Maydan and W. Pugh both showed that the sequences of exact tests were a practical solution.

In the area of uncovering and scheduling ILP, much of the early work was connected to the development of VLIW processors, described earlier. Lam [1988] developed algorithms for software pipelining and evaluated their use on Warp, a wide-instruction-word processor designed for special-purpose applications. Weiss and Smith [1987] compared software pipelining versus loop unrolling as techniques for scheduling code on a pipelined processor. Rau [1994] developed modulo scheduling to deal with the issues of software-pipelining loops and simultaneously handling register allocation.

Support for speculative code scheduling was explored in a variety of contexts, including several processors that provided a mode in which exceptions were ignored, allowing more aggressive scheduling of loads (e.g., the MIPS TFP processor [Hsu 1994]). Several groups explored ideas for more aggressive hardware support for speculative code scheduling. For example, Smith, Horowitz, and Lam

[1992] created a concept called boosting that contains a hardware facility for supporting speculation but provides a checking and recovery mechanism, similar to those in IA-64 and Crusoe. The sentinel scheduling idea, which is also similar to the speculate-and-check approach used in both Crusoe and the IA-64 architectures, was developed jointly by researchers at the University of Illinois and HP Laboratories (see Mahlke et al. [1992]).

In the early 1990s, Wen-Mei Hwu and his colleagues at the University of Illinois developed a compiler framework, called IMPACT (see Chang et al. [1991]), for exploring the interaction between multiple-issue architectures and compiler technology. This project led to several important ideas, including superblock scheduling (see Hwu et al. [1993]), extensive use of profiling for guiding a variety of optimizations (e.g., procedure inlining), and the use of a special buffer (similar to the ALAT or program-controlled store buffer) for compile-aided memory conflict detection (see Gallagher et al. [1994]). They also explored the performance trade-offs between partial and full support for predication in Mahlke et al. [1995].

The early RISC processors all had delayed branches, a scheme inspired from microprogramming, and several studies on compile time branch prediction were inspired by delayed branch mechanisms. McFarling and Hennessy [1986] did a quantitative comparison of a variety of compile time and runtime branch-prediction schemes. Fisher and Freudenberger [1992] evaluated a range of compile time branch-prediction schemes using the metric of distance between mispredictions. Ball and Larus [1993] and Calder et al. [1997] described static prediction schemes using collected program behavior.

## EPIC and the IA-64 Development

The roots of the EPIC approach lie in earlier attempts to build LIW and VLIW machines—especially those at Cydrome and Multiflow—and in a long history of compiler work that continued after these companies failed at HP, the University of Illinois, and elsewhere. Insights gained from that work led designers at HP to propose a VLIW-style, 64-bit architecture to follow the HP PA RISC architecture. Intel was looking for a new architecture to replace the x86 (now called IA-32) architecture and to provide 64-bit capability. In 1995, they formed a partnership to design a new architecture, IA-64 (see Huck et al. [2000]), and build processors based on it. Itanium (see Sharangpani and Arora [2000]) is the first such processor. In 2002, Intel introduced the second-generation IA-64 design, the Itanium 2 (see McNairy and Soltis [2003] and McCormick and Knies [2002]).

## Studies of ILP and Ideas to Increase ILP

A series of early papers, including Tjaden and Flynn [1970] and Riseman and Foster [1972], concluded that only small amounts of parallelism could be available at the instruction level without investing an enormous amount of hardware. These papers dampened the appeal of multiple instruction issue for more than 10 years.

Nicolau and Fisher [1984] published a paper based on their work with trace scheduling and asserted the presence of large amounts of potential ILP in scientific programs.

Since then there have been many studies of the available ILP. Such studies have been criticized because they presume some level of both hardware support and compiler technology. Nonetheless, the studies are useful to set expectations as well as to understand the sources of the limitations. Wall has participated in several such studies, including Jouppi and Wall [1989] and Wall [1991, 1993]. Although the early studies were criticized as being conservative (e.g., they didn't include speculation), the last study is by far the most ambitious study of ILP to date and the basis for the data in Section 3.10. Sohi and Vajapeyam [1989] provided measurements of available parallelism for wide-instruction-word processors. Smith, Johnson, and Horowitz [1989] also used a speculative superscalar processor to study ILP limits. At the time of their study, they anticipated that the processor they specified was an upper bound on reasonable designs. Recent and upcoming processors, however, are likely to be at least as ambitious as their processor. Skadron et al. [1999] examined the performance trade-offs and limitations in a processor comparable to the most aggressive processors in 2005, concluding that the larger window sizes will not make sense without significant improvements on branch prediction for integer programs.

Lam and Wilson [1992] looked at the limitations imposed by speculation and showed that additional gains are possible by allowing processors to speculate in multiple directions, which requires more than one PC. (Such schemes cannot exceed what perfect speculation accomplishes, but they help close the gap between realistic prediction schemes and perfect prediction.) Wall's 1993 study includes a limited evaluation of this approach (up to eight branches are explored).

## Going Beyond the Data Flow Limit

One other approach that has been explored in the literature is the use of value prediction. Value prediction can allow speculation based on data values. There have been a number of studies of the use of value prediction. Lipasti and Shen published two papers in 1996 evaluating the concept of value prediction and its potential impact on ILP exploitation. Calder, Reinman, and Tullsen [1999] explored the idea of selective value prediction. Sodani and Sohi [1997] approached the same problem from the viewpoint of reusing the values produced by instructions. Moshovos et al. [1997] showed that deciding when to speculate on values, by tracking whether such speculation has been accurate in the past, is important to achieving performance gains with value speculation. Moshovos and Sohi [1997] and Chrysos and Emer [1998] focused on predicting memory dependences and using this information to eliminate the dependence through memory. González and González [1998], Babbay and Mendelson [1998], and Calder, Reinman, and Tullsen [1999] are more recent studies of the use of value prediction. This area is currently highly active, with new results being published in every conference.

## Recent Advanced Microprocessors

The years 1994 and 1995 saw the announcement of wide superscalar processors (three or more issues per clock) by every major processor vendor: Intel Pentium Pro and Pentium II (these processors share the same core pipeline architecture, described by Colwell and Steck [1995]); AMD K-5, K-6, and Athlon; Sun UltraSPARC (see Lauterbach and Horel [1999]); Alpha 21164 (see Edmondson et al. [1995]) and 21264 (see Kessler [1999]); MIPS R10000 and R12000 (see Yeager [1996]); PowerPC 603, 604, and 620 (see Diep, Nelson, and Shen [1995]); and HP 8000 (Kumar [1997]). The latter part of the decade (1996–2000) saw second generations of many of these processors (Pentium III, AMD Athlon, and Alpha 21264, among others). The second generation, although similar in issue rate, could sustain a lower CPI and provided much higher clock rates. All included dynamic scheduling, and they almost universally supported speculation. In practice, many factors, including the implementation technology, the memory hierarchy, the skill of the designers, and the type of applications benchmarked, all play a role in determining which approach is best.

The period from 2000 to 2005 was dominated by three trends among superscalar processors: the introduction of higher clock rates achieved through deeper pipelining (e.g., in the Pentium 4; see Hinton et al. [2001]), the introduction of multithreading by IBM in the Power 4 and by Intel in the Pentium 4 Extreme, and the beginning of the movement to multicore by IBM in the Power 4, AMD in Opteron (see Keltcher et al. [2003]), and most recently by Intel (see Douglas [2005]).

## Multithreading and Simultaneous Multithreading

The concept of multithreading dates back to one of the earliest transistorized computers, the TX-2. TX-2 is also famous for being the computer on which Ivan Sutherland created Sketchpad, the first computer graphics system. TX-2 was built at MIT's Lincoln Laboratory and became operational in 1959. It used multiple threads to support fast context switching to handle I/O functions. Clark [1957] described the basic architecture, and Forgie [1957] described the I/O architecture. Multithreading was also used in the CDC 6600, where a fine-grained multithreading scheme with interleaved scheduling among threads was used as the architecture of the I/O processors. The HEP processor, a pipelined multiprocessor designed by Denelcor and shipped in 1982, used fine-grained multithreading to hide the pipeline latency as well as to hide the latency to a large memory shared among all the processors. Because the HEP had no cache, this hiding of memory latency was critical. Burton Smith, one of the primary architects, described the HEP architecture in a 1978 paper, and Jordan [1983] published a performance evaluation. The TERA processor extends the multithreading ideas and is described by Alverson et al. in a 1992 paper. The Niagara multithreading approach is similar to those of the HEP and TERA systems, although Niagara employs caches reducing the need for thread-based latency hiding.

In the late 1980s and early 1990s, researchers explored the concept of coarse-grained multithreading (also called *block multithreading*) as a way to tolerate

latency, especially in multiprocessor environments. The SPARCLE processor in the Alewife system used such a scheme, switching threads whenever a highlatency exceptional event, such as a long cache miss, occurred. Agarwal et al. described SPARCLE in a 1993 paper. The IBM Pulsar processor uses similar ideas.

By the early 1990s, several research groups had arrived at two key insights. First, they realized that fine-grained multithreading was needed to get the maximum performance benefit, since in a coarse-grained approach, the overhead of thread switching and thread start-up (e.g., filling the pipeline from the new thread) negated much of the performance advantage (see Laudon, Gupta, and Horowitz [1994]). Second, several groups realized that to effectively use large numbers of functional units would require both ILP and thread-level parallelism (TLP). These insights led to several architectures that used combinations of multithreading and multiple issue. Wolfe and Shen [1991] described an architecture called XIMD that statically interleaves threads scheduled for a VLIW processor. Hirata et al. [1992] described a proposed processor for media use that combines a static superscalar pipeline with support for multithreading; they reported speed-ups from combining both forms of parallelism. Keckler and Dally [1992] combined static scheduling of ILP and dynamic scheduling of threads for a processor with multiple functional units. The question of how to balance the allocation of functional units between ILP and TLP and how to schedule the two forms of parallelism remained open.

When it became clear in the mid-1990s that dynamically scheduled superscalars would be delivered shortly, several research groups proposed using the dynamic scheduling capability to mix instructions from several threads on the fly. Yamamoto et al. [1994] appear to have published the first such proposal, though the simulation results for their multithreaded superscalar architecture use simplistic assumptions. This work was quickly followed by Tullsen, Eggers, and Levy [1995], who provided the first realistic simulation assessment and coined the term *simultaneous multithreading*. Subsequent work by the same group together with industrial coauthors addressed many of the open questions about SMT. For example, Tullsen et al. [1996] addressed questions about the challenges of scheduling ILP versus TLP. Lo et al. [1997] provided an extensive discussion of the SMT concept and an evaluation of its performance potential, and Lo et al. [1998] evaluated database performance on an SMT processor. Tuck and Tullsen [2003] reviewed the performance of SMT on the Pentium 4.

The IBM Power4 introduced multithreading (see Tendler et al. [2002]), while the Power5 used simultaneous multithreading. Mathis et al. [2005] explored the performance of SMT in the Power5, while Sinharoy et al. [2005] described the system architecture.

## References

- Agarwal, A., J. Kubiawicz, D. Kranz, B.-H. Lim, D. Yeung, G. D'Souza, and M. Parkin [1993]. "Sparcle: An evolutionary processor design for large-scale multiprocessors," *IEEE Micro* 13 (June), 48–61.

- Agerwala, T., and J. Cocke [1987]. *High Performance Reduced Instruction Set Processors*, Tech. Rep. RC12434, IBM Thomas Watson Research Center, Yorktown Heights, N.Y.
- Alverson, G., R. Alverson, D. Callahan, B. Koblenz, A. Porterfield, and B. Smith [1992]. "Exploiting heterogeneous parallelism on a multithreaded multiprocessor," *Proc. ACM/IEEE Conf. on Supercomputing*, November 16–20, 1992, Minneapolis, Minn., 188–197.
- Anderson, D. W., F. J. Sparacio, and R. M. Tomasulo [1967]. "The IBM 360 Model 91: Processor philosophy and instruction handling," *IBM J. Research and Development* 11:1 (January), 8–24.
- Austin, T. M., and G. Sohi [1992]. "Dynamic dependency analysis of ordinary programs," *Proc. 19th Annual Int'l. Symposium on Computer Architecture (ISCA)*, May 19–21, 1992, Gold Coast, Australia, 342–351.
- Babbay, F., and A. Mendelson [1998]. "Using value prediction to increase the power of speculative execution hardware," *ACM Trans. on Computer Systems* 16:3 (August), 234–270.
- Bakoglu, H. B., G. F. Grohoski, L. E. Thatcher, J. A. Kaeli, C. R. Moore, D. P. Tattle, W. E. Male, W. R. Hardell, D. A. Hicks, M. Nguyen Phu, R. K. Montoye, W. T. Glover, and S. Dhawan [1989]. "IBM second-generation RISC processor organization," *Proc. IEEE Int'l. Conf. on Computer Design*, October, Rye Brook, N.Y., 138–142.
- Ball, T., and J. Larus [1993]. "Branch prediction for free," *Proc. ACM SIGPLAN'93 Conference on Programming Language Design and Implementation (PLDI)*, June 23–25, 1993, Albuquerque, N.M., 300–313.
- Bhandarkar, D., and D. W. Clark [1991]. "Performance from architecture: Comparing a RISC and a CISC with similar hardware organizations," *Proc. Fourth Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, April 8–11, 1991, Palo Alto, Calif., 310–319.
- Bhandarkar, D., and J. Ding [1997]. "Performance characterization of the Pentium Pro processor," *Proc. Third Int'l. Symposium on High Performance Computer Architecture*, February 1–5, 1997, San Antonio, Tex., 288–297.
- Bloch, E. [1959]. "The engineering design of the Stretch computer," *Proc. Eastern Joint Computer Conf.*, December 1–3, 1959, Boston, Mass., 48–59.
- Bucholtz, W. [1962]. *Planning a Computer System: Project Stretch*, McGraw-Hill, New York.
- Calder, B., D. Grunwald, M. Jones, D. Lindsay, J. Martin, M. Mozer, and B. Zorn [1997]. "Evidence-based static branch prediction using machine learning," *ACM Trans. Program. Lang. Syst.* 19:1, 188–222.
- Calder, B., G. Reinman, and D. M. Tullsen [1999]. "Selective value prediction," *Proc. 26th Annual Int'l. Symposium on Computer Architecture (ISCA)*, May 2–4, 1999, Atlanta, Ga.
- Chang, P. P., S. A. Mahlke, W. Y. Chen, N. J. Warter, and W. W. Hwu [1991]. "IMPACT: An architectural framework for multiple-instruction-issue processors," *Proc. 18th Annual Int'l. Symposium on Computer Architecture (ISCA)*, May 27–30, 1991, Toronto, Canada, 266–275.



- Charlesworth, A. E. [1981]. "An approach to scientific array processing: The architecture design of the AP-120B/FPS-164 family," *Computer* 14:9 (September), 18–27.
- Chen, T. C. [1980]. "Overlap and parallel processing," in *Introduction to Computer Architecture*, H. Stone, ed., Science Research Associates, Chicago, 427–486.
- Chrysos, G. Z., and J. S. Emer [1998]. "Memory dependence prediction using store sets," *Proc. 25th Annual Int'l. Symposium on Computer Architecture (ISCA)*, July 3–14, 1998, Barcelona, Spain, 142–153.
- Clark, D. W. [1987]. "Pipelining and performance in the VAX 8800 processor," *Proc. Second Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, October 5–8, 1987, Palo Alto, Calif., 173–177.
- Clark, W. A. [1957]. "The Lincoln TX-2 computer development," *Proc. Western Joint Computer Conference*, February 26–28, 1957, Los Angeles, 143–145.
- Colwell, R. P., and R. Steck [1995]. "A 0.6  $\mu\text{m}$  BiCMOS processor with dynamic execution," *Proc. of IEEE Int'l. Symposium on Solid State Circuits (ISSCC)*, February 15–17, 1995, San Francisco, 176–177.
- Colwell, R. P., R. P. Nix, J. J. O'Donnell, D. B. Papworth, and P. K. Rodman [1987]. "A VLIW architecture for a trace scheduling compiler," *Proc. Second Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, October 5–8, 1987, Palo Alto, Calif., 180–192.
- Cvetanovic, Z., and R. E. Kessler [2000]. "Performance analysis of the Alpha 21264-based Compaq ES40 system," *27th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 10–14, 2000, Vancouver, Canada, 192–202.
- Davidson, E. S. [1971]. "The design and control of pipelined function generators," *Proc. IEEE Conf. on Systems, Networks, and Computers*, January 19–21, 1971, Oaxtepec, Mexico, 19–21.
- Davidson, E. S., A. T. Thomas, L. E. Shar, and J. H. Patel [1975]. "Effective control for pipelined processors," *Proc. IEEE COMPCON*, February 25–27, 1975, San Francisco, 181–184.
- Dehnert, J. C., P. Y.-T. Hsu, and J. P. Bratt [1989]. "Overlapped loop support on the Cydra 5," *Proc. Third Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, April 3–6, 1989, Boston, Mass., 26–39.
- Diep, T. A., C. Nelson, and J. P. Shen [1995]. "Performance evaluation of the PowerPC 620 microarchitecture," *Proc. 22nd Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 22–24, 1995, Santa Margherita, Italy.
- Ditzel, D. R., and H. R. McLellan [1987]. "Branch folding in the CRISP microprocessor: Reducing the branch delay to zero," *Proc. 14th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 2–5, 1987, Pittsburgh, Penn., 2–7.
- Douglas, J. [2005]. "Intel 8xx series and Paxville Xeon-MP Microprocessors," paper presented at Hot Chips 17, August 14–16, 2005, Stanford University, Palo Alto, Calif.
- Eden, A., and T. Mudge [1998]. "The YAGS branch prediction scheme," *Proc. of the 31st Annual ACM/IEEE Int'l. Symposium on Microarchitecture*, November 30–December 2, 1998, Dallas, Tex., 69–80.

- Edmondson, J. H., P. I. Rubinfeld, R. Preston, and V. Rajagopalan [1995]. "Superscalar instruction execution in the 21164 Alpha microprocessor," *IEEE Micro* 15:2, 33–43.
- Ellis, J. R. [1986]. *Bulldog: A Compiler for VLIW Architectures*, MIT Press, Cambridge, Mass.
- Emer, J. S., and D. W. Clark [1984]. "A characterization of processor performance in the VAX-11/780," *Proc. 11th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 5–7, 1984, Ann Arbor, Mich., 301–310.
- Evers, M., S. J. Patel, R. S. Chappell, and Y. N. Patt [1998]. "An analysis of correlation and predictability: What makes two-level branch predictors work," *Proc. 25th Annual Int'l. Symposium on Computer Architecture (ISCA)*, July 3–14, 1998, Barcelona, Spain, 52–61.
- Fisher, J. A. [1981]. "Trace scheduling: A technique for global microcode compaction," *IEEE Trans. on Computers* 30:7 (July), 478–490.
- Fisher, J. A. [1983]. "Very long instruction word architectures and ELI-512," *10th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 5–7, 1982, Stockholm, Sweden, 140–150.
- Fisher, J. A., and S. M. Freudenberger [1992]. "Predicting conditional branches from previous runs of a program," *Proc. Fifth Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, October 12–15, 1992, Boston, 85–95.
- Fisher, J. A., and B. R. Rau [1993]. *Journal of Supercomputing*, January (special issue).
- Fisher, J. A., J. R. Ellis, J. C. Ruttenberg, and A. Nicolau [1984]. "Parallel processing: A smart compiler and a dumb processor," *Proc. SIGPLAN Conf. on Compiler Construction*, June 17–22, 1984, Montreal, Canada, 11–16.
- Forgie, J. W. [1957]. "The Lincoln TX-2 input-output system," *Proc. Western Joint Computer Conference*, February 26–28, 1957, Los Angeles, 156–160.
- Foster, C. C., and E. M. Riseman [1972]. "Percolation of code to enhance parallel dispatching and execution," *IEEE Trans. on Computers* C-21:12 (December), 1411–1415.
- Gallagher, D. M., W. Y. Chen, S. A. Mahlke, J. C. Gyllenhaal, and W. W. Hwu [1994]. "Dynamic memory disambiguation using the memory conflict buffer," *Proc. Sixth Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, October 4–7, Santa Jose, Calif., 183–193.
- González, J., and A. González [1998]. "Limits of instruction level parallelism with data speculation," *Proc. Vector and Parallel Processing (VECPAR) Conf.*, June 21–23, 1998, Porto, Portugal, 585–598.
- Heinrich, J. [1993]. *MIPS R4000 User's Manual*, Prentice Hall, Englewood Cliffs, N.J.
- Hinton, G., D. Sager, M. Upton, D. Boggs, D. Carmean, A. Kyker, and P. Roussel [2001]. "The microarchitecture of the Pentium 4 processor," *Intel Technology Journal*, February.
- Hirata, H., K. Kimura, S. Nagamine, Y. Mochizuki, A. Nishimura, Y. Nakase, and T. Nishizawa [1992]. "An elementary processor architecture with simultaneous instruction issuing from multiple threads," *Proc. 19th Annual Int'l. Symposium*



- on *Computer Architecture (ISCA)*, May 19–21, 1992, Gold Coast, Australia, 136–145.
- Hopkins, M. [2000]. “A critical look at IA-64: Massive resources, massive ILP, but can it deliver?” *Microprocessor Report*, February.
- Hsu, P. [1994]. “Designing the TFP microprocessor,” *IEEE Micro* 18:2 (April), 2333.
- Huck, J. et al. [2000]. “Introducing the IA-64 Architecture” *IEEE Micro*, 20:5 (September–October), 12–23.
- Hwu, W.-M., and Y. Patt [1986]. “HPSm, a high performance restricted data flow architecture having minimum functionality,” *13th Annual Int’l. Symposium on Computer Architecture (ISCA)*, June 2–5, 1986, Tokyo, 297–307.
- Hwu, W. W., S. A. Mahlke, W. Y. Chen, P. P. Chang, N. J. Warter, R. A. Bringmann, R. O. Ouellette, R. E. Hank, T. Kiyohara, G. E. Haab, J. G. Holm, and D. M. Lavery [1993]. “The superblock: An effective technique for VLIW and superscalar compilation,” *J. Supercomputing* 7:1, 2 (March), 229–248.
- IBM. [1990]. “The IBM RISC System/6000 processor” (collection of papers), *IBM J. Research and Development* 34:1 (January).
- Jimenez, D. A., and C. Lin [2002]. “Neural methods for dynamic branch prediction,” *ACM Trans. Computer Sys* 20:4 (November), 369–397.
- Johnson, M. [1990]. *Superscalar Microprocessor Design*, Prentice Hall, Englewood Cliffs, N.J.
- Jordan, H. F. [1983]. “Performance measurements on HEP—a pipelined MIMD computer,” *Proc. 10th Annual Int’l. Symposium on Computer Architecture (ISCA)*, June 5–7, 1982, Stockholm, Sweden, 207–212.
- Jouppi, N. P., and D. W. Wall [1989]. “Available instruction-level parallelism for superscalar and superpipelined processors,” *Proc. Third Int’l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, April 3–6, 1989, Boston, 272–282.
- Kaeli, D. R., and P. G. Emma [1991]. “Branch history table prediction of moving target branches due to subroutine returns,” *Proc. 18th Annual Int’l. Symposium on Computer Architecture (ISCA)*, May 27–30, 1991, Toronto, Canada, 34–42.
- Keckler, S. W., and W. J. Dally [1992]. “Processor coupling: Integrating compile time and runtime scheduling for parallelism,” *Proc. 19th Annual Int’l. Symposium on Computer Architecture (ISCA)*, May 19–21, 1992, Gold Coast, Australia, 202–213.
- Keller, R. M. [1975]. “Look-ahead processors,” *ACM Computing Surveys* 7:4 (December), 177–195.
- Keltcher, C. N., K. J. McGrath, A. Ahmed, and P. Conway [2003]. “The AMD Opteron processor for multiprocessor servers,” *IEEE Micro* 23:2 (March–April), 66–76.
- Kessler, R. [1999]. “The Alpha 21264 microprocessor,” *IEEE Micro* 19:2 (March/April) 24–36.
- Killian, E. [1991]. “MIPS R4000 technical overview—64 bits/100 MHz or bust,” *Hot Chips III Symposium Record*, August 26–27, 1991, Stanford University, Palo Alto, Calif., 1.6–1.19.

- Kogge, P. M. [1981]. *The Architecture of Pipelined Computers*, McGraw-Hill, New York.
- Kumar, A. [1997]. "The HP PA-8000 RISC CPU," *IEEE Micro* 17:2 (March/April).
- Kunkel, S. R., and J. E. Smith [1986]. "Optimal pipelining in supercomputers," *Proc. 13th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 2–5, 1986, Tokyo, 404–414.
- Lam, M. [1988]. "Software pipelining: An effective scheduling technique for VLIW processors," *SIGPLAN Conf. on Programming Language Design and Implementation*, June 22–24, 1988, Atlanta, Ga., 318–328.
- Lam, M. S., and R. P. Wilson [1992]. "Limits of control flow on parallelism," *Proc. 19th Annual Int'l. Symposium on Computer Architecture (ISCA)*, May 19–21, 1992, Gold Coast, Australia, 46–57.
- Laudon, J., A. Gupta, and M. Horowitz [1994]. "Interleaving: A multithreading technique targeting multiprocessors and workstations," *Proc. Sixth Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, October 4–7, San Jose, Calif., 308–318.
- Lauterbach, G., and T. Horel [1999]. "UltraSPARC-III: Designing third generation 64-bit performance," *IEEE Micro* 19:3 (May/June).
- Lipasti, M. H., and J. P. Shen [1996]. "Exceeding the dataflow limit via value prediction," *Proc. 29th Int'l. Symposium on Microarchitecture*, December 2–4, 1996, Paris, France.
- Lipasti, M. H., C. B. Wilkerson, and J. P. Shen [1996]. "Value locality and load value prediction," *Proc. Seventh Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, October 1–5, 1996, Cambridge, Mass., 138–147.
- Lo, J., L. Barroso, S. Eggers, K. Gharachorloo, H. Levy, and S. Parekh [1998]. "An analysis of database workload performance on simultaneous multithreaded processors," *Proc. 25th Annual Int'l. Symposium on Computer Architecture (ISCA)*, July 3–14, 1998, Barcelona, Spain, 39–50.
- Lo, J., S. Eggers, J. Emer, H. Levy, R. Stamm, and D. Tullsen [1997]. "Converting thread-level parallelism into instruction-level parallelism via simultaneous multithreading," *ACM Trans. on Computer Systems* 15:2 (August), 322–354.
- Mahlke, S. A., W. Y. Chen, W.-M. Hwu, B. R. Rau, and M. S. Schlansker [1992]. "Sentinel scheduling for VLIW and superscalar processors," *Proc. Fifth Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, October 12–15, 1992, Boston, 238–247.
- Mahlke, S. A., R. E. Hank, J. E. McCormick, D. I. August, and W. W. Hwu [1995]. "A comparison of full and partial predicated execution support for ILP processors," *Proc. 22nd Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 22–24, 1995, Santa Margherita, Italy, 138–149.
- Mathis, H. M., A. E. Mercias, J. D. McCalpin, R. J. Eickemeyer, and S. R. Kunkel [2005]. "Characterization of the multithreading (SMT) efficiency in Power5," *IBM J. of Research and Development*, 49:4/5 (July/September), 555–564.

- McCormick, J., and A. Knies [2002]. "A brief analysis of the SPEC CPU2000 benchmarks on the Intel Itanium 2 processor," paper presented at Hot Chips 14, August 18–20, 2002, Stanford University, Palo Alto, Calif.
- McFarling, S. [1993]. *Combining Branch Predictors*, WRL Technical Note TN-36, Digital Western Research Laboratory, Palo Alto, Calif.
- McFarling, S., and J. Hennessy [1986]. "Reducing the cost of branches," *Proc. 13th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 2–5, 1986, Tokyo, 396–403.
- McNairy, C., and D. Soltis [2003]. "Itanium 2 processor microarchitecture," *IEEE Micro* 23:2 (March–April), 44–55.
- Moshovos, A., and G. S. Sohi [1997]. "Streamlining inter-operation memory communication via data dependence prediction," *Proc. 30th Annual Int'l. Symposium on Microarchitecture*, December 1–3, Research Triangle Park, N.C., 235–245.
- Moshovos, A., S. Breach, T. N. Vijaykumar, and G. S. Sohi [1997]. "Dynamic speculation and synchronization of data dependences," *Proc. 24th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 2–4, 1997, Denver, Colo.
- Nicolau, A., and J. A. Fisher [1984]. "Measuring the parallelism available for very long instruction word architectures," *IEEE Trans. on Computers* C-33:11 (November), 968–976.
- Pan, S.-T., K. So, and J. T. Rameh [1992]. "Improving the accuracy of dynamic branch prediction using branch correlation," *Proc. Fifth Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, October 12–15, 1992, Boston, 76–84.
- Postiff, M.A., D. A. Greene, G. S. Tyson, and T. N. Mudge [1999]. "The limits of instruction level parallelism in SPEC95 applications," *Computer Architecture News* 27:1 (March), 31–40.
- Ramamoorthy, C. V., and H. F. Li [1977]. "Pipeline architecture," *ACM Computing Surveys* 9:1 (March), 61–102.
- Rau, B. R. [1994]. "Iterative modulo scheduling: An algorithm for software pipelining loops," *Proc. 27th Annual Int'l. Symposium on Microarchitecture*, November 30–December 2, 1994, San Jose, Calif., 63–74.
- Rau, B. R., C. D. Glaeser, and R. L. Picard [1982]. "Efficient code generation for horizontal architectures: Compiler techniques and architectural support," *Proc. Ninth Annual Int'l. Symposium on Computer Architecture (ISCA)*, April 26–29, 1982, Austin, Tex., 131–139.
- Rau, B. R., D. W. L. Yen, W. Yen, and R. A. Towle [1989]. "The Cydra 5 departmental supercomputer: Design philosophies, decisions, and trade-offs," *IEEE Computers* 22:1 (January), 12–34.
- Riseman, E. M., and C. C. Foster [1972]. "Percolation of code to enhance paralleled dispatching and execution," *IEEE Trans. on Computers* C-21:12 (December), 1411–1415.
- Rymarczyk, J. [1982]. "Coding guidelines for pipelined processors," *Proc. Symposium Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, March 1–3, 1982, Palo Alto, Calif., 12–19.

- Sharangpani, H., and K. Arora [2000]. "Itanium Processor Microarchitecture," *IEEE Micro*, 20:5 (September–October), 24–43.
- Sinharoy, B., R. N. Koala, J. M. Tandler, R. J. Eickemeyer, and J. B. Joyner [2005]. "POWER5 system microarchitecture," *IBM J. of Research and Development*, 49:4–5, 505–521.
- Sites, R. [1979]. *Instruction Ordering for the CRAY-1 Computer*, Tech. Rep. 78-CS-023, Dept. of Computer Science, University of California, San Diego.
- Skadron, K., P. S. Ahuja, M. Martonosi, and D. W. Clark [1999]. "Branch prediction, instruction-window size, and cache size: Performance tradeoffs and simulation techniques," *IEEE Trans. on Computers*, 48:11 (November).
- Smith, A., and J. Lee [1984]. "Branch prediction strategies and branch-target buffer design," *Computer* 17:1 (January), 6–22.
- Smith, B. J. [1978]. "A pipelined, shared resource MIMD computer," *Proc. Int'l. Conf. on Parallel Processing (ICPP)*, August, Bellaire, Mich., 6–8.
- Smith, J. E. [1981]. "A study of branch prediction strategies," *Proc. Eighth Annual Int'l. Symposium on Computer Architecture (ISCA)*, May 12–14, 1981, Minneapolis, Minn., 135–148.
- Smith, J. E. [1984]. "Decoupled access/execute computer architectures," *ACM Trans. on Computer Systems* 2:4 (November), 289–308.
- Smith, J. E. [1989]. "Dynamic instruction scheduling and the Astronautics ZS-1," *Computer* 22:7 (July), 21–35.
- Smith, J. E., and A. R. Pleszkun [1988]. "Implementing precise interrupts in pipelined processors," *IEEE Trans. on Computers* 37:5 (May), 562–573. (This paper is based on an earlier paper that appeared in *Proc. 12th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 17–19, 1985, Boston, Mass.)
- Smith, J. E., G. E. Dermer, B. D. Vanderwarn, S. D. Klinger, C. M. Rozewski, D. L. Fowler, K. R. Scidmore, and J. P. Laudon [1987]. "The ZS-1 central processor," *Proc. Second Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, October 5–8, 1987, Palo Alto, Calif., 199–204.
- Smith, M. D., M. Horowitz, and M. S. Lam [1992]. "Efficient superscalar performance through boosting," *Proc. Fifth Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, October 12–15, 1992, Boston, 248–259.
- Smith, M. D., M. Johnson, and M. A. Horowitz [1989]. "Limits on multiple instruction issue," *Proc. Third Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, April 3–6, 1989, Boston, 290–302.
- Sodani, A., and G. Sohi [1997]. "Dynamic instruction reuse," *Proc. 24th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 2–4, 1997, Denver, Colo.
- Sohi, G. S. [1990]. "Instruction issue logic for high-performance, interruptible, multiple functional unit, pipelined computers," *IEEE Trans. on Computers* 39:3 (March), 349–359.

- Sohi, G. S., and S. Vajapeyam [1989]. "Tradeoffs in instruction format design for horizontal architectures," *Proc. Third Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, April 3–6, 1989, Boston, 15–25.
- Sussenguth, E. [1999]. "IBM's ACS-1 Machine," *IEEE Computer* 22:11 (November).
- Tendler, J. M., J. S. Dodson, J. S. Fields, Jr., H. Le, and B. Sinharoy [2002]. "Power4 system microarchitecture," *IBM J. of Research and Development*, 46:1, 5–26.
- Thorlin, J. F. [1967]. "Code generation for PIE (parallel instruction execution) computers," *Proc. Spring Joint Computer Conf.*, April 18–20, 1967, Atlantic City, N.J., 27.
- Thornton, J. E. [1964]. "Parallel operation in the Control Data 6600," *Proc. AFIPS Fall Joint Computer Conf., Part II*, October 27–29, 1964, San Francisco, 26, 33–40.
- Thornton, J. E. [1970]. *Design of a Computer, the Control Data 6600*, Scott, Foresman, Glenview, Ill.
- Tjaden, G. S., and M. J. Flynn [1970]. "Detection and parallel execution of independent instructions," *IEEE Trans. on Computers* C-19:10 (October), 889–895.
- Tomasulo, R. M. [1967]. "An efficient algorithm for exploiting multiple arithmetic units," *IBM J. Research and Development* 11:1 (January), 25–33.
- Tuck, N., and D. Tullsen [2003]. "Initial observations of the simultaneous multithreading Pentium 4 processor," *Proc. 12th Int. Conf. on Parallel Architectures and Compilation Techniques (PACT'03)*, September 27–October 1, New Orleans, La., 26–34.
- Tullsen, D. M., S. J. Eggers, and H. M. Levy [1995]. "Simultaneous multithreading: Maximizing on-chip parallelism," *Proc. 22nd Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 22–24, 1995, Santa Margherita, Italy, 392–403.
- Tullsen, D. M., S. J. Eggers, J. S. Emer, H. M. Levy, J. L. Lo, and R. L. Stamm [1996]. "Exploiting choice: Instruction fetch and issue on an implementable simultaneous multithreading processor," *Proc. 23rd Annual Int'l. Symposium on Computer Architecture (ISCA)*, May 22–24, 1996, Philadelphia, Penn., 191–202.
- Wall, D. W. [1991]. "Limits of instruction-level parallelism," *Proc. Fourth Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, April 8–11, 1991, Palo Alto, Calif., 248–259.
- Wall, D. W. [1993]. *Limits of Instruction-Level Parallelism*, Research Rep. 93/6, Western Research Laboratory, Digital Equipment Corp., Palo Alto, Calif.
- Weiss, S., and J. E. Smith [1984]. "Instruction issue logic for pipelined supercomputers," *Proc. 11th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 5–7, 1984, Ann Arbor, Mich., 110–118.
- Weiss, S., and J. E. Smith [1987]. "A study of scalar compilation techniques for pipelined supercomputers," *Proc. Second Int'l. Conf. on Architectural Support*

- for Programming Languages and Operating Systems (ASPLOS)*, October 5–8, 1987, Palo Alto, Calif., 105–109.
- Wilson, R. P., and M. S. Lam [1995]. “Efficient context-sensitive pointer analysis for C programs,” *Proc. ACM SIGPLAN’95 Conf. on Programming Language Design and Implementation*, June 18–21, 1995, La Jolla, Calif., 1–12.
- Wolfe, A., and J. P. Shen [1991]. “A variable instruction stream extension to the VLIW architecture,” *Proc. Fourth Int’l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, April 8–11, 1991, Palo Alto, Calif., 2–14.
- Yamamoto, W., M. J. Serrano, A. R. Talcott, R. C. Wood, and M. Nemirosky [1994]. “Performance estimation of multistreamed, superscalar processors,” *Proc. 27th Annual Hawaii Int’l. Conf. on System Sciences*, January 4–7, 1994, Maui, 195–204.
- Yeager, K. [1996]. “The MIPS R10000 superscalar microprocessor,” *IEEE Micro* 16:2 (April), 28–40.
- Yeh, T., and Y. N. Patt [1992]. “Alternative implementations of two-level adaptive branch prediction,” *Proc. 19th Annual Int’l. Symposium on Computer Architecture (ISCA)*, May 19–21, 1992, Gold Coast, Australia, 124–134.
- Yeh, T., and Y. N. Patt [1993]. “A comparison of dynamic branch predictors that use two levels of branch history,” *Proc. 20th Annual Int’l. Symposium on Computer Architecture (ISCA)*, May 16–19, 1993, San Diego, Calif., 257–266.

---

## M.6

### The Development of SIMD Supercomputers, Vector Computers, Multimedia SIMD Instruction Extensions, and Graphical Processor Units (Chapter 4)

In this historical section, we start with perhaps the most infamous supercomputer, the Illiac IV, as a representative of the early SIMD (Single Instruction, Multiple Data) architectures and then move to perhaps the most famous supercomputer, the Cray-1, as a representative of vector architectures. The next step is Multimedia SIMD Extensions, which got its name in part due to an advertising campaign involving the “Bunny People,” a disco-dancing set of workers in cleansuits on a semiconductor fabrication line. We conclude with the history of GPUs, which is not quite as colorful.

#### SIMD Supercomputers

*The cost of a general multiprocessor is, however, very high and further design options were considered which would decrease the cost without seriously degrading the power or efficiency of the system. The options consist of recentralizing one of the three major components. ... Centralizing the [control unit] gives rise to the basic organization of [an] ... array processor such as the Illiac IV.*

**Bouknight et al. [1972]**



*... with Iliac IV, programming the machine was very difficult and the architecture probably was not very well suited to some of the applications we were trying to run. The key idea was that I did not think we had a very good match in Iliac IV between applications and architecture.*

**David Kuck**

*Software designer for the Illiac IV and  
early pioneer in parallel software*

**David Kuck**

*An oral history conducted in 1991 by Andrew Goldstein,  
IEEE History Center, New Brunswick, N.J.*

The SIMD model was one of the earliest models of parallel computing, dating back to the first large-scale multiprocessor, the Illiac IV. Rather than pipelining the data computation as in vector architectures, these machines had an array of functional units; hence, they might be considered array processors.

The earliest ideas on SIMD-style computers are from Unger [1958] and Slotnick, Borck, and McReynolds [1962]. Slotnick's Solomon design formed the basis of the Illiac IV, perhaps the most infamous of the supercomputer projects. Although successful in pushing several technologies that proved useful in later projects, it failed as a computer. Costs escalated from the \$8 million estimate in 1966 to \$31 million by 1972, despite construction of only a quarter of the planned multiprocessor. (In 2011 dollars, that was an increase from \$54M to \$152M.) Actual performance was at best 15 MFLOPS versus initial predictions of 1000 MFLOPS for the full system [Hord 1982]. Delivered to NASA Ames Research in 1972, the computer required three more years of engineering before it was usable. These events slowed investigation of SIMD, but Danny Hillis [1985] resuscitated this style in the Connection Machine, which had 65,536 1-bit processors.

The basic trade-off in SIMD multiprocessors is performance of a processor versus number of processors. SIMD supercomputers of the 1980s emphasized a large degree of parallelism over performance of the individual processors. The Connection Multiprocessor 2, for example, offered 65,536 single-bit-wide processors, while the Illiac IV planned for 64 64-bit processors. Massively parallel SIMD multiprocessors relied on interconnection or communication networks to exchange data between processing elements.

After being resurrected in the 1980s, first by Thinking Machines and then by MasPar, the SIMD model faded away as supercomputers for two main reasons. First, it is too inflexible. A number of important problems were not data parallel, and the architecture did not scale down in a competitive fashion; that is, small-scale SIMD multiprocessors often have worse cost-performance compared with that of the alternatives. Second, SIMD could not take advantage of the tremendous performance and cost advantages of SISD (Single Instruction, Single Data) microprocessor technology of the 1980s, which was doubling in performance every 18 months. Instead of leveraging this low-cost technology, designers of SIMD multiprocessors had to build custom processors for their multiprocessors.

## Vector Computers

*I'm certainly not inventing vector processors. There are three kinds that I know of existing today. They are represented by the Illiac-IV, the (CDC) Star processor, and the TI (ASC) processor. Those three were all pioneering processors. ... One of the problems of being a pioneer is you always make mistakes and I never, never want to be a pioneer. It's always best to come second when you can look at the mistakes the pioneers made.*

**Seymour Cray**

*Public lecture at Lawrence Livermore Laboratories  
on the introduction of the Cray-1 (1976)*

The first vector processors were the Control Data Corporation (CDC) STAR-100 (see Hintz and Tate [1972]) and the Texas Instruments ASC (see Watson [1972]), both announced in 1972. Both were memory-memory vector processors. They had relatively slow scalar units—the STAR used the same units for scalars and vectors—making the scalar pipeline extremely deep. Both processors had high start-up overhead and worked on vectors of several hundred to several thousand elements. The crossover between scalar and vector could be over 50 elements. It appears that not enough attention was paid to the role of Amdahl's law on these two processors.

Seymour Cray, who worked on the 6600 and the 7600 at CDC, founded Cray Research and introduced the Cray-1 in 1976 (see Russell [1978]). The Cray-1 used a vector-register architecture to lower start-up overhead significantly and to reduce memory bandwidth requirements. He also had efficient support for non-unit stride and invented chaining. Most importantly, the Cray-1 was the fastest scalar processor in the world at that time. This matching of good scalar and vector performance was probably the most significant factor in making the Cray-1 a success. Some customers bought the processor primarily for its outstanding scalar performance. Many subsequent vector processors are based on the architecture of this first commercially successful vector processor. Baskett and Keller [1977] provided a good evaluation of the Cray-1.

In 1981, CDC started shipping the CYBER 205 (see Lincoln [1982]). The 205 had the same basic architecture as the STAR but offered improved performance all around as well as expandability of the vector unit with up to four lanes, each with multiple functional units and a wide load-store pipe that provided multiple words per clock. The peak performance of the CYBER 205 greatly exceeded the performance of the Cray-1; however, on real programs, the performance difference was much smaller.

In 1983, Cray Research shipped the first Cray X-MP (see Chen [1983]). With an improved clock rate (9.5 ns versus 12.5 ns on the Cray-1), better chaining support (allowing vector operations with RAW dependencies to operate in parallel), and multiple memory pipelines, this processor maintained the Cray Research lead in supercomputers. The Cray-2, a completely new design configurable with up to four processors, was introduced later. A major feature of the Cray-2 was the use

of DRAM, which made it possible to have very large memories at the time. The first Cray-2, with its 256M word (64-bit words) memory, contained more memory than the total of all the Cray machines shipped to that point! The Cray-2 had a much faster clock than the X-MP, but also much deeper pipelines; however, it lacked chaining, had enormous memory latency, and had only one memory pipe per processor. In general, the Cray-2 was only faster than the Cray X-MP on problems that required its very large main memory.

That same year, processor vendors from Japan entered the supercomputer marketplace. First were the Fujitsu VP100 and VP200 (see Miura and Uchida [1983]), and later came the Hitachi S810 and the NEC SX/2 (see Watanabe [1987]). These processors proved to be close to the Cray X-MP in performance. In general, these three processors had much higher peak performance than the Cray X-MP. However, because of large start-up overhead, their typical performance was often lower than that of the Cray X-MP. The Cray X-MP favored a multiple-processor approach, first offering a two-processor version and later a four-processor version. In contrast, the three Japanese processors had expandable vector capabilities.

In 1988, Cray Research introduced the Cray Y-MP—a bigger and faster version of the X-MP. The Y-MP allowed up to eight processors and lowered the cycle time to 6 ns. With a full complement of eight processors, the Y-MP was generally the fastest supercomputer, though the single-processor Japanese supercomputers could be faster than a one-processor Y-MP. In late 1989, Cray Research was split into two companies, both aimed at building high-end processors available in the early 1990s. Seymour Cray headed the spin-off, Cray Computer Corporation, until its demise in 1995. Their initial processor, the Cray-3, was to be implemented in gallium arsenide, but they were unable to develop a reliable and cost-effective implementation technology. Shortly before his tragic death in a car accident in 1996, Seymour Cray started yet another company to develop high-performance systems but this time using commodity components.

Cray Research focused on the C90, a new high-end processor with up to 16 processors and a clock rate of 240 MHz. This processor was delivered in 1991. In 1993, Cray Research introduced their first highly parallel processor, the T3D, employing up to 2048 Digital Alpha21064 microprocessors. In 1995, they announced the availability of both a new low-end vector machine, the J90, and a high-end machine, the T90. The T90 was much like the C90, but with a clock that was twice as fast (460 MHz), using three-dimensional packaging and optical clock distribution.

In 1995, Cray Research was acquired by Silicon Graphics. In 1998, it released the SV1 system, which grafted considerably faster CMOS processors onto the J90 memory system. It also added a data cache for vectors to each CPU to help meet the increased memory bandwidth demands. Silicon Graphics sold Cray Research to Tera Computer in 2000, and the joint company was renamed Cray Inc.

The Japanese supercomputer makers continued to evolve their designs. In 2001, the NEC SX/5 was generally held to be the fastest available vector supercomputer, with 16 lanes clocking at 312 MHz and with up to 16 processors sharing the same memory. The NEC SX/6, released in 2001, was the first commercial single-chip vector microprocessor, integrating an out-of-order quad-issue

superscalar processor, scalar instruction and data caches, and an eight-lane vector unit on a single die [Kitagawa et al. 2003]. The Earth Simulator is constructed from 640 nodes connected with a full crossbar, where each node comprises eight SX-6 vector microprocessors sharing a local memory. The SX-8, released in 2004, reduces the number of lanes to four but increases the vector clock rate to 2 GHz. The scalar unit runs at a slower 1 GHz clock rate, a common pattern in vector machines where the lack of hazards simplifies the use of deeper pipelines in the vector unit.

In 2002, Cray Inc. released the X1 based on a completely new vector ISA. The X1 SSP processor chip integrates an out-of-order superscalar with scalar caches running at 400 MHz and a two-lane vector unit running at 800 MHz. When four SSP chips are ganged together to form an MSP, the resulting peak vector performance of 12.8 GFLOPS is competitive with the contemporary NEC SX machines. The X1E enhancement, delivered in 2004, raises the clock rates to 565 and 1130 MHz, respectively. Many of the ideas were borrowed from the Cray T3E design, which is a MIMD (Multiple Instruction, Multiple Data) computer that uses off-the-shelf microprocessors. X1 has a new instruction set with a larger number of registers and with memory distributed locally with the processor in shared address space. The out-of-order scalar unit and vector units are decoupled, so that the scalar unit can get ahead of the vector unit. Vectors become shorter when the data are blocked to utilize the MSP caches, which is not a good match to an eight-lane vector unit. To handle these shorter vectors, each processor with just two vector lanes can work on a different loop.

The Cray X2 was announced in 2007, and it may prove to be the last Cray vector architecture to be built, as it's difficult to justify the investment in new silicon given the size of the market. The processor has a 1.3 GHz clock rate and 8 vector lanes for a processor peak performance of 42 GFLOP/sec for single precision. It includes both L1 and L2 caches. Each node is a 4-way SMP with up to 128 GBytes of DRAM, and the maximum size is 8K nodes.

The NEC SX-9 has up to 16 processors per node, with each processor having 8 lanes and running at 3.2 GHz. It was announced in 2008. The peak double precision vector performance is 102 GFLOP/sec. The 16 processor SMP can have 1024 GBytes of DRAM. The maximum size is 512 nodes.

The basis for modern vectorizing compiler technology and the notion of data dependence was developed by Kuck and his colleagues [1974] at the University of Illinois. Padua and Wolfe [1986] gave a good overview of vectorizing compiler technology.

## Multimedia SIMD Instruction Extensions

*What could a computer hardware company ... possibly have in common with disco dancing. A lot, if one goes by an advertisement campaign released by the world's largest microprocessor company ... Intel, in 1997.*

**IBS Center for Management Research**  
*"Dancing Its Way Towards Leadership," 2002*

Going through the history books, the 1957 TX-2 had partitioned ALUs to support media of the time, but these ideas faded away to be rediscovered 30 years later in the personal computer era. Since every desktop microprocessor by definition has its own graphical displays, as transistor budgets increased it was inevitable that support would be added for graphics operations. Many graphics systems use 8 bits to represent each of the 3 primary colors plus 8 bits for a transparency of a pixel. The addition of speakers and microphones for teleconferencing and video games suggested support of sound as well. Audio samples need more than 8 bits of precision, but 16 bits are sufficient.

Every microprocessor has special support so that bytes and half words take up less space when stored in memory, but due to the infrequency of arithmetic operations on these data sizes in typical integer programs, there is little support beyond data transfers. The Intel i860 was justified as a graphical accelerator within the company. Its architects recognized that many graphics and audio applications would perform the same operation on vectors of these data [Atkins 1991; Kohn 1989]. Although a vector unit was beyond the transistor budget of the i860 in 1989, by partitioning the carry chains within a 64-bit ALU, it could perform simultaneous operations on short vectors of eight 8-bit operands, four 16-bit operands, or two 32-bit operands. The cost of such partitioned ALUs was small. Applications that lend themselves to such support include MPEG (video), video games (3D graphics), digital photography, and teleconferencing (audio and image processing).

Like a virus, over time such multimedia support has spread to nearly every desktop microprocessor. HP was the first successful desktop RISC to include such support, but soon every other manufacturer had their own take on the idea in the 1990s.

These extensions were originally called *subword parallelism* or *vector*. Since Intel marketing used SIMD to describe the MMX extension of the 80x86 announced in 1996, that became the popular name, due in part to a successful television advertising campaign involving disco dancers wearing clothing modeled after the cleansuits worn in semiconductor fabrication lines.

## Graphical Processor Units

*It's been almost three years since GPU computing broke into the mainstream of HPC with the introduction of NVIDIA's CUDA API in September 2007. Adoption of the technology since then has proceeded at a surprisingly strong and steady pace. Many organizations that began with small pilot projects a year or two ago have moved on to enterprise deployment, and GPU accelerated machines are now represented on the TOP500 list starting at position two. The relatively rapid adoption of CUDA by a community not known for the rapid adoption of much of anything is a noteworthy signal. Contrary to the accepted wisdom that GPU computing is more difficult, I believe its success thus far signals that it is no more complicated than good CPU programming. Further, it more clearly and succinctly expresses the parallelism of a large class of problems leading to code*

*that is easier to maintain, more scalable and better positioned to map to future many-core architectures.*

**Vincent Natol**

*"Kudos for CUDA," HPCwire (2010)*

3D graphics pipeline hardware evolved from the large expensive systems of the early 1980s to small workstations and then to PC accelerators in the mid- to late 1990s. During this period, three major transitions occurred:

- Performance-leading graphics subsystems declined in price from \$50,000 to \$200.
- Performance increased from 50 million pixels per second to 1 billion pixels per second and from 100,000 vertices per second to 10 million vertices per second.
- Native hardware capabilities evolved from wireframe (polygon outlines) to flat-shaded (constant color) filled polygons, to smooth-shaded (interpolated color) filled polygons, to full-scene anti-aliasing with texture mapping and rudimentary multitexturing.

## Scalable GPUs

Scalability has been an attractive feature of graphics systems from the beginning. Workstation graphics systems gave customers a choice in pixel horse-power by varying the number of pixel processor circuit boards installed. Prior to the mid-1990s PC graphics scaling was almost nonexistent. There was one option—the VGA controller. As 3D-capable accelerators appeared, the market had room for a range of offerings. 3dfx introduced multiboard scaling with the original SLI (Scan Line Interleave) on their Voodoo2, which held the performance crown for its time (1998). Also in 1998, NVIDIA introduced distinct products as variants on a single architecture with Riva TNT Ultra (high-performance) and Vanta (low-cost), first by speed binning and packaging, then with separate chip designs (GeForce 2 GTS and GeForce 2 MX). At present, for a given architecture generation, four or five separate GPU chip designs are needed to cover the range of desktop PC performance and price points. In addition, there are separate segments in notebook and workstation systems. After acquiring 3dfx, NVIDIA continued the multi-GPU SLI concept in 2004, starting with GeForce 6800—providing multi-GPU scalability transparently to the programmer and to the user. Functional behavior is identical across the scaling range; one application will run unchanged on any implementation of an architectural family.

## Graphics Pipelines

Early graphics hardware was configurable, but not programmable by the application developer. With each generation, incremental improvements were offered; however, developers were growing more sophisticated and asking for more new features than could be reasonably offered as built-in fixed functions. The NVIDIA



GeForce 3, described by Lindholm et al. [2001], took the first step toward true general shader programmability. It exposed to the application developer what had been the private internal instruction set of the floating-point vertex engine. This coincided with the release of Microsoft's DirectX 8 and OpenGL's vertex shader extensions. Later GPUs, at the time of DirectX 9, extended general programmability and floating-point capability to the pixel fragment stage and made texture available at the vertex stage. The ATI Radeon 9700, introduced in 2002, featured a programmable 24-bit floating-point pixel fragment processor programmed with DirectX 9 and OpenGL. The GeForce FX added 32-bit floating-point pixel processors. This was part of a general trend toward unifying the functionality of the different stages, at least as far as the application programmer was concerned. NVIDIA's GeForce 6800 and 7800 series were built with separate processor designs and separate hardware dedicated to the vertex and to the fragment processing. The Xbox 360 introduced an early unified processor GPU in 2005, allowing vertex and pixel shaders to execute on the same processor.

### **GPGPU: An Intermediate Step**

As DirectX 9-capable GPUs became available, some researchers took notice of the raw performance growth path of GPUs and began to explore the use of GPUs to solve complex parallel problems. DirectX 9 GPUs had been designed only to match the features required by the graphics API. To access the computational resources, a programmer had to cast their problem into native graphics operations. For example, to run many simultaneous instances of a pixel shader, a triangle had to be issued to the GPU (with clipping to a rectangle shape if that was what was desired). Shaders did not have the means to perform arbitrary scatter operations to memory. The only way to write a result to memory was to emit it as a pixel color value and configure the framebuffer operation stage to write (or blend, if desired) the result to a two-dimensional framebuffer. Furthermore, the only way to get a result from one pass of computation to the next was to write all parallel results to a pixel framebuffer, then use that framebuffer as a texture map as input to the pixel fragment shader of the next stage of the computation. Mapping general computations to a GPU in this era was quite awkward. Nevertheless, intrepid researchers demonstrated a handful of useful applications with painstaking efforts. This field was called "GPGPU" for general-purpose computing on GPUs.

### **GPU Computing**

While developing the Tesla architecture for the GeForce 8800, NVIDIA realized its potential usefulness would be much greater if programmers could think of the GPU as a processor. NVIDIA selected a programming approach in which programmers would explicitly declare the data-parallel aspects of their workload.

For the DirectX 10 generation, NVIDIA had already begun work on a high-efficiency floating-point and integer processor that could run a variety of

simultaneous workloads to support the logical graphics pipeline. This processor was designed to take advantage of the common case of groups of threads executing the same code path. NVIDIA added memory load and store instructions with integer byte addressing to support the requirements of compiled C programs. It introduced the thread block (cooperative thread array), grid of thread blocks, and barrier synchronization to dispatch and manage highly parallel computing work. Atomic memory operations were added. NVIDIA developed the CUDA C/C++ compiler, libraries, and runtime software to enable programmers to readily access the new data-parallel computation model and develop applications.

To create a vendor-neutral GPU programming language, a large number of companies are creating compilers for the OpenCL language, which has many of the features of CUDA but which runs on many more platforms. In 2011, the performance is much higher if you write CUDA code for GPUs than if you write OpenCL code.

AMD's acquisition of ATI, the second leading GPU vendor, suggests a spread of GPU computing. The AMD Fusion architecture, announced just as this edition was being finished, is an initial merger between traditional GPUs and traditional CPUs. NVIDIA also announced Project Denver, which combines an ARM scalar processor with NVIDIA GPUs in a single address space. When these systems are shipped, it will be interesting to learn just how tightly integrated they are and the impact of integration on performance and energy of both data parallel and graphics applications.

## References

### SIMD Supercomputers

- Bouknight, W. J., S. A. Deneberg, D. E. McIntyre, J. M. Randall, A. H. Sameh, and D. L. Slotnick [1972]. "The Illiac IV system," *Proc. IEEE* 60:4, 369–379. Also appears in D. P. Siewiorek, C. G. Bell, and A. Newell, *Computer Structures: Principles and Examples*, McGraw-Hill, New York, 1982, 306–316.
- Hillis, W. D. [1985]. *The Connection Multiprocessor*, MIT Press, Cambridge, Mass.
- Hord, R. M. [1982]. *The Illiac-IV, The First Supercomputer*, Computer Science Press, Rockville, Md.
- Slotnick, D. L., W. C. Borck, and R. C. McReynolds [1962]. "The Solomon computer," *Proc. AFIPS Fall Joint Computer Conf.*, December 4–6, 1962, Philadelphia, Penn., 97–107.
- Unger, S. H. [1958]. "A computer oriented towards spatial problems," *Proc. Institute of Radio Engineers* 46:10 (October), 1744–1750.

### Vector Architecture

- Asanovic, K. [1998]. "Vector Microprocessors," Ph.D. thesis, Computer Science Division, University of California, Berkeley.
- Baskett, F., and T. W. Keller [1977]. "An Evaluation of the Cray-1 Processor," in *High Speed Computer and Algorithm Organization*, D. J. Kuck, D. H. Lawrie, and A. H. Sameh, eds., Academic Press, San Diego, Calif., 71–84.

- Chen, S. [1983]. "Large-scale and high-speed multiprocessor system for scientific applications," *Proc. NATO Advanced Research Workshop on High Speed Computing*, June 20–22, Julich, West Germany. Also in K. Hwang, ed., "Superprocessors: Design and applications," *IEEE*, August, 59–73, 1984.
- Flynn, M. J. [1966]. "Very high-speed computing systems," *Proc. IEEE* 54:12 (December), 1901–1909.
- Gebis, J. and Patterson, D. [2007]. "Embracing and extending 20th-century instruction set architectures," *IEEE Computer*, 40:4 (April), 68–75.
- Hintz, R. G., and D. P. Tate [1972]. "Control data STAR-100 processor design," *Proc. IEEE COMPCON*, September 12–14, 1972, San Francisco, 1–4.
- Kitagawa, K., S. Tagaya, Y. Hagihara, and Y. Kanoh [2003]. "A hardware overview of SX-6 and SX-7 supercomputer," *NEC Research and Development Journal* 44:1 (January), 2–7.
- Kozyrakis, C., and D. Patterson [2002]. "Vector vs. superscalar and VLIW architectures for embedded multimedia benchmarks," *Proc. 35th Annual Intl. Symposium on Microarchitecture (MICRO)*, November 18–22, 2002, Istanbul, Turkey.
- Kuck, D., P. P. Budnik, S.-C. Chen, D. H. Lawrie, R. A. Towle, R. E. Strebendt, E. W. Davis, Jr., J. Han, P. W. Kraska, and Y. Muraoka [1974]. "Measurements of parallelism in ordinary Fortran programs," *Computer* 7:1 (January), 37–46.
- Lincoln, N. R. [1982]. "Technology and design trade offs in the creation of a modern supercomputer," *IEEE Trans. on Computers* C-31:5 (May), 363–376.
- Miura, K., and K. Uchida [1983]. "FACOM vector processing system: VP100/200," *Proc. NATO Advanced Research Workshop on High Speed Computing*, June 20–22, Julich, West Germany. Also in K. Hwang, ed., "Superprocessors: Design and applications," *IEEE*, August, 59–73, 1984.
- Padua, D., and M. Wolfe [1986]. "Advanced compiler optimizations for supercomputers," *Communications of the ACM* 29:12 (December), 1184–1201.
- Russell, R. M. [1978]. "The Cray-1 processor system," *Communications of the ACM* 21:1 (January), 63–72.
- Vajapeyam, S. [1991]. "Instruction-Level Characterization of the Cray Y-MP Processor," Ph.D. thesis, Computer Sciences Department, University of Wisconsin–Madison.
- Watanabe, T. [1987]. "Architecture and performance of the NEC supercomputer SX system," *Parallel Computing* 5, 247–255.
- Watson, W. J. [1972]. "The TI ASC—a highly modular and flexible super processor architecture," *Proc. AFIPS Fall Joint Computer Conf.*, December 5–7, 1972, Anaheim, Calif., 221–228.

## Multimedia SIMD

- Atkins, M. [1991]. "Performance and the i860 Microprocessor," *IEEE Micro*, 11:5 (September), 24–27, 72–78.
- Kohn, L., and N. Margulis [1989]. "Introducing the Intel i860 64-Bit Microprocessor," *IEEE Micro*, 9:4 (July), 15–30.

## GPU

- Akeley, K., and T. Jermoluk [1988]. “High-performance polygon rendering,” *Proc. SIGGRAPH 88*, August 1–5, 1988, Atlanta, Ga., 239–46.
- Hillis, W. D., and G. L. Steele [1986]. “Data parallel algorithms,” *Communications of the ACM* 29:12 (December), 1170–1183 (<http://doi.acm.org/10.1145/7902.7903>).
- IEEE 754-2008 Working Group. [2006]. *DRAFT Standard for Floating-Point Arithmetic*, 754-2008 (<https://doi.org/10.1109/IEEESTD.2008.4610935>).
- Lee, W. V., et al. [2010]. “Debunking the 100X GPU vs. CPU myth: an evaluation of throughput computing on CPU and GPU,” *Proc. ISCA '10*, June 19–23, 2010, Saint-Malo, France.
- Lindholm, E., M. J. Kligard, and H. Moreton [2001]. A user-programmable vertex engine. In *SIGGRAPH '01: Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, 149–158.
- Moore, G. E. [1965]. “Cramming more components onto integrated circuits,” *Electronics* 38:8 (April 19), 114–117.
- Williams, S., A. Waterman, and D. Patterson [2009]. “Roofline: An insightful visual performance model for multicore architectures,” *Communications of the ACM*, 52:4 (April), 65–76.

## M.7

### The History of Multiprocessors and Parallel Processing (Chapter 5 and Appendices F, G, and I)

There is a tremendous amount of history in multiprocessors; in this section, we divide our discussion by both time period and architecture. We start with the SIMD approach and the Illiac IV. We then turn to a short discussion of some other early experimental multiprocessors and progress to a discussion of some of the great debates in parallel processing. Next we discuss the historical roots of the present multiprocessors and conclude by discussing recent advances.

#### SIMD Computers: Attractive Idea, Many Attempts, No Lasting Successes

*The cost of a general multiprocessor is, however, very high and further design options were considered which would decrease the cost without seriously degrading the power or efficiency of the system. The options consist of recentralizing one of the three major components. ... Centralizing the [control unit] gives rise to the basic organization of [an] ... array processor such as the Illiac IV.*

**Bouknight et al. [1972]**

The SIMD model was one of the earliest models of parallel computing, dating back to the first large-scale multiprocessor, the Illiac IV. The key idea in that multiprocessor, as in more recent SIMD multiprocessors, is to have a single instruction that operates on many data items at once, using many functional units.

The earliest ideas on SIMD-style computers are from Unger [1958] and Slotnick, Borck, and McReynolds [1962]. Slotnick's Solomon design formed the basis of the Illiac IV, perhaps the most infamous of the supercomputer projects. Although successful in pushing several technologies that proved useful in later projects, it failed as a computer. Costs escalated from the \$8 million estimate in 1966 to \$31 million by 1972, despite construction of only a quarter of the planned multiprocessor. Actual performance was at best 15 MFLOPS versus initial predictions of 1000 MFLOPS for the full system [Hord 1982]. Delivered to NASA Ames Research in 1972, the computer took three more years of engineering before it was usable. These events slowed investigation of SIMD, but Danny Hillis [1985] resuscitated this style in the Connection Machine, which had 65,636 1-bit processors.

Real SIMD computers need to have a mixture of SISD and SIMD instructions. There is an SISD host computer to perform operations such as branches and address calculations that do not need parallel operation. The SIMD instructions are broadcast to all the execution units, each of which has its own set of registers. For flexibility, individual execution units can be disabled during an SIMD instruction. In addition, massively parallel SIMD multiprocessors rely on interconnection or communication networks to exchange data between processing elements.

SIMD works best in dealing with arrays in for loops; hence, to have the opportunity for massive parallelism in SIMD there must be massive amounts of data, or *data parallelism*. SIMD is at its weakest in case statements, where each execution unit must perform a different operation on its data, depending on what data it has. The execution units with the wrong data are disabled so that the proper units can continue. Such situations essentially run at  $1/n$ th performance, where  $n$  is the number of cases.

The basic trade-off in SIMD multiprocessors is performance of a processor versus number of processors. Recent multiprocessors emphasize a large degree of parallelism over performance of the individual processors. The Connection Multiprocessor 2, for example, offered 65,536 single-bit-wide processors, while the Illiac IV had 64 64-bit processors.

After being resurrected in the 1980s, first by Thinking Machines and then by MasPar, the SIMD model has once again been put to bed as a general-purpose multiprocessor architecture, for two main reasons. First, it is too inflexible. A number of important problems cannot use such a style of multiprocessor, and the architecture does not scale down in a competitive fashion; that is, small-scale SIMD multiprocessors often have worse cost-performance compared with that of the alternatives. Second, SIMD cannot take advantage of the tremendous performance and cost advantages of microprocessor technology. Instead of leveraging this low-cost technology, designers of SIMD multiprocessors must build custom processors for their multiprocessors.

Although SIMD computers have departed from the scene as general-purpose alternatives, this style of architecture will continue to have a role in special-purpose

designs. Many special-purpose tasks are highly data parallel and require a limited set of functional units. Thus, designers can build in support for certain operations, as well as hardwired interconnection paths among functional units. Such organizations are often called *array processors*, and they are useful for such tasks as image and signal processing.

### Other Early Experiments

It is difficult to distinguish the first MIMD multiprocessor. Surprisingly, the first computer from the Eckert-Mauchly Corporation, for example, had duplicate units to improve availability. Holland [1959] gave early arguments for multiple processors. Two of the best-documented multiprocessor projects were undertaken in the 1970s at Carnegie Mellon University. The first of these was C.mmp [Wulf and Bell 1972; Wulf and Harbison 1978], which consisted of 16 PDP-11s connected by a crossbar switch to 16 memory units. It was among the first multiprocessors with more than a few processors, and it had a shared-memory programming model. Much of the focus of the research in the C.mmp project was on software, especially in the OS area. A later multiprocessor, Cm\* [Swan et al. 1977], was a cluster-based multiprocessor with a distributed memory and a nonuniform access time. The absence of caches and a long remote access latency made data placement critical. This multiprocessor and a number of application experiments are well described by Gehring, Siewiorek, and Segall [1987]. Many of the ideas in these multiprocessors would be reused in the 1980s when the microprocessor made it much cheaper to build multiprocessors.

### Great Debates in Parallel Processing

*The turning away from the conventional organization came in the middle 1960s, when the law of diminishing returns began to take effect in the effort to increase the operational speed of a computer. ... Electronic circuits are ultimately limited in their speed of operation by the speed of light ... and many of the circuits were already operating in the nanosecond range.*

**Bouknight et al. [1972]**

*... sequential computers are approaching a fundamental physical limit on their potential computational power. Such a limit is the speed of light ...*

**Angel L. DeCegama**

*The Technology of Parallel Processing, Vol. I (1989)*

*... today's multiprocessors ... are nearing an impasse as technologies approach the speed of light. Even if the components of a sequential processor could be made to work this fast, the best that could be expected is no more than a few million instructions per second.*

**David Mitchell**

*The Transputer: The Time Is Now (1989)*



The quotes above give the classic arguments for abandoning the current form of computing, and Amdahl [1967] gave the classic reply in support of continued focus on the IBM 360 architecture. Arguments for the advantages of parallel execution can be traced back to the 19th century [Menabrea 1842]! Yet, the effectiveness of the multiprocessor for reducing latency of individual important programs is still being explored. Aside from these debates about the advantages and limitations of parallelism, several hot debates have focused on how to build multiprocessors.

It's hard to predict the future, yet in 1989 Gordon Bell made two predictions for 1995. We included these predictions in the first edition of the book, when the outcome was completely unclear. We discuss them in this section, together with an assessment of the accuracy of the prediction.

The first was that a computer capable of sustaining a teraFLOPS—one million MFLOPS—would be constructed by 1995, using either a multicomputer with 4K to 32K nodes or a Connection Multiprocessor with several million processing elements [Bell 1989]. To put this prediction in perspective, each year the Gordon Bell Prize acknowledges advances in parallelism, including the fastest real program (highest MFLOPS). In 1989, the winner used an eight-processor Cray Y-MP to run at 1680 MFLOPS. On the basis of these numbers, multiprocessors and programs would have to have improved by a factor of 3.6 each year for the fastest program to achieve 1 TFLOPS in 1995. In 1999, the first Gordon Bell prize winner crossed the 1 TFLOPS bar. Using a 5832-processor IBM RS/6000 SST system designed specially for Livermore Laboratories, they achieved 1.18 TFLOPS on a shock-wave simulation. This ratio represents a year-to-year improvement of 1.93, which is still quite impressive.

What has become recognized since the 1990s is that, although we may have the technology to build a TFLOPS multiprocessor, it is not clear that the machine is cost effective, except perhaps for a few very specialized and critically important applications related to national security. We estimated in 1990 that to achieve 1 TFLOPS would require a machine with about 5000 processors and would cost about \$100 million. The 5832-processor IBM system at Livermore cost \$110 million. As might be expected, improvements in the performance of individual microprocessors both in cost and performance directly affect the cost and performance of large-scale multiprocessors, but a 5000-processor system will cost more than 5000 times the price of a desktop system using the same processor. Since that time, much faster multiprocessors have been built, but the major improvements have increasingly come from the processors in the past five years, rather than fundamental breakthroughs in parallel architecture.

The second Bell prediction concerned the number of data streams in supercomputers shipped in 1995. Danny Hillis believed that, although supercomputers with a small number of data streams may be the best sellers, the biggest multiprocessors would be multiprocessors with many data streams, and these would perform the bulk of the computations. Bell bet Hillis that in the last quarter of calendar year 1995 more sustained MFLOPS would be shipped in multiprocessors using few data streams ( $\leq 100$ ) rather than many data streams ( $\geq 1000$ ). This bet concerned

only supercomputers, defined as multiprocessors costing more than \$1 million and used for scientific applications. Sustained MFLOPS was defined for this bet as the number of floating-point operations per *month*, so availability of multiprocessors affects their rating.

In 1989, when this bet was made, it was totally unclear who would win. In 1995, a survey of the current publicly known supercomputers showed only six multiprocessors in existence in the world with more than 1000 data streams, so Bell's prediction was a clear winner. In fact, in 1995, much smaller microprocessor-based multiprocessors ( $\leq 20$  processors) were becoming dominant. In 1995, a survey of the 500 highest-performance multiprocessors in use (based on Linpack ratings), called the TOP500, showed that the largest number of multiprocessors were bus-based shared-memory multiprocessors! By 2005, various clusters or multicomputers played a large role. For example, in the top 25 systems, 11 were custom clusters, such as the IBM Blue Gene system or the Cray XT3; 10 were clusters of shared-memory multiprocessors (both using distributed and centralized memory); and the remaining 4 were clusters built using PCs with an off-the-shelf interconnect.

## More Recent Advances and Developments

With the primary exception of the parallel vector multiprocessors (see Appendix G) and more recently of the IBM Blue Gene design, all other recent MIMD computers have been built from off-the-shelf microprocessors using a bus and logically central memory or an interconnection network and a distributed memory. A number of experimental multiprocessors built in the 1980s further refined and enhanced the concepts that form the basis for many of today's multiprocessors.

### *The Development of Bus-Based Coherent Multiprocessors*

Although very large mainframes were built with multiple processors in the 1960s and 1970s, multiprocessors did not become highly successful until the 1980s. Bell [1985] suggested that the key was that the smaller size of the microprocessor allowed the memory bus to replace the interconnection network hardware and that portable operating systems meant that multiprocessor projects no longer required the invention of a new operating system. In his paper, Bell defined the terms *multiprocessor* and *multicomputer* and set the stage for two different approaches to building larger scale multiprocessors.

The first bus-based multiprocessor with snooping caches was the Synapse N+1 described by Frank [1984]. Goodman [1983] wrote one of the first papers to describe snooping caches. The late 1980s saw the introduction of many commercial bus-based, snooping cache architectures, including the Silicon Graphics 4D/240 [Baskett, Jermoluk, and Solomon 1988], the Encore Multimax [Wilson 1987], and the Sequent Symmetry [Lovett and Thakkar 1988]. The mid-1980s

Name	Protocol type	Memory write policy	Unique feature	Multiprocessors using
Write Once	Write invalidate	Write-back after first write	First snooping protocol described in literature	
Synapse N+1	Write invalidate	Write-back	Explicit state where memory is the owner	Synapse multiprocessors; first cache-coherent multiprocessors available
Berkeley (MOESI)	Write invalidate	Write-back	Owned shared state	Berkeley SPUR multiprocessor; Sun Enterprise servers
Illinois (MESI)	Write invalidate	Write-back	Clean private state; can supply data from any cache with a clean copy	SGI Power and Challenge series
“Firefly”	Write broadcast	Write-back when private, write through when shared	Memory updated on broadcast	No current multiprocessors; SPARCCenter 2000 closest

**Figure M.2 Five snooping protocols summarized.** Archibald and Baer [1986] use these names to describe the five protocols, and Eggers [1989] summarizes the similarities and differences as shown in this figure. The Firefly protocol was named for the experimental DEC Firefly multiprocessor, in which it appeared. The alternative names for protocols are based on the states they support: M = Modified, E = Exclusive (private clean), S = Shared, I = Invalid, O = Owner (shared dirty).

saw an explosion in the development of alternative coherence protocols, and Archibald and Baer [1986] provided a good survey and analysis, as well as references to the original papers. Figure M.2 summarizes several snooping cache coherence protocols and shows some multiprocessors that have used or are using that protocol.

The early 1990s saw the beginning of an expansion of such systems with the use of very wide, high-speed buses (the SGI Challenge system used a 256-bit, packet-oriented bus supporting up to 8 processor boards and 32 processors) and later the use of multiple buses and crossbar interconnects—for example, in the Sun SPARCCenter and Enterprise systems (Charlesworth [1998] discussed the interconnect architecture of these multiprocessors). In 2001, the Sun Enterprise servers represented the primary example of large-scale (>16 processors), symmetric multiprocessors in active use. Today, most bus-based machines offer only four or so processors and switches, or alternative designs are used for eight or more.

## Toward Large-Scale Multiprocessors

In the effort to build large-scale multiprocessors, two different directions were explored: message-passing multicomputers and scalable shared-memory multiprocessors. Although there had been many attempts to build mesh and hypercube-connected multiprocessors, one of the first multiprocessors to successfully bring together all the pieces was the Cosmic Cube built at Caltech [Seitz 1985]. It introduced important advances in routing and interconnect technology and substantially

reduced the cost of the interconnect, which helped make the multicomputer viable. The Intel iPSC 860, a hypercube-connected collection of i860s, was based on these ideas. More recent multiprocessors, such as the Intel Paragon, have used networks with lower dimensionality and higher individual links. The Paragon also employed a separate i860 as a communications controller in each node, although a number of users have found it better to use both i860 processors for computation as well as communication. The Thinking Multiprocessors CM-5 made use of off-the-shelf microprocessors and a fat tree interconnect (see Appendix F). It provided user-level access to the communication channel, thus significantly improving communication latency. In 1995, these two multiprocessors represented the state of the art in message-passing multicomputers.

Early attempts at building a scalable shared-memory multiprocessor include the IBM RP3 [Pfister et al. 1985], the NYU Ultracomputer [Elder et al. 1985; Schwartz 1980], the University of Illinois Cedar project [Gajski et al. 1983], and the BBN Butterfly and Monarch [BBN Laboratories 1986; Rettberg et al. 1990]. These multiprocessors all provided variations on a nonuniform distributed-memory model and, hence, are distributed shared-memory (DSM) multiprocessors, but they did not support cache coherence, which substantially complicated programming. The RP3 and Ultracomputer projects both explored new ideas in synchronization (fetch-and-operate) as well as the idea of combining references in the network. In all four multiprocessors, the interconnect networks turned out to be more costly than the processing nodes, raising problems for smaller versions of the multiprocessor. The Cray T3D/E (see Arpaci et al. [1995] for an evaluation of the T3D and Scott [1996] for a description of the T3E enhancements) builds on these ideas, using a noncoherent shared address space but building on the advances in interconnect technology developed in the multicomputer domain (see Scott and Thorson [1996]).

Extending the shared-memory model with scalable cache coherence was done by combining a number of ideas. Directory-based techniques for cache coherence were actually known before snooping cache techniques. In fact, the first cache coherence protocols actually used directories, as described by Tang [1976] and implemented in the IBM 3081. Censier and Feautrier [1978] described a directory coherence scheme with tags in memory. The idea of distributing directories with the memories to obtain a scalable implementation of cache coherence was first described by Agarwal et al. [1988] and served as the basis for the Stanford DASH multiprocessor (see Lenoski et al. [1990, 1992]), which was the first operational cache-coherent DSM multiprocessor. DASH was a “plump” node cc-NUMA machine that used four-processor SMPs as its nodes, interconnecting them in a style similar to that of Wildfire but using a more scalable two-dimensional grid rather than a crossbar for the interconnect.

The Kendall Square Research KSR-1 [Burkhardt et al. 1992] was the first commercial implementation of scalable coherent shared memory. It extended the basic DSM approach to implement a concept called *cache-only memory architecture* (COMA), which makes the main memory a cache. In the KSR-1, memory blocks could be replicated in the main memories of each node with hardware support to

handle the additional coherence requirements for these replicated blocks. (The KSR-1 was not strictly a pure COMA because it did not migrate the home location of a data item but always kept a copy at home. Essentially, it implemented only replication.) Many other research proposals [Falsafi and Wood 1997; Hagersten, Landin, and Haridi 1992; Saulsbury et al. 1995; Stenström, Joe, and Gupta 1992] for COMA-style architectures and similar approaches that reduce the burden of nonuniform memory access through migration [Chandra et al. 1994; Soundararajan et al. 1998] were developed, but there have been no further commercial implementations.

The Convex Exemplar implemented scalable coherent shared memory using a two-level architecture: At the lowest level, eight-processor modules are built using a crossbar. A ring can then connect up to 32 of these modules, for a total of 256 processors (see Thekkath et al. [1997] for an evaluation). Laudon and Lenoski [1997] described the SGI Origin, which was first delivered in 1996 and is closely based on the original Stanford DASH machine, although including a number of innovations for scalability and ease of programming. Origin uses a bit vector for the directory structure, which is either 16 or 32 bits long. Each bit represents a node, which consists of two processors; a coarse bit vector representation allows each bit to represent up to 8 nodes for a total of 1024 processors. As Galles [1996] described, a high-performance fat hypercube is used for the global interconnect. Hristea, Lenoski, and Keen [1997] have provided a thorough evaluation of the performance of the Origin memory system.

Several research prototypes were undertaken to explore scalable coherence with and without multithreading. These include the MIT Alewife machine [Agarwal et al. 1995] and the Stanford FLASH multiprocessor [Gibson et al. 2000; Kuskina et al. 1994].

## Clusters

Clusters were probably “invented” in the 1960s by customers who could not fit all their work on one computer or who needed a backup machine in case of failure of the primary machine [Pfister 1998]. Tandem introduced a 16-node cluster in 1975. Digital followed with VAX clusters, introduced in 1984. They were originally independent computers that shared I/O devices, requiring a distributed operating system to coordinate activity. Soon they had communication links between computers, in part so that the computers could be geographically distributed to increase availability in case of a disaster at a single site. Users log onto the cluster and are unaware of which machine they are running on. DEC (now HP) sold more than 25,000 clusters by 1993. Other early companies were Tandem (now HP) and IBM (still IBM). Today, virtually every company has cluster products. Most of these products are aimed at availability, with performance scaling as a secondary benefit.

Scientific computing on clusters emerged as a competitor to MPPs. In 1993, the Beowulf project started with the goal of fulfilling NASA’s desire for a 1 GFLOPS computer for under \$50,000. In 1994, a 16-node cluster built from off-the-shelf

PCs using 80486s achieved that goal [Bell and Gray 2001]. This emphasis led to a variety of software interfaces to make it easier to submit, coordinate, and debug large programs or a large number of independent programs.

Efforts were made to reduce latency of communication in clusters as well as to increase bandwidth, and several research projects worked on that problem. (One commercial result of the low-latency research was the VI interface standard, which has been embraced by Infiniband, discussed below.) Low latency then proved useful in other applications. For example, in 1997 a cluster of 100 Ultra-SPARC desktop computers at the University of California–Berkeley, connected by 160 MB/sec per link Myrinet switches, was used to set world records in database sort—sorting 8.6 GB of data originally on disk in 1 minute—and in cracking an encrypted message—taking just 3.5 hours to decipher a 40-bit DES key.

This research project, called Network of Workstations [Anderson, Culler, and Patterson 1995], also developed the Inktomi search engine, which led to a startup company with the same name. Google followed the example of Inktomi to build search engines from clusters of desktop computers rather large-scale SMPs, which was the strategy of the leading search engine Alta Vista that Google overtook [Brin and Page 1998]. In 2011, nearly all Internet services rely on clusters to serve their millions of customers.

Clusters are also very popular with scientists. One reason is their low cost, so individual scientists or small groups can own a cluster dedicated to their programs. Such clusters can get results faster than waiting in the long job queues of the shared MPPs at supercomputer centers, which can stretch to weeks. For those interested in learning more, Pfister [1998] wrote an entertaining book on clusters.

### *Recent Trends in Large-Scale Multiprocessors*

In the mid- to late 1990s, it became clear that the hoped for growth in the market for ultralarge-scale parallel computing was unlikely to occur. Without this market growth, it became increasingly clear that the high-end parallel computing market could not support the costs of highly customized hardware and software designed for a small market. Perhaps the most important trend to come out of this observation was that clustering would be used to reach the highest levels of performance. There are now four general classes of large-scale multiprocessors:

- Clusters that integrate standard desktop motherboards using interconnection technology such as Myrinet or Infiniband.
- Multicomputers built from standard microprocessors configured into processing elements and connected with a custom interconnect. These include the Cray XT3 (which used an earlier version of Cray interconnect with a simple cluster architecture) and IBM Blue Gene (more on this unique machine momentarily).
- Clusters of small-scale shared-memory computers, possibly with vector support, which includes the Earth Simulator (which has its own journal available online).



- Large-scale shared-memory multiprocessors, such as the Cray X1 [Dunigan et al. 2005] and SGI Origin and Altix systems. The SGI systems have also been configured into clusters to provide more than 512 processors, although only message passing is supported across the clusters.

The IBM Blue Gene is the most interesting of these designs since its rationale parallels the underlying causes of the recent trend toward multicore in uniprocessor architectures. Blue Gene started as a research project within IBM aimed at the protein sequencing and folding problem. The Blue Gene designers observed that power was becoming an increasing concern in large-scale multiprocessors and that the performance/watt of processors from the embedded space was much better than those in the high-end uniprocessor space. If parallelism was the route to high performance, why not start with the most efficient building block and simply have more of them?

Thus, Blue Gene is constructed using a custom chip that includes an embedded PowerPC microprocessor offering half the performance of a high-end PowerPC, but at a much smaller fraction of the area of power. This allows more system functions, including the global interconnect, to be integrated onto the same die. The result is a highly replicable and efficient building block, allowing Blue Gene to reach much larger processor counts more efficiently. Instead of using stand-alone microprocessors or standard desktop boards as building blocks, Blue Gene uses processor cores. There is no doubt that such an approach provides much greater efficiency. Whether the market can support the cost of a customized design and special software remains an open question.

In 2006, a Blue Gene processor at Lawrence Livermore with 32K processors (and scheduled to go to 65K in late 2005) holds a factor of 2.6 lead in Linpack performance over the third-place system consisting of 20 SGI Altix 512-processor systems interconnected with Infiniband as a cluster.

Blue Gene's predecessor was an experimental machine, QCDOD, which pioneered the concept of a machine using a lower-power embedded microprocessor and tightly integrated interconnect to drive down the cost and power consumption of a node.

### *Developments in Synchronization and Consistency Models*

A wide variety of synchronization primitives have been proposed for shared-memory multiprocessors. Mellor-Crummey and Scott [1991] provided an overview of the issues as well as efficient implementations of important primitives, such as locks and barriers. An extensive bibliography supplies references to other important contributions, including developments in spin locks, queuing locks, and barriers. Lamport [1979] introduced the concept of sequential consistency and what correct execution of parallel programs means. Dubois, Scheurich, and Briggs [1988] introduced the idea of weak ordering (originally in 1986). In 1990, Adve and Hill provided a better definition of weak ordering and also defined the concept of data-race-free; at the same conference, Gharachorloo and his colleagues [1990] introduced release consistency and provided the first data on the performance of

relaxed consistency models. More relaxed consistency models have been widely adopted in microprocessor architectures, including the Sun SPARC, Alpha, and IA-64. Adve and Gharachorloo [1996] have provided an excellent tutorial on memory consistency and the differences among these models.

## Other References

The concept of using virtual memory to implement a shared address space among distinct machines was pioneered in Kai Li's Ivy system in 1988. There have been subsequent papers exploring hardware support issues, software mechanisms, and programming issues. Amza et al. [1996] described a system built on workstations using a new consistency model, Kontothanassis et al. [1997] described a software shared-memory scheme using remote writes, and Erlichson et al. [1996] described the use of shared virtual memory to build large-scale multiprocessors using SMPs as nodes.

There is an almost unbounded amount of information on multiprocessors and multicomputers: Conferences, journal papers, and even books seem to appear faster than any single person can absorb the ideas. No doubt many of these papers will go unnoticed—not unlike the past. Most of the major architecture conferences contain papers on multiprocessors. An annual conference, Supercomputing *XY* (where *X* and *Y* are the last two digits of the year), brings together users, architects, software developers, and vendors, and the proceedings are published in book, CD-ROM, and online (see [www.scXY.org](http://www.scXY.org)) form. Two major journals, *Journal of Parallel and Distributed Computing* and the *IEEE Transactions on Parallel and Distributed Systems*, contain papers on all aspects of parallel processing. Several books focusing on parallel processing are included in the following references, with Culler, Singh, and Gupta [1999] being the most recent, large-scale effort. For years, Eugene Miya of NASA Ames Research Center has collected an online bibliography of parallel-processing papers. The bibliography, which now contains more than 35,000 entries, is available online at [linwww.ira.uka.de/bibliography/Parallel/Eugene/index.html](http://linwww.ira.uka.de/bibliography/Parallel/Eugene/index.html).

In addition to documenting the discovery of concepts now used in practice, these references also provide descriptions of many ideas that have been explored and found wanting, as well as ideas whose time has just not yet come. Given the move toward multicore and multiprocessors as the future of high-performance computer architecture, we expect that many new approaches will be explored in the years ahead. A few of them will manage to solve the hardware and software problems that have been the key to using multiprocessing for the past 40 years!

## References

- Adve, S. V., and K. Gharachorloo [1996]. "Shared memory consistency models: A tutorial," *IEEE Computer* 29:12 (December), 66–76.
- Adve, S. V., and M. D. Hill [1990]. "Weak ordering—a new definition," *Proc. 17th Annual Int'l. Symposium on Computer Architecture (ISCA)*, May 28–31, 1990, Seattle, Wash., 2–14.

- Agarwal, A., R. Bianchini, D. Chaiken, K. Johnson, and D. Kranz [1995]. "The MIT Alewife machine: Architecture and performance," *22nd Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 22–24, 1995, Santa Margherita, Italy, 2–13.
- Agarwal, A., J. L. Hennessy, R. Simoni, and M. A. Horowitz [1988]. "An evaluation of directory schemes for cache coherence," *Proc. 15th Annual Int'l. Symposium on Computer Architecture*, May 30–June 2, 1988, Honolulu, Hawaii, 280–289.
- Agarwal, A., J. Kubiawicz, D. Kranz, B.-H. Lim, D. Yeung, G. D'Souza, and M. Parkin [1993]. "Sparcle: An evolutionary processor design for large-scale multiprocessors," *IEEE Micro* 13 (June), 48–61.
- Alles, A. [1995]. "ATM Internetworking," White Paper (May), Cisco Systems, Inc., San Jose, Calif. ([www.cisco.com/warp/public/614/12.html](http://www.cisco.com/warp/public/614/12.html)).
- Almasi, G. S., and A. Gottlieb [1989]. *Highly Parallel Computing*, Benjamin/Cummings, Redwood City, Calif.
- Alverson, G., R. Alverson, D. Callahan, B. Koblenz, A. Porterfield, and B. Smith [1992]. "Exploiting heterogeneous parallelism on a multithreaded multiprocessor," *Proc. ACM/IEEE Conf. on Supercomputing*, November 16–20, 1992, Minneapolis, Minn., 188–197.
- Amdahl, G. M. [1967]. "Validity of the single processor approach to achieving large scale computing capabilities," *Proc. AFIPS Spring Joint Computer Conf.*, April 18–20, 1967, Atlantic City, N.J., 483–485.
- Amza C., A. L. Cox, S. Dwarkadas, P. Keleher, H. Lu, R. Rajamony, W. Yu, and W. Zwaenepoel [1996]. "Treadmarks: Shared memory computing on networks of workstations," *IEEE Computer* 29:2 (February), 18–28.
- Anderson, T. E., D. E. Culler, and D. Patterson [1995]. "A case for NOW (networks of workstations)," *IEEE Micro* 15:1 (February), 54–64.
- Ang, B., D. Chiou, D. Rosenband, M. Ehrlich, L. Rudolph, and Arvind [1998]. "StarT-Voyager: A flexible platform for exploring scalable SMP issues," *Proc. ACM/IEEE Conf. on Supercomputing*, November 7–13, 1998, Orlando, FL.
- Archibald, J., and J.-L. Baer [1986]. "Cache coherence protocols: Evaluation using a multiprocessor simulation model," *ACM Trans. on Computer Systems* 4:4 (November), 273–298.
- Arpaci, R. H., D. E. Culler, A. Krishnamurthy, S. G. Steinberg, and K. Yelick [1995]. "Empirical evaluation of the CRAY-T3D: A compiler perspective," *Proc. 22nd Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 22–24, 1995, Santa Margherita, Italy.
- Baer, J.-L., and W.-H. Wang [1988]. "On the inclusion properties for multi-level cache hierarchies," *Proc. 15th Annual Int'l. Symposium on Computer Architecture*, May 30–June 2, 1988, Honolulu, Hawaii, 73–80.
- Balakrishnan, H. V., N. Padmanabhan, S. Seshan, and R. H. Katz [1997]. "A comparison of mechanisms for improving TCP performance over wireless links," *IEEE/ACM Trans. on Networking* 5:6 (December), 756–769.
- Barroso, L. A., K. Gharachorloo, and E. Bugnion [1998]. "Memory system characterization of commercial workloads," *Proc. 25th Annual Int'l. Symposium on Computer Architecture (ISCA)*, July 3–14, 1998, Barcelona, Spain, 3–14.

- Baskett, F., T. Jermoluk, and D. Solomon [1988]. "The 4D-MP graphics super-workstation: Computing + graphics = 40 MIPS + 40 MFLOPS and 10,000 lighted polygons per second," *Proc. IEEE COMPCON*, February 29–March 4, 1988, San Francisco, 468–471.
- BBN Laboratories. [1986]. *Butterfly Parallel Processor Overview*, Tech. Rep. 6148, BBN Laboratories, Cambridge, Mass.
- Bell, C. G. [1985]. "Multis: A new class of multiprocessor computers," *Science* 228 (April 26), 462–467.
- Bell, C. G. [1989]. "The future of high performance computers in science and engineering," *Communications of the ACM* 32:9 (September), 1091–1101.
- Bell, C. G., and J. Gray [2001]. *Crays, Clusters and Centers*, Tech. Rep. MSR-TR-2001-76, Microsoft Research, Redmond, Wash.
- Bell, C. G., and J. Gray [2002]. "What's next in high performance computing," *CACM*, 45:2 (February), 91–95.
- Bouknight, W. J., S. A. Deneberg, D. E. McIntyre, J. M. Randall, A. H. Sameh, and D. L. Slotnick [1972]. "The Illiac IV system," *Proc. IEEE* 60:4, 369–379. Also appears in D. P. Siewiorek, C. G. Bell, and A. Newell, *Computer Structures: Principles and Examples*, McGraw-Hill, New York, 1982, 306–316.
- Brain, M. [2000]. *Inside a Digital Cell Phone*, [www.howstuffworks.com/inside-cell-phone.htm](http://www.howstuffworks.com/inside-cell-phone.htm).
- Brewer, E. A., and B. C. Kuszmaul [1994]. "How to get good performance from the CM-5 data network," *Proc. Eighth Int'l. Parallel Processing Symposium (IPPS)*, April 26–29, 1994, Cancun, Mexico.
- Brin, S., and L. Page [1998]. "The anatomy of a large-scale hypertextual Web search engine," *Proc. 7th Int'l. World Wide Web Conf.*, April 14–18, 1998, Brisbane, Queensland, Australia, 107–117.
- Burkhardt III, H., S. Frank, B. Knobe, and J. Rothnie [1992]. *Overview of the KSR1 Computer System*, Tech. Rep. KSR-TR-9202001, Kendall Square Research, Boston.
- Censier, L., and P. Feautrier [1978]. "A new solution to coherence problems in multicache systems," *IEEE Trans. on Computers* C-27:12 (December), 1112–1118.
- Chandra, R., S. Devine, B. Verghese, A. Gupta, and M. Rosenblum [1994]. "Scheduling and page migration for multiprocessor compute servers," *Proc. Sixth Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, October 4–7, 1994, San Jose, Calif., 12–24.
- Charlesworth, A. [1998]. "Starfire: Extending the SMP envelope," *IEEE Micro* 18:1 (January/February), 39–49.
- Clark, W. A. [1957]. "The Lincoln TX-2 computer development," *Proc. Western Joint Computer Conference*, February 26–28, 1957, Los Angeles, 143–145.
- Comer, D. [1993]. *Internetworking with TCP/IP*, 2nd ed., Prentice Hall, Englewood Cliffs, N.J.
- Culler, D. E., J. P. Singh, and A. Gupta [1999]. *Parallel Computer Architecture: A Hardware/Software Approach*, Morgan Kaufmann, San Francisco.
- Dally, W. J., and C. I. Seitz [1986]. "The torus routing chip," *Distributed Computing* 1:4, 187–196.

- Davie, B. S., L. L. Peterson, and D. Clark [1999]. *Computer Networks: A Systems Approach*, 2nd ed., Morgan Kaufmann, San Francisco.
- Desurvire, E. [1992]. "Lightwave communications: The fifth generation," *Scientific American* (International Edition) 266:1 (January), 96–103.
- Dongarra, J., T. Sterling, H. Simon, and E. Strohmaier [2005]. "High-performance computing: Clusters, constellations, MPPs, and future directions," *Computing in Science & Engineering*, 7:2 (March/April), 51–59.
- Dubois, M., C. Scheurich, and F. Briggs [1988]. "Synchronization, coherence, and event ordering," *IEEE Computer* 21:2 (February), 9–21.
- Dunigan, W., K. Vetter, K. White, and P. Worley [2005]. "Performance evaluation of the Cray X1 distributed shared memory architecture," *IEEE Micro*, January/February, 30–40.
- Eggers, S. [1989]. "Simulation Analysis of Data Sharing in Shared Memory Multiprocessors," Ph.D. thesis, Computer Science Division, University of California, Berkeley.
- Elder, J., A. Gottlieb, C. K. Kruskal, K. P. McAuliffe, L. Randolph, M. Snir, P. Teller, and J. Wilson [1985]. "Issues related to MIMD shared-memory computers: The NYU Ultracomputer approach," *Proc. 12th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 17–19, 1985, Boston, Mass., 126–135.
- Erlichson, A., N. Nuckolls, G. Chesson, and J. L. Hennessy [1996]. "SoftFLASH: Analyzing the performance of clustered distributed virtual shared memory," *Proc. Seventh Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, October 1–5, 1996, Cambridge, Mass., 210–220.
- Falsafi, B., and D. A. Wood [1997]. "Reactive NUMA: A design for unifying S-COMA and CC-NUMA," *Proc. 24th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 2–4, 1997, Denver, Colo., 229–240.
- Flynn, M. J. [1966]. "Very high-speed computing systems," *Proc. IEEE* 54:12 (December), 1901–1909.
- Forgie, J. W. [1957]. "The Lincoln TX-2 input-output system," *Proc. Western Joint Computer Conference*, February 26–28, 1957, Los Angeles, 156–160.
- Frank, S. J. [1984]. "Tightly coupled multiprocessor systems speed memory access time," *Electronics* 57:1 (January), 164–169.
- Gajski, D., D. Kuck, D. Lawrie, and A. Sameh [1983]. "CEDAR—a large scale multiprocessor," *Proc. Int'l. Conf. on Parallel Processing (ICPP)*, August, Columbus, Ohio, 524–529.
- Galles, M. [1996]. "Scalable pipelined interconnect for distributed endpoint routing: The SGI SPIDER chip," *Proc. IEEE HOT Interconnects '96*, August 15–17, 1996, Stanford University, Palo Alto, Calif.
- Gehringer, E. F., D. P. Siewiorek, and Z. Segall [1987]. *Parallel Processing: The Cm\* Experience*, Digital Press, Bedford, Mass.
- Gharachorloo, K., A. Gupta, and J. L. Hennessy [1992]. "Hiding memory latency using dynamic scheduling in shared-memory multiprocessors," *Proc. 19th Annual Int'l. Symposium on Computer Architecture (ISCA)*, May 19–21, 1992, Gold Coast, Australia.

- Gharachorloo, K., D. Lenoski, J. Laudon, P. Gibbons, A. Gupta, and J. L. Hennessy [1990]. "Memory consistency and event ordering in scalable shared-memory multiprocessors," *Proc. 17th Annual Int'l. Symposium on Computer Architecture (ISCA)*, May 28–31, 1990, Seattle, Wash., 15–26.
- Gibson, J., R. Kunz, D. Ofelt, M. Horowitz, J. Hennessy, and M. Heinrich [2000]. "FLASH vs. (simulated) FLASH: Closing the simulation loop," *Proc. Ninth Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, November 12–15, Cambridge, Mass., 49–58.
- Goodman, J. R. [1983]. "Using cache memory to reduce processor memory traffic," *Proc. 10th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 5–7, 1982, Stockholm, Sweden, 124–131.
- Goralski, W. [1997]. *SONET: A Guide to Synchronous Optical Network*, McGraw-Hill, New York.
- Grice, C., and M. Kanellos [2000]. "Cell phone industry at crossroads: Go high or low?" *CNET News* (August 31), [technews.netscape.com/news/0-1004-201-2518386-0.html?tag=st.ne.1002.gif.sf](http://technews.netscape.com/news/0-1004-201-2518386-0.html?tag=st.ne.1002.gif.sf).
- Groe, J. B., and L. E. Larson [2000]. *CDMA Mobile Radio Design*, Artech House, Boston.
- Hagersten E., and M. Koster [1998]. "WildFire: A scalable path for SMPs," *Proc. Fifth Int'l. Symposium on High-Performance Computer Architecture*, January 9–12, 1999, Orlando, Fla.
- Hagersten, E., A. Landin, and S. Haridi [1992]. "DDM—a cache-only memory architecture," *IEEE Computer* 25:9 (September), 44–54.
- Hill, M. D. [1998]. "Multiprocessors should support simple memory consistency models," *IEEE Computer* 31:8 (August), 28–34.
- Hillis, W. D. [1985]. *The Connection Multiprocessor*, MIT Press, Cambridge, Mass.
- Hirata, H., K. Kimura, S. Nagamine, Y. Mochizuki, A. Nishimura, Y. Nakase, and T. Nishizawa [1992]. "An elementary processor architecture with simultaneous instruction issuing from multiple threads," *Proc. 19th Annual Int'l. Symposium on Computer Architecture (ISCA)*, May 19–21, 1992, Gold Coast, Australia, 136–145.
- Hockney, R. W., and C. R. Jesshope [1988]. *Parallel Computers 2: Architectures, Programming and Algorithms*, Adam Hilger, Ltd., Bristol, England.
- Holland, J. H. [1959]. "A universal computer capable of executing an arbitrary number of subprograms simultaneously," *Proc. East Joint Computer Conf.* 16, 108–113.
- Hord, R. M. [1982]. *The Illiac-IV, The First Supercomputer*, Computer Science Press, Rockville, Md.
- Hristea, C., D. Lenoski, and J. Keen [1997]. "Measuring memory hierarchy performance of cache-coherent multiprocessors using micro benchmarks," *Proc. ACM/IEEE Conf. on Supercomputing*, November 15–21, 1997, San Jose, Calif.
- Hwang, K. [1993]. *Advanced Computer Architecture and Parallel Programming*, McGraw-Hill, New York.



- IBM. [2005]. "Blue Gene," *IBM J. of Research and Development*, 49:2/3 (special issue).
- Infiniband Trade Association. [2001]. *InfiniBand Architecture Specifications Release 1.0.a*, [www.infinibandta.org](http://www.infinibandta.org).
- Jordan, H. F. [1983]. "Performance measurements on HEP—a pipelined MIMD computer," *Proc. 10th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 5–7, 1982, Stockholm, Sweden, 207–212.
- Kahn, R. E. [1972]. "Resource-sharing computer communication networks," *Proc. IEEE* 60:11 (November), 1397–1407.
- Keckler, S. W., and W. J. Dally [1992]. "Processor coupling: Integrating compile time and runtime scheduling for parallelism," *Proc. 19th Annual Int'l. Symposium on Computer Architecture (ISCA)*, May 19–21, 1992, Gold Coast, Australia, 202–213.
- Kontothanassis, L., G. Hunt, R. Stets, N. Hardavellas, M. Cierniak, S. Parthasarathy, W. Meira, S. Dwarkadas, and M. Scott [1997]. "VM-based shared memory on low-latency, remotememory-access networks," *Proc. 24th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 2–4, 1997, Denver, Colo.
- Kurose, J. F., and K. W. Ross [2001]. *Computer Networking: A Top-Down Approach Featuring the Internet*, Addison-Wesley, Boston.
- Kuskin, J., D. Ofelt, M. Heinrich, J. Heinlein, R. Simoni, K. Gharachorloo, J. Chapin, D. Nakahira, J. Baxter, M. Horowitz, A. Gupta, M. Rosenblum, and J. L. Hennessy [1994]. "The Stanford FLASH multiprocessor," *Proc. 21st Annual Int'l. Symposium on Computer Architecture (ISCA)*, April 18–21, 1994, Chicago.
- Lamport, L. [1979]. "How to make a multiprocessor computer that correctly executes multiprocess programs," *IEEE Trans. on Computers* C-28:9 (September), 241–248.
- Laudon, J., A. Gupta, and M. Horowitz [1994]. "Interleaving: A multithreading technique targeting multiprocessors and workstations," *Proc. Sixth Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, October 4–7, 1994, San Jose, Calif., 308–318.
- Laudon, J., and D. Lenoski [1997]. "The SGI Origin: A ccNUMA highly scalable server," *Proc. 24th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 2–4, 1997, Denver, Colo., 241–251.
- Lenoski, D., J. Laudon, K. Gharachorloo, A. Gupta, and J. L. Hennessy [1990]. "The Stanford DASH multiprocessor," *Proc. 17th Annual Int'l. Symposium on Computer Architecture (ISCA)*, May 28–31, 1990, Seattle, Wash., 148–159.
- Lenoski, D., J. Laudon, K. Gharachorloo, W.-D. Weber, A. Gupta, J. L. Hennessy, M. A. Horowitz, and M. Lam [1992]. "The Stanford DASH multiprocessor," *IEEE Computer* 25:3 (March), 63–79.
- Li, K. [1988]. "IVY: A shared virtual memory system for parallel computing," *Proc. Int'l. Conf. on Parallel Processing (ICPP)*, August, The Pennsylvania State University, University Park, Penn.
- Lo, J., L. Barroso, S. Eggers, K. Gharachorloo, H. Levy, and S. Parekh [1998]. "An analysis of database workload performance on simultaneous multithreaded

- processors,” *Proc. 25th Annual Int’l. Symposium on Computer Architecture (ISCA)*, July 3–14, 1998, Barcelona, Spain, 39–50.
- Lo, J., S. Eggers, J. Emer, H. Levy, R. Stamm, and D. Tullsen [1997]. “Converting thread-level parallelism into instruction-level parallelism via simultaneous multithreading,” *ACM Trans. on Computer Systems* 15:2 (August), 322–354.
- Lovett, T., and S. Thakkar [1988]. “The Symmetry multiprocessor system,” *Proc. Int’l. Conf. on Parallel Processing (ICPP)*, August, The Pennsylvania State University, University Park, Penn., 303–310.
- Mellor-Crummey, J. M., and M. L. Scott [1991]. “Algorithms for scalable synchronization on shared-memory multiprocessors,” *ACM Trans. on Computer Systems* 9:1 (February), 21–65.
- Menabrea, L. F. [1842]. “Sketch of the analytical engine invented by Charles Babbage,” *Bibliothèque Universelle de Genève*, 82 (October).
- Metcalf, R. M. [1993]. “Computer/network interface design: Lessons from Arpanet and Ethernet,” *IEEE J. on Selected Areas in Communications* 11:2 (February), 173–180.
- Metcalf, R. M., and D. R. Boggs [1976]. “Ethernet: Distributed packet switching for local computer networks,” *Communications of the ACM* 19:7 (July), 395–404.
- Mitchell, D. [1989]. “The Transputer: The time is now,” *Computer Design (RISC suppl.)*, 40–41.
- Miya, E. N. [1985]. “Multiprocessor/distributed processing bibliography,” *Computer Architecture News* 13:1, 27–29.
- National Research Council. [1997]. *The Evolution of Untethered Communications*, Computer Science and Telecommunications Board, National Academy Press, Washington, D.C.
- Nikhil, R. S., G. M. Papadopoulos, and Arvind [1992]. “\*T: A multithreaded massively parallel architecture,” *Proc. 19th Annual Int’l. Symposium on Computer Architecture (ISCA)*, May 19–21, 1992, Gold Coast, Australia, 156–167.
- Noordergraaf, L., and R. van der Pas [1999]. “Performance experiences on Sun’s WildFire prototype,” *Proc. ACM/IEEE Conf. on Supercomputing*, November 13–19, 1999, Portland, Ore.
- Partridge, C. [1994]. *Gigabit Networking*, Addison-Wesley, Reading, Mass.
- Pfister, G. F. [1998]. *In Search of Clusters*, 2nd ed., Prentice Hall, Upper Saddle River, N.J.
- Pfister, G. F., W. C. Brantley, D. A. George, S. L. Harvey, W. J. Kleinfekder, K. P. McAuliffe, E. A. Melton, V. A. Norton, and J. Weiss [1985]. “The IBM research parallel processor prototype (RP3): Introduction and architecture,” *Proc. 12th Annual Int’l. Symposium on Computer Architecture (ISCA)*, June 17–19, 1985, Boston, Mass., 764–771.
- Reinhardt, S. K., J. R. Larus, and D. A. Wood [1994]. “Tempest and Typhoon: User-level shared memory,” *Proc. 21st Annual Int’l. Symposium on Computer Architecture (ISCA)*, April 18–21, 1994, Chicago, 325–336.
- Rettberg, R. D., W. R. Crowther, P. P. Carvey, and R. S. Towlinson [1990]. “The Monarch parallel processor hardware design,” *IEEE Computer* 23:4 (April), 18–30.

- Rosenblum, M., S. A. Herrod, E. Witchel, and A. Gupta [1995]. "Complete computer simulation: The SimOS approach," in *IEEE Parallel and Distributed Technology* (now called *Concurrency*) 4:3, 34–43.
- Saltzer, J. H., D. P. Reed, and D. D. Clark [1984]. "End-to-end arguments in system design," *ACM Trans. on Computer Systems* 2:4 (November), 277–288.
- Satran, J., D. Smith, K. Meth, C. Sapuntzakis, M. Wakeley, P. Von Stamwitz, R. Haagens, E. Zeidner, L. Dalle Ore, and Y. Klein [2001]. "iSCSI," IPS Working Group of IETF, Internet draft [www.ietf.org/internet-drafts/draft-ietf-ips-iscsi-07.txt](http://www.ietf.org/internet-drafts/draft-ietf-ips-iscsi-07.txt).
- Saulsbury, A., T. Wilkinson, J. Carter, and A. Landin [1995]. "An argument for Simple COMA," *Proc. First IEEE Symposium on High-Performance Computer Architectures*, January 22–25, 1995, Raleigh, N.C., 276–285.
- Schwartz, J. T. [1980]. "Ultracomputers," *ACM Trans. on Programming Languages and Systems* 4:2, 484–521.
- Scott, S. L. [1996]. "Synchronization and communication in the T3E multiprocessor," *Seventh Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, October 1–5, 1996, Cambridge, Mass., 26–36.
- Scott, S. L., and G. M. Thorson [1996]. "The Cray T3E network: Adaptive routing in a high-performance 3D torus," *Proc. IEEE HOT Interconnects '96*, August 15–17, 1996, Stanford University, Palo Alto, Calif., 14–156.
- Seitz, C. L. [1985]. "The Cosmic Cube (concurrent computing)," *Communications of the ACM* 28:1 (January), 22–33.
- Singh, J. P., J. L. Hennessy, and A. Gupta [1993]. "Scaling parallel programs for multiprocessors: Methodology and examples," *Computer* 26:7 (July), 22–33.
- Slotnick, D. L., W. C. Borck, and R. C. McReynolds [1962]. "The Solomon computer," *Proc. AFIPS Fall Joint Computer Conf.*, December 4–6, 1962, Philadelphia, Penn., 97–107.
- Smith, B. J. [1978]. "A pipelined, shared resource MIMD computer," *Proc. Int'l. Conf. on Parallel Processing (ICPP)*, August, Bellaire, Mich., 6–8.
- Soundararajan, V., M. Heinrich, B. Verghese, K. Gharachorloo, A. Gupta, and J. L. Hennessy [1998]. "Flexible use of memory for replication/migration in cache-coherent DSM multiprocessors," *Proc. 25th Annual Int'l. Symposium on Computer Architecture (ISCA)*, July 3–14, 1998, Barcelona, Spain, 342–355.
- Spurgeon, C. [2001]. "Charles Spurgeon's Ethernet Web site," [www.host.ots.utexas.edu/ethernet/ethernet-home.html](http://www.host.ots.utexas.edu/ethernet/ethernet-home.html).
- Stenström, P., T. Joe, and A. Gupta [1992]. "Comparative performance evaluation of cachecoherent NUMA and COMA architectures," *Proc. 19th Annual Int'l. Symposium on Computer Architecture (ISCA)*, May 19–21, 1992, Gold Coast, Australia, 80–91.
- Sterling, T. [2001]. *Beowulf PC Cluster Computing with Windows and Beowulf PC Cluster Computing with Linux*, MIT Press, Cambridge, Mass.
- Stevens, W. R. [1994–1996]. *TCP/IP Illustrated* (three volumes), Addison-Wesley, Reading, Mass.
- Stone, H. [1991]. *High Performance Computers*, Addison-Wesley, New York.

- Swan, R. J., A. Bechtolsheim, K. W. Lai, and J. K. Ousterhout [1977]. "The implementation of the Cm\* multi-microprocessor," *Proc. AFIPS National Computing Conf.*, June 13–16, 1977, Dallas, Tex., 645–654.
- Swan, R. J., S. H. Fuller, and D. P. Siewiorek [1977]. "Cm\*—a modular, multi-microprocessor," *Proc. AFIPS National Computing Conf.*, June 13–16, 1977, Dallas, Tex., 637–644.
- Tanenbaum, A. S. [1988]. *Computer Networks*, 2nd ed., Prentice Hall, Englewood Cliffs, N.J.
- Tang, C. K. [1976]. "Cache design in the tightly coupled multiprocessor system," *Proc. AFIPS National Computer Conf.*, June 7–10, 1976, New York, 749–753.
- Thacker, C. P., E. M. McCreight, B. W. Lampson, R. F. Sproull, and D. R. Boggs [1982]. "Alto: A personal computer," in D. P. Siewiorek, C. G. Bell, and A. Newell, eds., *Computer Structures: Principles and Examples*, McGraw-Hill, New York, 549–572.
- Thekkath, R., A. P. Singh, J. P. Singh, S. John, and J. L. Hennessy [1997]. "An evaluation of a commercial CC-NUMA architecture—the CONVEX Exemplar SPP1200," *Proc. 11th Int'l. Parallel Processing Symposium (IPPS)*, April 1–7, 1997, Geneva, Switzerland.
- Tullsen, D. M., S. J. Eggers, J. S. Emer, H. M. Levy, J. L. Lo, and R. L. Stamm [1996]. "Exploiting choice: Instruction fetch and issue on an implementable simultaneous multithreading processor," *Proc. 23rd Annual Int'l. Symposium on Computer Architecture (ISCA)*, May 22–24, 1996, Philadelphia, Penn., 191–202.
- Tullsen, D. M., S. J. Eggers, and H. M. Levy [1995]. "Simultaneous multithreading: Maximizing on-chip parallelism," *Proc. 22nd Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 22–24, 1995, Santa Margherita, Italy, 392–403.
- Unger, S. H. [1958]. "A computer oriented towards spatial problems," *Proc. Institute of Radio Engineers* 46:10 (October), 1744–1750.
- Walrand, J. [1991]. *Communication Networks: A First Course*, Aksen Associates: Irwin, Homewood, Ill.
- Wilson, A. W., Jr. [1987]. "Hierarchical cache/bus architecture for shared-memory multiprocessors," *Proc. 14th Annual Int'l. Symposium on Computer Architecture (ISCA)*, June 2–5, 1987, Pittsburgh, Penn., 244–252.
- Wolfe, A., and J. P. Shen [1991]. "A variable instruction stream extension to the VLIW architecture," *Proc. Fourth Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, April 8–11, 1991, Palo Alto, Calif., 2–14.
- Wood, D. A., and M. D. Hill [1995]. "Cost-effective parallel computing," *IEEE Computer* 28:2 (February), 69–72.
- Wulf, W., and C. G. Bell [1972]. "C.mmp—A multi-mini-processor," *Proc. AFIPS Fall Joint Computer Conf.*, December 5–7, 1972, Anaheim, Calif., 765–777.
- Wulf, W., and S. P. Harbison [1978]. "Reflections in a pool of processors—an experience report on C.mmp/Hydra," *Proc. AFIPS National Computing Conf.* June 5–8, 1978, Anaheim, Calif., 939–951.

Yamamoto, W., M. J. Serrano, A. R. Talcott, R. C. Wood, and M. Nemirosky [1994]. "Performance estimation of multistreamed, superscalar processors," *Proc. 27th Hawaii Int'l. Conf. on System Sciences*, January 4–7, 1994, Wailea, 195–204.

---

## M.8

### The Development of Clusters (Chapter 6)

In this section, we cover the development of clusters that were the foundation of warehouse-scale computers (WSCs) and of utility computing. (Readers interested in learning more should start with Barroso and Hölzle [2009] and the blog postings and talks of James Hamilton at <http://perspectives.mvdirona.com>.)

#### Clusters, the Forerunner of WSCs

Clusters were probably "invented" in the 1960s by customers who could not fit all their work on one computer or who needed a backup machine in case of failure of the primary machine [Pfister 1998]. Tandem introduced a 16-node cluster in 1975. Digital followed with VAX clusters, introduced in 1984. They were originally independent computers that shared I/O devices, requiring a distributed operating system to coordinate activity. Soon they had communication links between computers, in part so that the computers could be geographically distributed to increase availability in case of a disaster at a single site. Users log onto the cluster and are unaware of which machine they are running on. DEC (now HP) sold more than 25,000 clusters by 1993. Other early companies were Tandem (now HP) and IBM (still IBM). Today, virtually every company has cluster products. Most of these products are aimed at availability, with performance scaling as a secondary benefit.

Scientific computing on clusters emerged as a competitor to MPPs. In 1993, the Beowulf project started with the goal of fulfilling NASA's desire for a 1 GFLOPS computer for under \$50,000. In 1994, a 16-node cluster built from off-the-shelf PCs using 80486s achieved that goal [Bell and Gray 2001]. This emphasis led to a variety of software interfaces to make it easier to submit, coordinate, and debug large programs or a large number of independent programs.

Efforts were made to reduce latency of communication in clusters as well as to increase bandwidth, and several research projects worked on that problem. (One commercial result of the low-latency research was the VI interface standard, which has been embraced by Infiniband, discussed below.) Low latency then proved useful in other applications. For example, in 1997 a cluster of 100 UltraSPARC desktop computers at the University of California–Berkeley, connected by 160 MB/sec per link Myrinet switches, was used to set world records in database sort—sorting 8.6 GB of data originally on disk in 1 minute—and in cracking an encrypted message—taking just 3.5 hours to decipher a 40-bit DES key.

This research project, called Network of Workstations [Anderson, Culler, and Patterson 1995], also developed the Inktomi search engine, which led to a start-up

company with the same name. Eric Brewer led the Inktomi effort at Berkeley and then at the company to demonstrate the use of commodity hardware to build computing infrastructure for Internet services. Using standardized networks within a rack of PC servers gave Inktomi better scalability. In contrast, the strategy of the prior leading search engine Alta Vista was to build from large-scale SMPs. Compared to the high-performance computing work in clusters, the emphasis was on a relatively large number of low-cost nodes and a clear programming model. Hence, the NOW project and Inktomi are considered the foundation of WSCs and Cloud Computing. Google followed the example of Inktomi technology when it took the leading search engine mantle from Inktomi just as Inktomi had taken it from Alta Vista [Brin and Page 1998]. (Google's initial innovation was search quality; the WSC innovations came much later.) For many years now, all Internet services have relied on cluster technology to serve their millions of customers.

### Utility Computing, the Forerunner of Cloud Computing

As stated in the text, the earliest version of utility computing was timesharing. Although timesharing faded away over time with the creation of smaller and cheaper personal computers, in the last decade there have been many less than fully successful attempts to resuscitate utility computing. Sun began selling time on Sun Cloud at \$1 per hour in 2000, HP offered a Utility Data Center in 2001, and Intel tried selling time on internal supercomputers in the early 2000s. Although they were commercially available, few customers used them.

A related topic is *grid computing*, which was originally invented so that scientific programs could be run across geographically distributed computing facilities. At the time, some questioned the wisdom of this goal, setting aside how difficult it would be to achieve. Grid computing tended to require very large systems running very large programs, using multiple datacenters for the tasks. Single applications did not really run well when geographically distributed, given the long latencies inherent with long distance. This first step eventually led to some conventions for data access, but the grid computing community did not develop APIs that were useful beyond the high-performance computing community, so the cloud computing effort shares little code or history with grid computing.

Armbrust et al [2009] argued that, once the Internet service companies solved the operational problems to work at large scale, the significant economies of scale that they uncovered brought their costs down below those of smaller datacenters. Amazon recognized that if this cost advantage was true then Amazon should be able to make a profit selling this service. In 2006, Amazon announced Elastic Cloud Computing (EC2) at \$0.10 per hour per instance. The subsequent popularity of EC2 led other Internet companies to offer cloud computing services, such as Google App Engine and Microsoft Azure, albeit at higher abstraction levels than the x86 virtual machines of Amazon Web Services. Hence, the current popularity of pay-as-you go computing isn't because someone recently came up with the idea;



it's because the technology and business models have aligned so that companies can make money offering a service that many people want to use. Time will tell whether there will be many successful utility computing models or whether the industry will converge around a single standard. It will certainly be interesting to watch.

## Containers

In the fall of 2003, many people were thinking about using containers to hold servers. Brewster Kahle, director and founder of the Internet Archive, gave talks about how he could fit the whole archive in a single 40-foot container. His interest was making copies of the Archive and distributing it around the world to ensure its survivability, thereby avoiding the fate of the Library of Alexandria that was destroyed by fire in 48 B.C.E. People working with Kahle wrote a white paper based on his talk in November 2003 to get more detail about what a container design would look like.

That same year, engineers at Google were also looking at building datacenters using containers and submitted a patent on aspects of it in December 2003. The first container for a datacenter was delivered in January 2005, and Google received the patent in October 2007. Google publicly revealed the use of containers in April 2009.

Greg Papadopolous of Sun Microsystems and Danny Hillis of Applied Minds heard Kahle's talk and designed a product called the Sun Modular Datacenter that debuted in October 2006. (The project code name was Black Box, a term many people still use.) This half-length (20-foot) container could hold 280 servers. This product release combined with Microsoft's announcement that they were building a datacenter designed to hold 220 40-foot containers inspired many other companies to offer containers and servers designed to be placed in them.

In a nice turn of events, in 2009 the Internet Archive migrated its data to a Sun Modular Datacenter. A copy of the Internet Archive is now at the New Library of Alexandria in Egypt, near the site of the original library.

## References

- Anderson, T. E., D. E. Culler, and D. Patterson [1995]. "A case for NOW (networks of workstations)," *IEEE Micro* 15:1 (February), 54–64.
- Apache Software Foundation. [2011]. Apache Hadoop project, <http://hadoop.apache.org>.
- Armbrust, M., A. Fox, R. Griffith, A.D. Joseph, R. Katz, A. Konwinski, G. Lee, D. Patterson, A. Rabkin, I. Stoica, and M. Zaharia [2009]. *Above the Clouds: A Berkeley View of Cloud Computing*, Tech. Rep. UCB/EECS-2009-28, University of California, Berkeley (<http://www.eecs.berkeley.edu/Pubs/TechRpts/2009/EECS-2009-28.html>).

- Barroso, L. A. [2010]. “Warehouse scale computing [keynote address],” *Proc. ACM SIG-MOD*, June 8–10, 2010, Indianapolis, Ind.
- Barroso, L. A., and U. Hözlze [2007]. “The case for energy-proportional computing,” *IEEE Computer* 40:12 (December), 33–37.
- Barroso, L.A., and U. Hözlze [2009]. “The datacenter as a computer: An introduction to the design of warehouse-scale machines,” in M. D. Hill, ed., *Synthesis Lectures on Computer Architecture*, Morgan & Claypool, San Rafael, Calif.
- Barroso, L.A., Clidaras, J. and Hözlze, U., 2013. *The datacenter as a computer: An introduction to the design of warehouse-scale machines*. Synthesis lectures on computer architecture, 8(3), pp.1–154.
- Barroso, L.A., Marty, M., Patterson, D., and Ranganathan, P. 2017. Attack of the Killer Microseconds. *Communications of the ACM*, 56(2).
- Bell, C. G., and J. Gray [2002]. “What’s next in high performance computing,” *Communications of the ACM* 45:2 (February), 91–95.
- Brady, J.T., 1986. A theory of productivity in the creative process. *IEEE Computer Graphics and Applications*, 6(5), pp.25–34.
- Brin, S., and L. Page [1998]. “The anatomy of a large-scale hypertextual Web search engine,” *Proc. 7th Int’l. World Wide Web Conf.*, April 14–18, 1998, Brisbane, Queensland, Australia, 107–117.
- Carter, J., and K. Rajamani [2010]. “Designing energy-efficient servers and data centers,” *IEEE Computer* 43:7 (July), 76–78.
- Chang, F., J. Dean, S. Ghemawat, W. C. Hsieh, D. A. Wallach, M. Burrows, T. Chandra, A. Fikes, and R. E. Gruber [2006]. “Bigtable: A distributed storage system for structured data,” in *Proc. Operating Systems Design and Implementation (OSDI ’06)*, November 6–8, 2006, Seattle, Wash.
- Chang, J., J. Meza, P. Ranganathan, C. Bash, and A. Shah [2010]. “Green server design: Beyond operational energy to sustainability,” *Workshop on Power Aware Computing and Systems (HotPower ’10)*, October 4–6, 2010, Vancouver, British Columbia.
- Clark, J., 2014 Five Numbers That Illustrate the Mind-Bending Size of Amazon’s Cloud, Bloomberg, <https://www.bloomberg.com/news/2014-11-14/5-numbers-that-illustrate-the-mind-bending-size-of-amazon-s-cloud.html>.
- Clidaras, J., C. Johnson, and B. Felderman [2010]. Private communication.
- Climate Savers Computing. [2007]. Efficiency specs, <http://www.climatesaverscomputing.org/>.
- Clos, C., 1953. A Study of Non-Blocking Switching Networks. *Bell Labs Technical Journal*, 32(2), pp.406–424.
- Dean, J. [2009]. “Designs, lessons and advice from building large distributed systems [keynote address],” *Proc. 3rd ACM SIGOPS International Workshop on Large Scale Distributed Systems and Middleware, Co-located with the 22nd ACM Symposium on Operating Systems Principles (SOSP 2009)*, October 10–11, 2009, Big Sky, Mont.
- Dean, J. and Barroso, L.A., 2013. The tail at scale. *Communications of the ACM*, 56(2), pp.74–80.

- Dean, J., and S. Ghemawat [2004]. “MapReduce: Simplified data processing on large clusters.” In *Proc. Operating Systems Design and Implementation (OSDI '04)*, December 6–8, 2004, San Francisco, 137–150.
- Dean, J., and S. Ghemawat [2008]. “MapReduce: simplified data processing on large clusters,” *Communications of the ACM* 51:1, 107–113.
- DeCandia, G., D. Hastorun, M. Jampani, G. Kakulapati, A. Lakshman, A. Pilchin, S. Sivasubramanian, P. Voshall, and W. Vogels [2007]. “Dynamo: Amazon’s highly available key-value store,” in *Proc. 21st ACM Symposium on Operating Systems Principles*, October 14–17, 2007, Stevenson, Wash.
- Doherty, W.J. and Thadhani, A.J., 1982. The economic value of rapid response time. IBM Report.
- Fan, X., W. Weber, and L. A. Barroso [2007]. “Power provisioning for a warehouse-sized computer,” in *Proc. 34th Annual Int’l. Symposium on Computer Architecture (ISCA)*, June 9–13, 2007, San Diego, Calif.
- A. Fikes, “Storage architecture and challenges,” in Google Faculty Summit, 2010.
- Ghemawat, S., H. Gobioff, and S.-T. Leung (2003). “The Google file system,” in *Proc. 19th ACM Symposium on Operating Systems Principles*, October 19–22, 2003, Lake George, N.Y.
- Greenberg, A., N. Jain, S. Kandula, C. Kim, P. Lahiri, D. Maltz, P. Patel, and S. Sengupta [2009]. “VL2: A scalable and flexible data center network,” in *Proc. SIGCOMM*, August 17–21, Barcelona, Spain.
- González, A. and Day, M. April 27, 2016, “Amazon, Microsoft invest billions as computing shifts to cloud,” *The Seattle Times*. <http://www.seattletimes.com/business/technology/amazon-microsoft-invest-billions-as-computing-shifts-to-cloud/>
- Hamilton, J. [2009]. “Data center networks are in my way,” Stanford Clean Slate CTO Summit, October 23, 2009, [http://mvdirona.com/jrh/TalksAndPapers/JamesHamilton\\_CleanSlateCTO2009.pdf](http://mvdirona.com/jrh/TalksAndPapers/JamesHamilton_CleanSlateCTO2009.pdf).
- Hamilton, J. [2010]. “Cloud computing economies of scale,” *Proc. AWS Workshop on Genomics & Cloud Computing*, June 8, 2010, Seattle, Wash. ([http://mvdirona.com/jrh/TalksAndPapers/JamesHamilton\\_GenomicsCloud20100608.pdf](http://mvdirona.com/jrh/TalksAndPapers/JamesHamilton_GenomicsCloud20100608.pdf)).
- Hamilton, J., 2014. AWS Innovation at Scale, AWS Re-invent conference. [https://www.youtube.com/watch?v=JIQETrFC\\_SQ](https://www.youtube.com/watch?v=JIQETrFC_SQ)
- Hamilton, J., May 2015. The Return to the Cloud, <http://perspectives.mvdirona.com/2015/05/the-return-to-the-cloud/>
- Hamilton, J., April 2017. How Many Data Centers Needed World-Wide, <http://perspectives.mvdirona.com/2017/04/how-many-data-centers-needed-worldwide/>
- Hölzle, U. [2010]. “Brawny cores still beat wimpy cores, most of the time,” *IEEE Micro*, July/August.
- Kanev, S., Darago, J.P., Hazelwood, K., Ranganathan, P., Moseley, T., Wei, G.Y. and Brooks, D., 2015, June. Profiling a warehouse-scale computer. *ACM/IEEE 42nd Annual International Symposium on Computer Architecture (ISCA)*.

- Lang, W., J. M. Patel, and S. Shankar [2010]. “Wimpy node clusters: What about non-wimpy workloads?” *Proc. Sixth Int’l. Workshop on Data Management on New Hardware*, June 7, 2010, Indianapolis, Ind.
- Lim, K., P. Ranganathan, J. Chang, C. Patel, T. Mudge, and S. Reinhardt [2008]. “Understanding and designing new system architectures for emerging warehouse-computing environments,” *Proc. 35th Annual Int’l. Symposium on Computer Architecture (ISCA)*, June 21–25, 2008, Beijing, China.
- Narayanan, D., E. Thereska, A. Donnelly, S. Elnikety, and A. Rowstron [2009]. “Migrating server storage to SSDs: Analysis of trade-offs,” *Proc. 4th ACM European Conf. on Computer Systems*, April 1–3, 2009, Nuremberg, Germany.
- Pfister, G. F. [1998]. *In Search of Clusters*, 2nd ed., Prentice Hall, Upper Saddle River, N.J.
- Pinheiro, E., W.-D. Weber, and L. A. Barroso [2007]. “Failure trends in a large disk drive population,” *Proc. 5th USENIX Conference on File and Storage Technologies (FAST ’07)*, February 13–16, 2007, San Jose, Calif.
- Ranganathan, P., P. Leech, D. Irwin, and J. Chase [2006]. “Ensemble-level power management for dense blade servers,” *Proc. 33rd Annual Int’l. Symposium on Computer Architecture (ISCA)*, June 17–21, 2006, Boston, Mass., 66–77.
- Reddi, V. J., B. C. Lee, T. Chilimbi, and K. Vaid [2010]. “Web search using mobile cores: Quantifying and mitigating the price of efficiency,” *Proc. 37th Annual Int’l. Symposium on Computer Architecture (ISCA)*, June 19–23, 2010, Saint-Malo, France.
- Schroeder, B., and G. A. Gibson [2007]. “Understanding failures in petascale computers,” *Journal of Physics: Conference Series* 78, 188–198.
- Schroeder, B., E. Pinheiro, and W.-D. Weber [2009]. “DRAM errors in the wild: A large-scale field study,” *Proc. Eleventh Int’l. Joint Conf. on Measurement and Modeling of Computer Systems (SIGMETRICS)*, June 15–19, 2009, Seattle, Wash.
- Schurman, E. and J. Brutlag [2009]. “The User and Business Impact of Server Delays,” *Proc. Velocity: Web Performance and Operations Conf.*, June 22–24, 2009, San Jose, Calif.
- Tezzaron Semiconductor. [2004]. “*Soft Errors in Electronic Memory—A White Paper*, Tezzaron Semiconductor, Naperville, Ill. ([http://www.tezzaron.com/about/papers/soft\\_errors\\_1\\_1\\_secure.pdf](http://www.tezzaron.com/about/papers/soft_errors_1_1_secure.pdf)).
- Vahdat, A., M. Al-Fares, N. Farrington, R. N. Mysore, G. Porter, and S. Radhakrishnan [2010]. “Scale-out networking in the data center,” *IEEE Micro* July/August 2010.

As architects experiment with DSAs, knowing architecture history may help. There are likely older architecture ideas that were unsuccessful for general-purpose computing that could nevertheless make eminent sense for domain-specific architectures. After all, they probably did some things well, and either they might match

your domain, or, conversely, your domain might omit features that were challenges for these architectures. For example, both the Illiac IV (Barnes et al., 1968) from the 1960s and the FPS 120a (Charlesworth, 1981) from the 1970s had two-dimensional arrays of processing elements, so they are proper ancestors to the TPU and Paintbox. Similarly, while VLIW architectures of the Multiflow (Rau and Fisher, 1993) and Itanium (Sharangpani and Arora, 2000) were not commercial successes for general-purpose computing, Paintbox does not have the erratic data cache misses, unpredictable branches, or large code footprint that were difficult for VLIW architectures.

Two survey articles document that custom neural network ASICs go back at least 25 years (Ienne et al., 1996; Asanović, 2002). For example, CNAPS chips contained a 64 SIMD array of 16-bit by 8-bit multipliers, and several CNAPS chips could be connected together with a sequencer (Hammerstrom, 1990). The Synapse-1 system was based on a custom systolic multiply-accumulate chip called the MA-16, which performed sixteen 16-bit multiplications at a time (Ramacher et al., 1991). The system concatenated MA-16 chips and had custom hardware to do activation functions.

Twenty-five SPERT-II workstations, accelerated by the T0 custom ASIC, were deployed starting in 1995 to do both NN training and inference for speech recognition (Asanović et al., 1998). The 40-MHz T0 added vector instructions to the MIPS instruction set architecture. The eight-lane vector unit could produce up to sixteen 32-bit arithmetic results per clock cycle based on 8-bit and 16-bit inputs, making it 25 times faster at inference and 20 times faster at training than a SPARC-20 workstation. They found that 16 bits were insufficient for training, so they used two 16-bit words instead, which doubled training time. To overcome that drawback, they introduced “bunches” (batches) of 32–1000 data sets to reduce time spent updating weights, which made it faster than training with one word but no batches.

We use the phrase *Image Processing Unit* for Paintbox to identify this emerging class of processor, but this is not the first use of the term. The earliest use we can find is 1999, when the Sony Playstation put the name on a chip that was basically an MPEG2 decoder (Sony/Toshiba, 1999). In 2006, Freescale used IPU to name part of the i.MX31 Applications Processor, which is closer to the more generic way we interpret it (Freescale as part of i.MX31 Applications Processor, 2006).

## References

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., 2016. Tensor-flow: large-scale machine learning on heterogeneous distributed systems. arXiv preprint arXiv:1603.04467.
- Adolf, R., Rama, S., Reagen, B., Wei, G.Y., Brooks, D., 2016. Fathom: reference workloads for modern deep learning methods. In: IEEE International Symposium on Workload Characterization (IISWC).
- Amodei, D., et al., 2015. Deep speech 2: end-to-end speech recognition in English and mandarin, arXiv:1512.02595.

- Asanović, K., 2002. Programmable neurocomputing. In: Arbib, M.A. (Ed.), *The Handbook of Brain Theory and Neural Networks*, second ed. MIT Press, Cambridge, MA. ISBN: 0-262-01197-2. <https://people.eecs.berkeley.edu/~krste/papers/neurocomputing.pdf>.
- Asanović, K., Beck, A., Johnson, J., Wawrzynek, J., Kingsbury, B., Morgan, N., 1998. Training neural networks with Spert-II. In: Sundararajan, N., Saratchandran, P. (Eds.), *Parallel Architectures for Artificial Networks: Paradigms and Implementations*. IEEE Computer Society Press. ISBN: 0-8186-8399-6. (Chapter 11) <https://people.eecs.berkeley.edu/~krste/papers/annbook.pdf>.
- Bachrach, J., Vo, H., Richards, B., Lee, Y., Waterman, A., Avižienis, R., Wawrzynek, J., Asanović, K., 2012. Chisel: constructing hardware in a Scala embedded language. In: *Proceedings of the 49th Annual Design Automation Conference*, pp. 1216–1225.
- Barnes, G.H., Brown, R.M., Kato, M., Kuck, D.J., Slotnick, D.L., Stokes, R., 1968. The ILLIAC IV computer. *IEEE Trans. Comput.* 100 (8), 746–757.
- Bhattacharya, S., Lane, N.D., 2016. Sparsification and separation of deep learning layers for constrained resource inference on wearables. In: *Proceedings of the 14th ACM Conference on Embedded Network Sensor Systems CD-ROM*, pp. 176–189.
- Brunhaver, J., 2014. PhD thesis. Stanford.
- Canis, A., Choi, J., Aldham, M., Zhang, V., Kammoona, A., Czajkowski, T., Brown, S.D., Anderson, J.H., 2013. LegUp: an open-source high-level synthesis tool for FPGA-based processor/accelerator systems. *ACM Trans. Embed. Comput. Syst.* 13 (2).
- Canny, J., et al., 2015. Machine learning at the limit. In: *IEEE International Conference on Big Data*.
- Caulfield, A.M., Chung, E.S., Putnam, A., Haselman, H.A.J.F.M., Humphrey, S.H.M., Daniel, P.K.J.Y.K., Ovtcharov, L.T.M.K., Lanka, M.P.L.W.S., Burger, D.C.D., 2016. A cloud-scale acceleration architecture. In: *MICRO Conference*.
- Charlesworth, A.E., 1981. An approach to scientific array processing: the architectural design of the AP-120B/FPS-164 family. *Computer* 9, 18–27.
- Clark, J., October 26, 2015. Google Turning Its Lucrative Web Search Over to AI Machines. Bloomberg Technology, [www.bloomberg.com](http://www.bloomberg.com).
- Dally, W.J., 2002. Computer architecture is all about interconnect. In: *Proceedings of the 8th International Symposium High Performance Computer Architecture*.
- Freescale as part of i.MX31 Applications Processor, 2006. [http://cache.freescale.com/files/32bit/doc/white\\_paper/IMX31MULTIWP.pdf](http://cache.freescale.com/files/32bit/doc/white_paper/IMX31MULTIWP.pdf).
- Galal, S., Shacham, O., Brunhaver II, J.S., Pu, J., Vassiliev, A., Horowitz, M., 2013. FPU generator for design space exploration. In: *21st IEEE Symposium on Computer Arithmetic (ARITH)*.
- Hameed, R., Qadeer, W., Wachs, M., Azizi, O., Solomatnikov, A., Lee, B.C., Richardson, S., Kozyrakis, C., Horowitz, M., 2010. Understanding sources of inefficiency in general-purpose chips. *ACM SIGARCH Comput. Architect. News* 38 (3), 37–47.



- Hammerstrom, D., 1990. A VLSI architecture for high-performance, low-cost, on-chip learning. In: IJCNN International Joint Conference on Neural Networks.
- He, K., Zhang, X., Ren, S., Sun, J., 2016. Identity mappings in deep residual networks. Also in arXiv preprint arXiv:1603.05027.
- Huang, M., Wu, D., Yu, C.H., Fang, Z., Interlandi, M., Condie, T., Cong, J., 2016. Programming and runtime support to blaze FPGA accelerator deployment at datacenter scale. In: Proceedings of the Seventh ACM Symposium on Cloud Computing. ACM, pp. 456–469.
- Iandola, F., 2016. Exploring the Design Space of Deep Convolutional Neural Networks at Large Scale (Ph.D. dissertation). UC Berkeley.
- Ienne, P., Cornu, T., Kuhn, G., 1996. Special-purpose digital hardware for neural networks: an architectural survey. *J. VLSI Signal Process. Syst. Signal Image Video Technol.* 13 (1).
- Jouppi, N., 2016. Google supercharges machine learning tasks with TPU custom chip. <https://cloudplatform.googleblog.com>.
- Jouppi, N., Young, C., Patil, N., Patterson, D., Agrawal, G., et al., 2017. Datacenter performance analysis of a matrix processing unit. In: 44th International Symposium on Computer Architecture.
- Karpathy, A., et al., 2014. Large-scale video classification with convolutional neural networks. *CVPR*.
- Krizhevsky, A., Sutskever, I., Hinton, G., 2012. Imagenet classification with deep convolutional neural networks. *Adv. Neural Inf. Process. Syst.*
- Kung, H.T., Leiserson, C.E., 1980. Algorithms for VLSI processor arrays. Introduction to VLSI systems.
- Lee, Y., Waterman, A., Cook, H., Zimmer, B., Keller, B., Puggelli, A., Kwak, J., Jevtic, R., Bailey, S., Blagojevic, M., Chiu, P.-F., Avizienis, R., Richards, B., Bachrach, J., Patterson, D., Alon, E., Nikolic, B., Asanovic, K., 2016. An agile approach to building RISC-V microprocessors. *IEEE Micro* 36 (2), 8–20.
- Lewis-Kraus, G., 2016. The Great A.I. Awakening. *New York Times Magazine*.
- Nielsen, M., 2016. Neural Networks and Deep Learning. <http://neuralnetworksanddeeplearning.com/>.
- Nvidia, 2016. Tesla GPU Accelerators For Servers. <http://www.nvidia.com/object/teslaservers.html>.
- Olofsson, A., 2011. Debunking the myth of the \$100 M ASIC. *EE Times*. [http://www.eetimes.com/author.asp?section\\_id=36&doc\\_id=1266014](http://www.eetimes.com/author.asp?section_id=36&doc_id=1266014).
- Ovtcharov, K., Ruwase, O., Kim, J.Y., Fowers, J., Strauss, K., Chung, E.S., 2015a. Accelerating deep convolutional neural networks using specialized hardware. Microsoft Research Whitepaper. <https://www.microsoft.com/en-us/research/publication/accelerating-deepconvolutional-neural-networks-using-specialized-hardware/>.
- Ovtcharov, K., Ruwase, O., Kim, J.Y., Fowers, J., Strauss, K., Chung, E.S., 2015b. Toward accelerating deep learning at scale using specialized hardware in the datacenter. In: 2015 IEEE Hot Chips 27 Symposium.

- Patterson, D., Nikolić, B., 7/25/2015, Agile Design for Hardware, Parts I, II, and III. EE Times, [http://www.eetimes.com/author.asp?doc\\_id=1327239](http://www.eetimes.com/author.asp?doc_id=1327239).
- Patterson, D.A., Ditzel, D.R., 1980. The case for the reduced instruction set computer. *ACM SIGARCH Comput. Architect. News* 8 (6), 25–33.
- Prabhakar, R., Koeplinger, D., Brown, K.J., Lee, H., De Sa, C., Kozyrakis, C., Olukotun, K., 2016. Generating configurable hardware from parallel patterns. In: *Proceedings of the Twenty-First International Conference on Architectural Support for Programming Languages and Operating Systems*. ACM, pp. 651–665.
- Putnam, A., Caulfield, A.M., Chung, E.S., Chiou, D., Constantinides, K., Demme, J., Esmaeilzadeh, H., Fowers, J., Gopal, G.P., Gray, J., Haselman, M., Hauck, S., Heil, S., Hormati, A., Kim, J.-Y., Lanka, S., Larus, J., Peterson, E., Pope, S., Smith, A., Thong, J., Xiao, P.Y., Burger, D., 2014. A reconfigurable fabric for accelerating large-scale datacenter services. In: *41st International Symposium on Computer Architecture*.
- Putnam, A., Caulfield, A.M., Chung, E.S., Chiou, D., Constantinides, K., Demme, J., Esmaeilzadeh, H., Fowers, J., Gopal, G.P., Gray, J., Haselman, M., Hauck, S., Heil, S., Hormati, A., Kim, J.-Y., Lanka, S., Larus, J., Peterson, E., Pope, S., Smith, A., Thong, J., Xiao, P.Y., Burger, D., 2015. A reconfigurable fabric for accelerating large-scale datacenter services. *IEEE Micro*. 35 (3).
- Putnam, A., Caulfield, A.M., Chung, E.S., Chiou, D., Constantinides, K., Demme, J., Esmaeilzadeh, H., Fowers, J., Gopal, G.P., Gray, J., Haselman, M., Hauck, S., Heil, S., Hormati, A., Kim, J.-Y., Lanka, S., Larus, J., Peterson, E., Pope, S., Smith, A., Thong, J., Xiao, P.Y., Burger, D., 2016. A reconfigurable fabric for accelerating large-scale datacenter services. *Commun. ACM*.
- Qadeer, W., Hameed, R., Shacham, O., Venkatesan, P., Kozyrakis, C., Horowitz, M.A., 2015. Convolution engine: balancing efficiency & flexibility in specialized computing. *Commun. ACM* 58 (4).
- Ragan-Kelley, J., Barnes, C., Adams, A., Paris, S., Durand, F., Amarasinghe, S., 2013. Halide: a language and compiler for optimizing parallelism, locality, and recomputation in image processing pipelines. *ACM SIGPLAN Not.* 48 (6), 519–530.
- Ramacher, U., Beichter, J., Raab, W., Anlauf, J., Bruels, N., Hachmann, A., Wesseling, M., 1991. Design of a 1st generation neurocomputer. *VLSI Design of Neural Networks*. Springer, USA.
- Rau, B.R., Fisher, J.A., 1993. Instruction-level parallelism. *J. Supercomput.* 235, Springer Science & Business Media.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., 2015. Imagenet large scale visual recognition challenge. *Int. J. Comput. Vis.* 115 (3).
- Sergio Guadarrama, 2015. BVLC googlenet. [https://github.com/BVLC/caffe/tree/master/models/bvlc\\_googlenet](https://github.com/BVLC/caffe/tree/master/models/bvlc_googlenet).
- Shao, Y.S., Brooks, D., 2015. Research infrastructures for hardware accelerators. *Synth. Lect. Comput. Architect.* 10 (4), 1–99.

- Sharangpani, H., Arora, K., 2000. Itanium processor microarchitecture. *IEEE Micro* 20 (5), 24–43.
- Silver, D., Huang, A., Maddison, C.J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., Dieleman, S., 2016. Mastering the game of Go with deep neural networks and tree search. *Nature* 529 (7587).
- Smith, J.E., 1982. Decoupled access/execute computer architectures. In: *Proceedings of the 11th International Symposium on Computer Architecture*.
- Sony/Toshiba, 1999. ‘Emotion Engine’ in PS2 (“IPU is basically an MPEG2 decoder...”). <http://www.cpu-collection.de/?l0=co&l1=Sony&l2=Emotion+Engine>, <http://arstechnica.com/gadgets/2000/02/ee/3/>.
- Steinberg, D., 2015. Full-Chip Simulations, Keys to Success. In: *Proceedings of the Synopsys Users Group (SNUG) Silicon Valley 2015*.
- Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., Rabinovich, A., 2015. Going deeper with convolutions. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- TensorFlow Tutorials, 2016. <https://www.tensorflow.org/versions/r0.12/tutorials/index.html>.
- Tung, L., 2016. Google Translate: ‘This landmark update is our biggest single leap in 10 years’, *ZDNet*. <http://www.zdnet.com/article/google-translate-this-landmarkupdate-is-our-biggest-single-leap-in-10years/>.
- Vanhoucke, V., Senior, A., Mao, M.Z., 2011. Improving the speed of neural networks on CPUs. <https://static.googleusercontent.com/media/research.google.com/en//pubs/archive/37631.pdf>.
- Wu, Y., Schuster, M., Chen, Z., Le, Q., Norouzi, M., Macherey, W., Krikun, M., Cao, Y., Gao, Q., Macherey, K., Klingner, J., Shah, A., Johnson, M., Liu, X., Kaiser, Ł., Gouws, S., Kato, Y., Kudo, T., Kazawa, H., Stevens, K., Kurian, G., Patil, N., Wang, W., Young, C., Smith, J., Riesa, J., Rudnick, A., Vinyals, O., Corrado, G., Hughes, M., Dean, J., 2016. Google’s Neural Machine Translation System: Bridging the Gap between Human and Machine Translation. <http://arxiv.org/abs/1609.08144>.

## M.10

## The History of Magnetic Storage, RAID, and I/O Buses (Appendix D)

*Mass storage is a term used there to imply a unit capacity in excess of one million alphanumeric characters ...*

**Hoagland [1963]**

The variety of storage I/O and issues leads to a varied history for the rest of the story. (Smotherman [1989] explored the history of I/O in more depth.) This section discusses magnetic storage, RAID, and I/O buses and controllers. Jain [1991] and Lazowska et al. [1984] are books for those interested in learning more about queuing theory.

## Magnetic Storage

Magnetic recording was invented to record sound, and by 1941 magnetic tape was able to compete with other storage devices. It was the success of the ENIAC in 1947 that led to the push to use tapes to record digital information. Reels of magnetic tapes dominated removable storage through the 1970s. In the 1980s, the IBM 3480 cartridge became the *de facto* standard, at least for mainframes. It can transfer at 3 MB/sec by reading 18 tracks in parallel. The capacity is just 200 MB for this 1/2-inch tape. The 9840 cartridge, used by StorageTek in the Powder-Horn, transfers at 10 MB/sec and stores 20,000 MB. This device records the tracks in a zigzag fashion rather than just longitudinally, so that the head reverses direction to follow the track. This technique is called *serpentine recording*. Another 1/2-inch tape is Digital Linear Tape; the DLT7000 stores 35,000 MB and transfers at 5 MB/sec. Its competitor is helical scan, which rotates the head to get the increased recording density. In 2001, the 8-mm helical-scan tapes contain 20,000 MB and transfer at about 3 MB/sec. Whatever their density and cost, the serial nature of tapes creates an appetite for storage devices with random access.

In 1953, Reynold B. Johnson of IBM picked a staff of 15 scientists with the goal of building a radically faster random access storage system than tape. The goal was to have the storage equivalent of 50,000 standard IBM punch cards and to fetch the data in a single second. Johnson's disk drive design was simple but untried: The magnetic read/write sensors would have to float a few thousandths of an inch above the continuously rotating disk. Twenty-four months later the team emerged with the functional prototype. It weighed 1 ton and occupied about 300 cubic feet of space. The RAMAC-350 (Random Access Method of Accounting Control) used 50 platters that were 24 inches in diameter, rotated at 1200 RPM, with a total capacity of 5 MB and an access time of 1 second.

Starting with the RAMAC, IBM maintained its leadership in the disk industry, with its storage headquarters in San Jose, California, where Johnson's team did its work. Many of the future leaders of competing disk manufacturers started their careers at IBM, and many disk companies are located near San Jose.

Although RAMAC contained the first disk, a major breakthrough in magnetic recording was found in later disks with air-bearing read/write heads, where the head would ride on a cushion of air created by the fast-moving disk surface. This cushion meant the head could both follow imperfections in the surface and yet be very close to the surface. Subsequent advances have come largely from improved quality of components and higher precision. In 2001, heads flew 2 to 3 microinches above the surface, whereas in the RAMAC drive they were 1000 microinches away.

Moving-head disks quickly became the dominant high-speed magnetic storage, although their high cost meant that magnetic tape continued to be used extensively until the 1970s. The next important development for hard disks was the removable hard disk drive developed by IBM in 1962; this made it possible to share the expensive drive electronics and helped disks overtake tapes as the preferred storage medium. The IBM 1311 disk in 1962 had an areal density of 50,000

bits per square inch and a cost of about \$800 per megabyte. IBM also invented the floppy disk drive in 1970, originally to hold microcode for the IBM 370 series. Floppy disks became popular with the PC about 10 years later.

The second major disk breakthrough was the so-called Winchester disk design in about 1973. Winchester disks benefited from two related properties. First, integrated circuits lowered the costs of not only CPUs but also of disk controllers and the electronics to control disk arms. Reductions in the cost of the disk electronics made it unnecessary to share the electronics and thus made nonremovable disks economical. Since the disk was fixed and could be in a sealed enclosure, both the environmental and control problems were greatly reduced. Sealing the system allowed the heads to fly closer to the surface, which in turn enabled increases in areal density. The first sealed disk that IBM shipped had two spindles, each with a 30 MB disk; the moniker “30-30” for the disk led to the name Winchester. (America’s most popular sporting rifle, the Winchester 94, was nicknamed the “30-30” after the caliber of its cartridge.) Winchester disks grew rapidly in popularity in the 1980s, completely replacing removable disks by the middle of that decade. Before this time, the cost of the electronics to control the disk meant that the media had to be removable.

As mentioned in Appendix D, as DRAMs started to close the areal density gap and appeared to be catching up with disk storage, internal meetings at IBM called into question the future of disk drives. Disk designers concluded that disks must improve at 60% per year to forestall the DRAM threat, in contrast to the historical 29% per year. The essential enabler was magnetoresistive heads, with giant magnetoresistive heads enabling the current densities. Because of this competition, the gap in time between when a density record is achieved in the lab and when a disk is shipped with that density has closed considerably.

The personal computer created a market for small form factor (SFF) disk drives, since the 14-inch disk drives used in mainframes were bigger than the PC. In 2006, the 3.5-inch drive was the market leader, although the smaller 2.5-inch drive required for laptop computers was significant in sales volume. It remains to be seen whether handheld devices such as iPods or video cameras, which require even smaller disks, will remain significant in sales volume. For example, 1.8-inch drives were developed in the early 1990s for palmtop computers, but that market chose Flash instead and 1.8-inch drives disappeared.

## RAID

The SFF hard disks for PCs in the 1980s led a group at Berkeley to propose redundant arrays of inexpensive disks (RAID). This group had worked on the reduced instruction set computer effort and so expected much faster CPUs to become available. They asked: What could be done with the small disks that accompanied their PCs? and What could be done in the area of I/O to keep up with much faster processors? They argued to replace one mainframe drive with 50 small drives to gain much greater performance from that many independent arms. The many small

drives even offered savings in power consumption and floor space. The downside of many disks was much lower mean time to failure (MTTF). Hence, on their own they reasoned out the advantages of redundant disks and rotating parity to address how to get greater performance with many small drives yet have reliability as high as that of a single mainframe disk.

The problem they experienced when explaining their ideas was that some researchers had heard of disk arrays with some form of redundancy, and they didn't understand the Berkeley proposal. Hence, the first RAID paper [Patterson, Gibson, and Katz 1987] is not only a case for arrays of SFF disk drives but also something of a tutorial and classification of existing work on disk arrays. Mirroring (RAID 1) had long been used in fault-tolerant computers such as those sold by Tandem. Thinking Machines had arrays with 32 data disks and 7 check disks using ECC for correction (RAID 2) in 1987, and Honeywell Bull had a RAID 2 product even earlier. Also, disk arrays with a single parity disk had been used in scientific computers in the same time frame (RAID 3). Their paper then described a single parity disk with support for sector accesses (RAID 4) and rotated parity (RAID 5). Chen et al. [1994] surveyed the original RAID ideas, commercial products, and more recent developments.

Unknown to the Berkeley group, engineers at IBM working on the AS/400 computer also came up with rotated parity to give greater reliability for a collection of large disks. IBM filed a patent on RAID 5 before the Berkeley group wrote their paper. Patents for RAID 1, RAID 2, and RAID 3 from several companies predate the IBM RAID 5 patent, which has led to plenty of courtroom action.

The Berkeley paper was written before the World Wide Web, but it captured the imagination of many engineers, as copies were faxed around the world. One engineer at what is now Seagate received seven copies of the paper from friends and customers. EMC had been a supplier of DRAM boards for IBM computers, but around 1988 new policies from IBM made it nearly impossible for EMC to continue to sell IBM memory boards. Apparently, the Berkeley paper also crossed the desks of EMC executives, and they decided to go after the market dominated by IBM disk storage products instead. As the paper advocated, their model was to use many small drives to compete with mainframe drives, and EMC announced a RAID product in 1990. It relied on mirroring (RAID 1) for reliability; RAID 5 products came much later for EMC. Over the next year, Micropolis offered a RAID 3 product, Compaq offered a RAID 4 product, and Data General, IBM, and NCR offered RAID 5 products.

The RAID ideas soon spread to the rest of the workstation and server industry. An article explaining RAID in *Byte* magazine (see Anderson [1990]) led to RAID products being offered on desktop PCs, which was something of a surprise to the Berkeley group. They had focused on performance with good availability, but higher availability was attractive to the PC market.

Another surprise was the cost of the disk arrays. With redundant power supplies and fans, the ability to “hot swap” a disk drive, the RAID hardware controller itself, the redundant disks, and so on, the first disk arrays cost many times the cost of the disks. Perhaps as a result, the “inexpensive” in RAID morphed into



“independent.” Many marketing departments and technical writers today know of RAID only as “redundant arrays of independent disks.”

The EMC transformation was successful; in 2006, EMC was the leading supplier of storage systems, and NetApp was the leading supplier of Network-Attached Storage systems. RAID was a \$30 billion industry in 2006, and more than 80% of the non-PC drive sales were found in RAIDs. In recognition of their role, in 1999 Garth Gibson, Randy Katz, and David Patterson received the IEEE Reynold B. Johnson Information Storage Award “for the development of Redundant Arrays of Inexpensive Disks (RAID).”

## I/O Buses and Controllers

The ubiquitous microprocessor inspired not only the personal computers of the 1970s but also the trend in the late 1980s and 1990s of moving controller functions into I/O devices. I/O devices have continued this trend by moving controllers into the devices themselves. These devices are called *intelligent devices*, and some bus standards (e.g., SCSI) have been created specifically for them. Intelligent devices can relax the timing constraints by handling many low-level tasks themselves and queuing the results. For example, many SCSI-compatible disk drives include a track buffer on the disk itself, supporting read ahead and connect/disconnect. Thus, on a SCSI string some disks can be seeking and others loading their track buffer while one is transferring data from its buffer over the SCSI bus. The controller in the original RAMAC, built from vacuum tubes, only needed to move the head over the desired track, wait for the data to pass under the head, and transfer data with calculated parity. SCSI, which stands for *small computer systems interface*, is an example of one company inventing a bus and generously encouraging other companies to build devices that would plug into it. Shugart created this bus, originally called SASI. It was later standardized by the IEEE.

There have been several candidates to be the successor to SCSI, with the current leading contender being Fibre Channel Arbitrated Loop (FC-AL). The SCSI committee continues to increase the clock rate of the bus, giving this standard a new life, and SCSI is lasting much longer than some of its proposed successors. With the creation of serial interfaces for SCSI (“Serial Attach SCSI”) and ATA (“Serial ATA”), they may have very long lives.

Perhaps the first multivendor bus was the PDP-11 Unibus in 1970 from DEC. Alas, this open-door policy on buses is in contrast to companies with proprietary buses using patented interfaces, thereby preventing competition from plug-compatible vendors. Making a bus proprietary also raises costs and lowers the number of available I/O devices that plug into it, since such devices must have an interface designed just for that bus. The PCI bus pushed by Intel represented a return to open, standard I/O buses inside computers. Its immediate successor is PCI-X, with Infiniband under development in 2000. Both were standardized by multicompany trade associations.

The machines of the RAMAC era gave us I/O interrupts as well as storage devices. The first machine to extend interrupts from detecting arithmetic abnormalities to detecting asynchronous I/O events is credited as the NBS DYSEAC in 1954 [Leiner and Alexander 1954]. The following year, the first machine with DMA was operational, the IBM SAGE. Just as today's DMA has, the SAGE had address counters that performed block transfers in parallel with CPU operations.

The early IBM 360s pioneered many of the ideas that we use in I/O systems today. The 360 was the first commercial machine to make heavy use of DMA, and it introduced the notion of I/O programs that could be interpreted by the device. Chaining of I/O programs was an important feature. The concept of channels introduced in the 360 corresponds to the I/O bus of today.

Myer and Sutherland [1968] wrote a classic paper on the trade-off of complexity and performance in I/O controllers. Borrowing the religious concept of the "wheel of reincarnation," they eventually noticed they were caught in a loop of continuously increasing the power of an I/O processor until it needed its own simpler coprocessor. Their quote in Appendix D captures their cautionary tale.

The IBM mainframe I/O channels, with their I/O processors, can be thought of as an inspiration for Infiniband, with their processors on their Host Channel Adaptor cards.

## References

- Anderson, D. [2003]. "You don't know jack about disks," *Queue* 1:4 (June), 20–30.
- Anderson, D., J. Dykes, and E. Riedel [2003]. "SCSI vs. ATA—more than an interface," *Proc. 2nd USENIX Conf. on File and Storage Technology (FAST '03)*, March 31–April 2, 2003, San Francisco.
- Anderson, M. H. [1990]. "Strength (and safety) in numbers (RAID, disk storage technology)," *Byte* 15:13 (December), 337–339.
- Anon. et al. [1985]. *A Measure of Transaction Processing Power*, Tandem Tech. Rep. TR 85.2. Also appeared in *Datamation*, 31:7 (April), 112–118.
- Bashe, C. J., W. Buchholz, G. V. Hawkins, J. L. Ingram, and N. Rochester [1981]. "The architecture of IBM's early computers," *IBM J. Research and Development* 25:5 (September), 363–375.
- Bashe, C. J., L. R. Johnson, J. H. Palmer, and E. W. Pugh [1986]. *IBM's Early Computers*, MIT Press, Cambridge, Mass.
- Blaum, M., J. Brady, J. Bruck, and J. Menon [1994]. "EVENODD: An optimal scheme for tolerating double disk failures in RAID architectures," *Proc. 21st Annual Int'l. Symposium on Computer Architecture (ISCA)*, April 18–21, 1994, Chicago, 245–254.
- Blaum, M., J. Brady, J. Bruck, and J. Menon [1995]. "EVENODD: An optimal scheme for tolerating double disk failures in RAID architectures," *IEEE Trans. on Computers* 44:2 (February), 192–202.

- Blaum, M., J. Brady, J., Bruck, J. Menon, and A. Vardy [2001]. "The EVENODD code and its generalization," in H. Jin, T. Cortes, and R. Buyya, eds., *High Performance Mass Storage and Parallel I/O: Technologies and Applications*, IEEE & Wiley Press, New York, 187–208.
- Blaum, M., J. Bruck, and A. Vardy [1996]. "MDS array codes with independent parity symbols," *IEEE Trans. on Information Theory*, IT-42 (March), 529–542.
- Brady, J. T. [1986]. "A theory of productivity in the creative process," *IEEE CG&A* (May), 25–34.
- Brown, A., and D. A. Patterson [2000]. "Towards maintainability, availability, and growth benchmarks: A case study of software RAID systems," *Proc. 2000 USE-NIX Annual Technical Conf.*, June 18–23, San Diego, Calif.
- Bucher, I. V., and A. H. Hayes [1980]. "I/O performance measurement on Cray-1 and CDC 7000 computers," *Proc. Computer Performance Evaluation Users Group, 16th Meeting*, October 20–23, 1980, Orlando, Fl., 245–254.
- Chen, P. M., G. A. Gibson, R. H. Katz, and D. A. Patterson [1990]. "An evaluation of redundant arrays of inexpensive disks using an Amdahl 5890," *Proc. ACM SIGMETRICS Conf. on Measurement and Modeling of Computer Systems*, May 22–25, 1990, Boulder, Colo.
- Chen, P. M., and E. K. Lee [1995]. "Striping in a RAID level 5 disk array," *Proc. ACM SIGMETRICS Conf. on Measurement and Modeling of Computer Systems*, May 15–19, 1995, Ottawa, Canada, 136–145.
- Chen, P. M., E. K. Lee, G. A. Gibson, R. H. Katz, and D. A. Patterson [1994]. "RAID: High-performance, reliable secondary storage," *ACM Computing Surveys* 26:2 (June), 145–188.
- Corbett, P., B. English, A. Goel, T. Grcanac, S. Kleiman, J. Leong, and S. Sankar [2004]. "Row-diagonal parity for double disk failure correction," *Proc. 3rd USENIX Conf. on File and Storage Technology (FAST '04)*, March 31–April 2, 2004, San Francisco.
- Denehy, T. E., J. Bent, F. I. Popovici, A. C. Arpaci-Dusseau, and R. H. Arpaci-Dusseau [2004]. "Deconstructing storage arrays," *Proc. 11th Int'l. Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, October 7–13, 2004, Boston, Mass., 59–71.
- Doherty, W. J., and R. P. Kelisky [1979]. "Managing VM/CMS systems for user effectiveness," *IBM Systems J.* 18:1, 143–166.
- Douceur, J. R., and W. J. Bolosky [1999]. "A large scale study of file-system contents," *Proc. ACM SIGMETRICS Conf. on Measurement and Modeling of Computer Systems*, May 1–9, 1999, Atlanta, Ga., 59–69.
- Enriquez, P. [2001]. "What happened to my dial tone? A study of FCC service disruption reports," poster, *Richard Tapia Symposium on the Celebration of Diversity in Computing*, October 18–20, 2001, Houston, Tex.
- Friesenborg, S. E., and R. J. Wicks [1985]. *DASD Expectations: The 3380, 3380-23, and MVS/XA*, Tech. Bulletin GG22-9363-02, IBM Washington Systems Center, Gaithersburg, Md.

- Gibson, G. A. [1992]. *Redundant Disk Arrays: Reliable, Parallel Secondary Storage*, ACM Distinguished Dissertation Series, MIT Press, Cambridge, Mass.
- Goldstein, S. [1987]. *Storage Performance—An Eight Year Outlook*, Tech. Rep. TR 03.308-1, IBM Santa Teresa Laboratory, San Jose, Calif.
- Gray, J. [1990]. “A census of Tandem system availability between 1985 and 1990,” *IEEE Trans. on Reliability*, 39:4 (October), 409–418.
- Gray, J. (ed.) [1993]. *The Benchmark Handbook for Database and Transaction Processing Systems*, 2nd ed., Morgan Kaufmann, San Francisco.
- Gray, J., and A. Reuter [1993]. *Transaction Processing: Concepts and Techniques*, Morgan Kaufmann, San Francisco.
- Gray, J., and D. P. Siewiorek [1991]. “High-availability computer systems.” *Computer* 24:9 (September), 39–48.
- Gray, J., and C. van Ingen [2005]. *Empirical Measurements of Disk Failure Rates and Error Rates*, MSR-TR-2005-166, Microsoft Research, Redmond, Wash.
- Gurumurthi, S., A. Sivasubramaniam, and V. Natarajan [2005]. Disk Drive Roadmap from the Thermal Perspective: A Case for Dynamic Thermal Management, *Proceedings of the International Symposium on Computer Architecture (ISCA)*, June, 38–49.
- Henly, M., and B. McNutt [1989]. *DASD I/O Characteristics: A Comparison of MVS to VM*, Tech. Rep. TR 02.1550, IBM General Products Division, San Jose, Calif.
- Hewlett-Packard. [1998]. “HP’s ‘5NINES:5MINUTES’ vision extends leadership and re-defines high availability in mission-critical environments,” February 10, [www.future.enterprisecomputing.hp.com/ia64/news/5nines\\_vision\\_pr.html](http://www.future.enterprisecomputing.hp.com/ia64/news/5nines_vision_pr.html).
- Hoagland, A. S. [1963]. *Digital Magnetic Recording*, Wiley, New York.
- Hospodor, A. D., and A. S. Hoagland [1993]. “The changing nature of disk controllers.” *Proc. IEEE* 81:4 (April), 586–594.
- IBM. [1982]. *The Economic Value of Rapid Response Time*, GE20-0752-0, IBM, White Plains, N.Y., 11–82.
- Imprimis. [1989]. *Imprimis Product Specification, 97209 Sabre Disk Drive IPI-2 Interface 1.2 GB*, Document No. 64402302, Imprimis, Dallas, Tex.
- Jain, R. [1991]. *The Art of Computer Systems Performance Analysis: Techniques for Experimental Design, Measurement, Simulation, and Modeling*, Wiley, New York.
- Katz, R. H., D. A. Patterson, and G. A. Gibson [1989]. “Disk system architectures for high performance computing,” *Proc. IEEE* 77:12 (December), 1842–1858.
- Kim, M. Y. [1986]. “Synchronized disk interleaving,” *IEEE Trans. on Computers* C-35:11 (November), 978–988.
- Kuhn, D. R. [1997]. “Sources of failure in the public switched telephone network,” *IEEE Computer* 30:4 (April), 31–36.
- Lambright, D. [2000]. “Experiences in measuring the reliability of a cache-based storage system,” *Proc. of First Workshop on Industrial Experiences with Systems Software (WIESS 2000), Co-Located with the 4th Symposium on Operating Systems Design and Implementation (OSDI)*, October 22, 2000, San Diego, Calif.

- Laprie, J.-C. [1985]. "Dependable computing and fault tolerance: Concepts and terminology," *Proc. 15th Annual Int'l. Symposium on Fault-Tolerant Computing*, June 19–21, 1985, Ann Arbor, Mich., 2–11.
- Lazowska, E. D., J. Zahorjan, G. S. Graham, and K. C. Sevcik [1984]. *Quantitative System Performance: Computer System Analysis Using Queueing Network Models*, Prentice Hall, Englewood Cliffs, N.J. (Although out of print, it is available online at [www.cs.washington.edu/homes/lazowska/qsp/](http://www.cs.washington.edu/homes/lazowska/qsp/).)
- Leiner, A. L. [1954]. "System specifications for the DYSEAC," *J. ACM* 1:2 (April), 57–81.
- Leiner, A. L., and S. N. Alexander [1954]. "System organization of the DYSEAC," *IRE Trans. of Electronic Computers* EC-3:1 (March), 1–10.
- Maberly, N. C. [1966]. *Mastering Speed Reading*, New American Library, New York.
- Major, J. B. [1989]. "Are queuing models within the grasp of the unwashed?" *Proc. Int'l. Conf. on Management and Performance Evaluation of Computer Systems*, December 11–15, 1989, Reno, Nev., 831–839.
- Mueller, M., L. C. Alves, W. Fischer, M. L. Fair, and I. Modi [1999]. "RAS strategy for IBM S/390 G5 and G6," *IBM J. Research and Development*, 43:5–6 (September–November), 875–888.
- Murphy, B., and T. Gent [1995]. "Measuring system and software reliability using an automated data collection process," *Quality and Reliability Engineering International*, 11:5 (September–October), 341–353.
- Myer, T. H., and I. E. Sutherland [1968]. "On the design of display processors," *Communications of the ACM*, 11:6 (June), 410–414.
- National Storage Industry Consortium. [1998]. "Tape Roadmap," [www.nsic.org](http://www.nsic.org).
- Nelson, V. P. [1990]. "Fault-tolerant computing: Fundamental concepts," *Computer* 23:7 (July), 19–25.
- Nyberg, C. R., T. Barclay, Z. Cvetanovic, J. Gray, and D. Lomet [1994]. "Alpha-Sort: A RISC machine sort," *Proc. ACM SIGMOD*, May 24–27, 1994, Minneapolis, Minn.
- Okada, S., S. Okada, Y. Matsuda, T. Yamada, and A. Kobayashi [1999]. "System on a chip for digital still camera," *IEEE Trans. on Consumer Electronics* 45:3 (August), 584–590.
- Patterson, D. A., G. A. Gibson, and R. H. Katz [1987]. *A Case for Redundant Arrays of Inexpensive Disks (RAID)*, Tech. Rep. UCB/CSD 87/391, University of California, Berkeley. Also appeared in *Proc. ACM SIGMOD*, June 1–3, 1988, Chicago, 109–116.
- Pavan, P., R. Bez, P. Olivo, and E. Zandoni [1997]. "Flash memory cells—an overview," *Proc. IEEE* 85:8 (August), 1248–1271.
- Robinson, B., and L. Blount [1986]. *The VM/HPO 3880-23 Performance Results*, IBM Tech. Bulletin GG66-0247-00, IBM Washington Systems Center, Gaithersburg, Md.
- Salem, K., and H. Garcia-Molina [1986]. "Disk striping," *Proc. 2nd Int'l. IEEE Conf. on Data Engineering*, February 5–7, 1986, Washington, D.C., 249–259.

- Scranton, R. A., D. A. Thompson, and D. W. Hunter [1983]. *The Access Time Myth*, Tech. Rep. RC 10197 (45223), IBM, Yorktown Heights, N.Y.
- Seagate. [2000]. *Seagate Cheetah 73 Family: ST173404LW/LWV/LC/LCV Product Manual*, Vol. 1, Seagate, Scotts Valley, Calif. ([www.seagate.com/support/disc/manuals/scsi/29478b.pdf](http://www.seagate.com/support/disc/manuals/scsi/29478b.pdf)).
- Smotherman, M. [1989]. "A sequencing-based taxonomy of I/O systems and review of historical machines," *Computer Architecture News* 17:5 (September), 5–15. Reprinted in *Computer Architecture Readings*, M. D. Hill, N. P. Jouppi, and G. S. Sohi, eds., Morgan Kaufmann, San Francisco, 1999, 451–461.
- Talagala, N. [2000]. "Characterizing Large Storage Systems: Error Behavior and Performance Benchmarks," Ph.D. dissertation, Computer Science Division, University of California, Berkeley.
- Talagala, N., and D. Patterson [1999]. *An Analysis of Error Behavior in a Large Storage System*, Tech. Report UCB//CSD-99-1042, Computer Science Division, University of California, Berkeley.
- Talagala, N., R. Arpaci-Dusseau, and D. Patterson [2000]. *Micro-Benchmark Based Extraction of Local and Global Disk Characteristics*, CSD-99-1063, Computer Science Division, University of California, Berkeley.
- Talagala, N., S. Asami, D. Patterson, R. Futernick, and D. Hart [2000]. "The art of massive storage: A case study of a Web image archive," *IEEE Computer* (November), 22–28.
- Thadhani, A. J. [1981]. "Interactive user productivity," *IBM Systems J.* 20:4, 407–423.