# A Split at the Core

Physics is forcing the microchip industry to redesign its most lucrative products. That is bad news for software companies

By W. Wayt Gibbs

It was never a question of whether, but only of when and why. When would microprocessor manufacturers be forced to slow the primary propulsive force for their industry—namely, the biennial release of new chips with much smaller transistors and much higher clock speeds that has made it so attractive for computer users to periodically trade up to a faster machine? And would it be fundamental physics or simple economics that raised the barrier to further scaling? The answers are: in 2004, and for both reasons.

Production difficulties bedeviled almost every major semiconductor firm this year, but none were more apparent than the travails of Intel, the flagship of the microchip business. The company delayed the release of "Prescott," a faster version of its Pentium 4 processor, by more than six months as it worked out

glitches in the fabrication of the 125-million-transistor chip. When Prescott did finally arrive, analysts were generally unimpressed by its performance, which was only marginally superior to the previous, 55-million-transistor Pentium 4. The company recalled defective batches of another microchip, postponed the introduction of new notebook processors, and pushed to next year a four-gigahertz Pentium model that it had promised to deliver this autumn.

The decision of greatest portent, however, was the one that Intel took in May to halt work on its next-generation Pentium 4 and Xeon processors. "They were probably a couple of years in design," estimates William M. Siu, manager of the company's desktop platforms group and the executive who proposed the cancellation. "It was obviously a significant decision," he says—not just because of the lost investment but because

it means that the Pentium microarchitecture, the central engine both of Intel's business and of about three quarters of the world's computers, has reached the end of its life earlier than planned.

Beginning next year, all new Intel microprocessor designs for desktop and server computers will have not one but two "cores," or computational engines, on the same chip. Some high-end machines already have two or more microprocessors working side by side, as separate chips on a circuit board. But integrating multiple processors into one "multicore" chip involves a much more dramatic design change.

"When you bring those processors onto a single chip and reduce their interaction time to fractions of a nanosecond, that changes the whole equation," observes Justin R. Rattner, who joined Intel in 1973 and now directs its systems technology lab. "This is a major inflec-

tion point" in computer architecture, he emphasizes. The shift to multicore processing has considerable ramifications for how computers are sold, how they are upgraded and—most significantly—how they are programmed.

## Dodging the Danger Zone

"WE ARE NOT THE FIRST to do multicore," acknowledges Bob Liang, head of Intel's architecture research lab. In 2001 IBM introduced a dual-core processor, the Power4. "But we will be the first to bring it to a mass market," Liang claims. To make good on that promise, Intel will have to beat AMD, which in August demonstrated a dual-core version of its fast-selling Opteron processor and promised to have the chips in volume production by mid-2005. Meanwhile Sun Microsystems is rushing to develop "Niagara," a new microprocessor for network servers that boasts eight identical cores.

"The basic idea is to run them slower [than the single cores in today's processors] and make them simpler but to use more of them," says Stephen S. Pawlowski, who runs Intel's microprocessor technology lab. "Slower" and "simpler" are words rarely heard in the mi-

crochip industry—they give marketers migraines—yet that strategy may offer the only practical course around serious technical and economic obstacles.
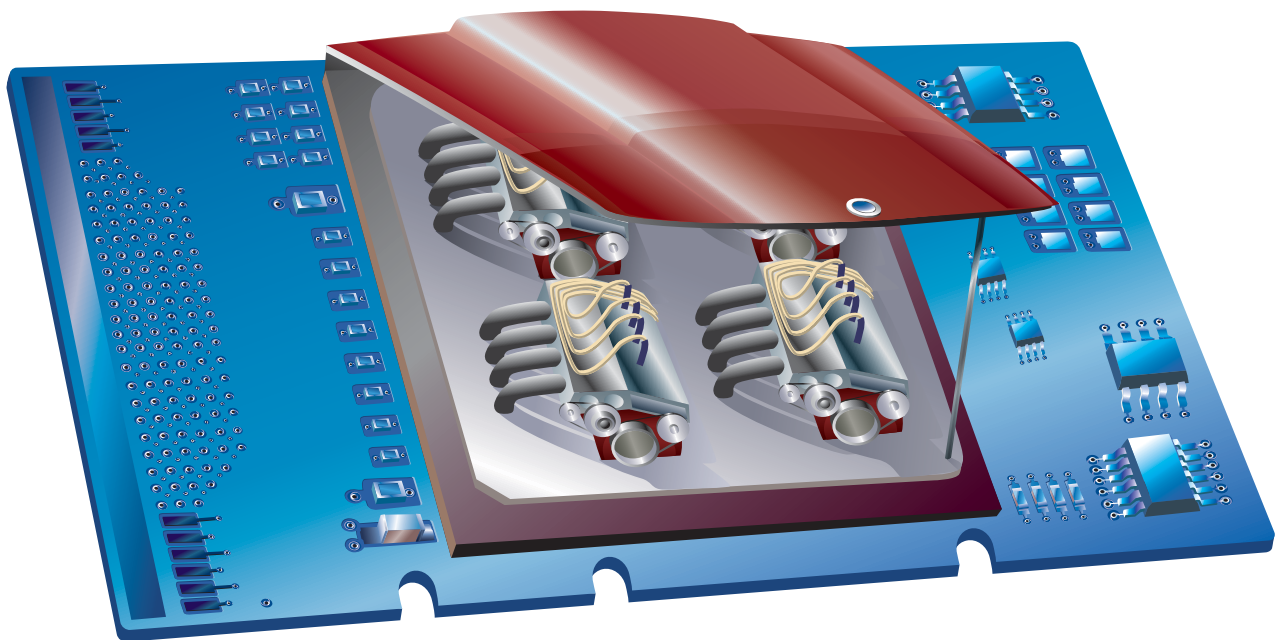
Engineers have been able to continue shrinking the smallest logic gates, albeit with difficulty and some delays. The next generation of processors, expected in 2006, will knock the length of logic gates down from 50 nanometers (billionths of a meter) to 35. "We're now making test chips with half a billion transistors on that process," reports Mark T. Bohr, Intel's director of process technologies. Bohr says the industry is on track to produce chips with 18-nanometer-long gates by the end of the decade. So the number of switches that fit on a chip—the so-called transistor budget—is rising as quickly as ever [see "The First Nanochips," by G. Dan Hutcheson; SCIENTIFIC AMERICAN, April].

Heat and power budgets are tightening rapidly, however. The peak energy consumption of a microprocessor has soared to well over 110 watts in recent years as chipmakers have cranked up the clock frequencies at which processors run [*see lower chart on next page*]. Most of that energy ends up as heat; a new Pentium 4 can generate more heat, per
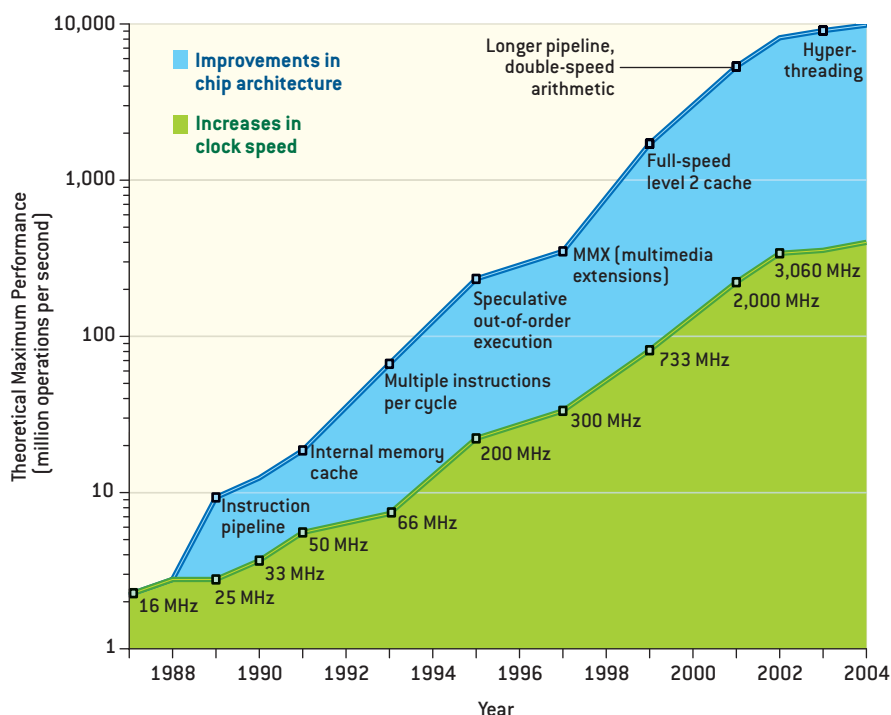
square centimeter, than an iron set on "cotton." As engineers scale down transistor sizes and pack them more densely onto the thumbnail-size processor die, operating temperatures will rise further unless clock speeds stabilize.

"One limit we face is the threshold voltage of transistors, which is determined by their ability to shut off current," Bohr explains. He compares a transistor to a valve. "We used to have to turn the wheel three or four times to get it fully open or fully closed. Now we're dealing with valves that turn off if you move them just a few degrees to the left or right." The sensitivity of the transistors makes them leaky. Even when turned off, each typically draws 100 nanoamperes of current, Bohr says. That may not sound like much. "But multiply that times 100 million transistors, and it adds up."
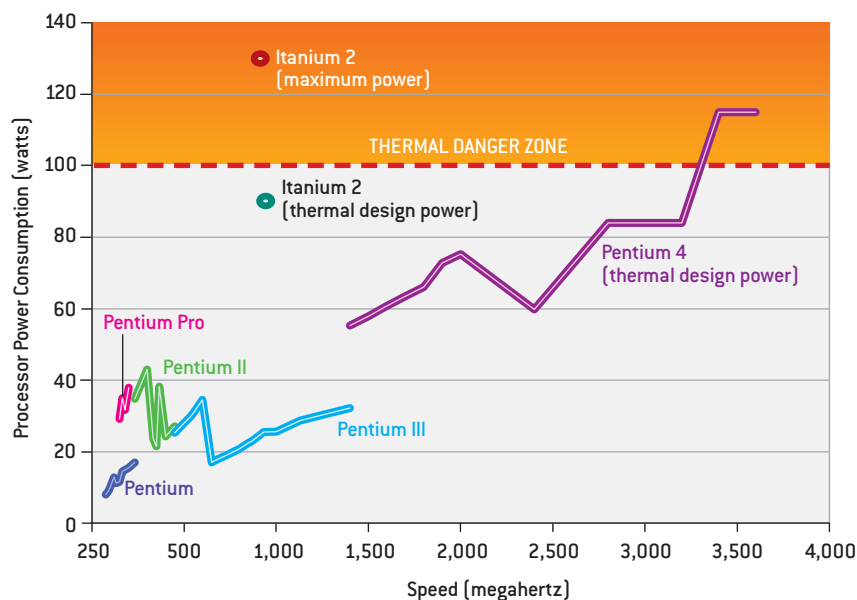
"When you push a certain thermal density, probably 100 watts per square centimeter, then you are really stressing the silicon," Pawlowski notes. Thanks to elaborate metal heat sinks and multiple fans, the hottest microprocessors still operate safely below that physical limit. "But the real issue is the cost of getting the heat out of the box, which is becom-



PEEK UNDER THE HOOD of a state-of-the-art microprocessor late next year, and you are likely to find two or more separate computing engines, running in parallel on a single slab of silicon. Such "multicore" designs alleviate old headaches for microchip engineers—but create new ones for software developers.

STEADILY INCREASING PERFORMANCE of Intel microprocessors has come mainly from quickening the clock pulse that sets the pace for the devices' transistors. But the company's engineers say that they will not be able to raise chip frequencies as quickly as in the past; further performance gains will have to come from innovations in the "architecture," or internal layout, of future processors. In fact, the next major generation of Intel processors will most likely use a radically different multicore architecture that may initially run at clock speeds below the current maximum frequency of 3,600 megahertz (MHz).



LATEST GENERATIONS OF INTEL PROCESSORS, the Pentium 4 and Itanium 2, are capable of burning more power than current heat sinks are designed to dissipate. To keep the chips from damaging themselves, Intel designers added circuits that monitor the processor temperature and throttle back the clock speed if a device runs too hot. Future generations of Intel processors will have multiple "cores" that run at lower speeds, generating less heat and spreading it more broadly across the surface of the chip.

ing prohibitive," he continues. "That is what creates this 'power wall.'"

Intel could require computer makers to switch from air-cooled to liquid-cooled machines. Apple Computer took that approach this year with its Power Mac G5 systems. But it adds to the cost. Intel sells roughly 50 times as many Pentium 4 systems as Apple sells G5s, in large part because Apple demands a premium price.

In any case, heat and power are not the only concerns. "When you have transistors on opposite corners of a chip, and you need to send a signal from one to the other, then those electrons have to flow through a copper wire," Bohr says. "The speed at which the electrons can flow is limited by the resistance and the capacitance of that wire. And while most wires in a chip are getting shorter, which is helpful, the wires are also getting thinner, which increases the delays caused by resistance and capacitance. So interconnects are inherently becoming more and more of a bottleneck."

That goes double for the relatively slow connection between the processor and the main memory bank. A microprocessor running at 3.6 gigahertz can execute several instructions each time its clock ticks, once every 277 trillionths of a second. But the system typically takes about 400 times that long to fetch information from the main memory. "The processor is just sitting there, waiting an eternity for each piece of data to come back from memory," Rattner observes.

Microprocessor architects have used on-chip memory caches and a technique called instruction-level parallelism to keep the processor busy working on instructions B and C while instruction A is waiting for its data to arrive. But that technique is nearly exhausted. "We're on the wrong side of a square law," Rattner says. "It is taking an exponential increase in transistors—and dramatic increases in the amount of power and chip area—to get even a modest increase in instruction-level parallelism."

Hence the dramatic change of strategy. Because the transistor budget is still rising by tens of millions of switches with each generation, engineers can ex-

## A NEW STRATEGY: E PLURIBUS UNUM

| COMPANY | DEVELOPMENT OF MULTICORE PROCESSORS |
|---------|-------------------------------------|
| AMD | Demonstrated its first working dual-core processor in August; the chip is scheduled to reach market next summer |
| Cisco Systems | Released a new network router in May that uses 192-core processors to handle 1.2 trillion bits per second of Internet traffic |
| IBM | Was first to sell a dual-core processor, the Power4, in 2001; Power5 processor introduced in May also sports two cores |
| Intel | Has prototype dual-core Itanium 2 chips; announced in May that new desktop and server microprocessor designs will have multiple cores |
| Sun Microsystems | UltraSparc-IV processor unveiled in February is dual-core; "Niagara" chip scheduled to appear by early 2006 will have eight cores |

ploit higher levels of parallelism by divvying the chip into multiple cores.

Intel's first dual-core chips will probably run at lower frequencies than the fastest Pentiums. But clock speeds will still continue to rise, asserts Patrick P. Gelsinger, Intel's chief technical officer, "just much more gradually than in the past." Intel recently relabeled its chips with abstract model numbers instead of the megahertz ratings it has used for 15 years. Gelsinger predicts that from now on, 70 percent of performance gains will come from architectural improvements—mainly parallelism—rather than from additional megahertz.

### Life in a Parallel Universe

IN PRINCIPLE, multicore processors could work more efficiently and more flexibly than today's single-core chips do. A notebook processor might have eight cores; a program customized for such a chip could divide itself into many "threads," each running simultaneously on a different core. Alternatively, the operating system might turn off some of the cores to extend battery life.

"The cores don't have to be identical," Siu points out. Building a variety of different cores could help deal with the fact that most existing software has no idea how to exploit a multicore processor. "You could have a big single-threaded core that can run legacy applications and also a bunch of small cores sitting on the side to run new [multicore-savvy] applications," Pawlowski elaborates.

But then he pauses to think about that prospect: "Quite frankly, it is going to take the software community a long time to start working on that. Unfortunately, very few people have played in this space."

"One of the big problems with parallelism for 40 years has been that it is hard to think about it and hard to do," says David J. Kuck, director of the KAI Software Lab, a company that Intel bought to help it make this transition. "When these threads are handing things back and forth, everyone gets lost."

A parallel processor deprives the programmer of one of the most valuable tools for debugging: repeatability. "A threaded [parallel] program is not a deterministic thing," Kuck explains. "It may execute one way one time and a different way the next time, just because of subtle timing differences in the machine's state. So most [software executives] think: 'Oh, my God, I just don't want to face this.' That holds all the way up and down the line from Oracle

### MORE TO EXPLORE

**International Technology Roadmap for Semiconductors.** International Sematech, 2003. Available at **http://public.itrs.net/Files/2003ITRS/Home2003.htm**

**Architecting the Era of Tera.** Intel Research and Development, 2004. Available at **ftp://download.intel.com/labs/nextnet/download/Tera_Era.pdf**

and Microsoft to the little guys," he says.

Even setting aside the difficulty of rewriting software in parallel form, "there are some applications where you won't get any boost from multicore. So it's just lost," acknowledges Glenn J. Hinton, director of microarchitecture development at Intel.

But certainly many kinds of tasks could run dramatically faster when redesigned for multicore chips. When converting a home movie to DVD format, for example, several frames can be processed in parallel. Rendering 3-D scenes, manipulating photographs, running scientific models, searching through databases and similar tasks can all be more quickly conquered once divided. A few specialized tasks could exploit as many cores as chipmakers can throw at them.

For general-purpose computing, however, "there is a point of diminishing returns," Pawlowski avers. "Sixteen cores is not clearly better than eight."

The most worrisome long-term question for the microprocessor industry may be whether the shift to multicore processors will discourage its customers from upgrading to newer computers. Today's computers are more than fast enough to handle most popular software. Demand for speedier machines has already begun to flag. In July, Intel reported that its inventory of unsold products had risen by 15 percent; the bad news knocked 11 percent off its share price.

A major design change adds uncertainty to apathy as a reason that computer owners might postpone an upgrade. It is not yet clear whether customers who buy the first dual-core machines and replace much of their software to suit the new architecture will have to repeat the process three years later to take advantage of "quad-core" machines. Faced with that prospect, many users—and for that matter, many software companies—could decide that the new architecture is simply not worth the hassle. On the other hand, the most obvious lesson from the history of computing is that every leap in performance is never enough for long.  SA

*W. Wayt Gibbs is senior writer.*