# Xeon Replaces Pentium Pro

## Intel Targets Servers and Workstations

*by Keith Diefendorff*

Intel has plugged the gaping hole at the top end of its product line—previously served by the aging Pentium Pro—with a Deschutes-based processor module the company labels Pentium II Xeon. As Figure 1 shows, the new processor family will serve the midrange to high-end server and workstation markets until the 64-bit Merced processor enters service in 2000.

Pentium Pro was previously the only processor in Intel's lineup capable of addressing this high-end segment, because it's the only processor that supports four-way multiprocessing (MP), memories larger than four gigabytes, and fast ECC L2 caches larger than 512K—all minimum requirements of the high-end market. But Pentium Pro, still in 0.35-micron technology, tops out at a paltry 200 MHz and supports only a 66-MHz system bus. These factors saddle it with anemic performance relative to its RISC competitors.

Xeon matches Pentium Pro's MP, memory, and L2 capabilities, plus it raises the system-bus speed to 100 MHz and dramatically boosts CPU speed to 400 MHz and beyond. Xeon is supported by two new chip sets (see MPR 7/13/98, p. 11) and by an Intel-built L2 SRAM that runs at the processor's full core frequency. Xeon and its companion SRAM(s) are packaged together on the new Slot 2 module (see MPR 3/30/98, p. 14), which uses the same bus protocol as Slot 1 but supports four-way SMP systems, has greater cooling capacity, and is physically larger (to accommodate more L2 cache).

The new processor is available now at 400 MHz. With 512K and 1M L2 caches, it is offered at essentially the same $1,100 and $2,800 price points as the Pentium Pro it replaces. Later this year, a 450-MHz version will become available, probably at similar price points. We expect the 2M cache version of this processor, however, to command a stratospheric price (relative to PC processors) close to $3,500!

### Xeon Challenges RISCs at the High End

Lately, Intel's overall average selling price (ASP) has taken a beating from declining PC prices and from competitive pressure applied on the low end by AMD, Cyrix, and IDT. The problem for Intel is that it needs a high ASP to fuel the semiconductor R&D and fab improvements that keep it ahead of its competitors.

Having so far failed to stimulate demand for higher performance (and higher priced) processors in PCs, Intel will try to take a larger share of the higher-margin workstation and server markets. While these markets are each about only 1% of the size of the PC market in unit volume, they can easily bear 10 times the processor price. This fact makes these markets immensely profitable and gives Intel an opportunity to increase revenue and ASP.

Beyond the desire to prop up revenue and ASP, Intel realizes that strategically it needs to own the markets on both sides of its core desktop-PC business to guard against competitors gaining a foothold and attacking from above or below. Having already made the mistake of allowing that to

| | 3Q98 | 4Q98 | 1H99 | 2H99 |
|---|---|---|---|---|
| **Xeon** Midrange to High-end Workstations & Servers | 400 MHz 512K, 1M Deschutes | 450 MHz 512K–2M Deschutes | 500 MHz 512K–2M Tanner | 700 MHz 1M, 2M 0.18μ Tanner |
| **Pentium II** Performance Desktop | 400 MHz 512K Deschutes | 450 MHz 512K Deschutes | 500 MHz 512K Katmai | 600 MHz 512K 0.18μ Katmai |
| **Celeron** Basic PC Desktop | 300 MHz Covington | 333 MHz 128K* Mendocino | 333 MHz 128K* Mendocino | 400 MHz 128K* Mendocino |
| **Mobile Pentium II** Mobile PC | 266 MHz 512K Deschutes | 300 MHz 512K Deschutes | 333 MHz 256K* Dixon | 366 MHz 256K* 0.18μ Katmai |

☐ Available  ▨ Source: Intel  ▨ Source: MDR

**Figure 1.** Intel divides the market into four segments, with Xeon taking the high ground. Shown are the fastest processors Intel will field in each segment over the next year. (*indicates on-chip L2)

*Inside:* **Xeon Chip Sets** ◇ **IBM C54xDSP** ◇ **Savage3D** ◇ **MGA-G200** ◇ **Banshee**

happen on the low end, Intel cannot afford to repeat the error on the high end. Actually, it may already have done so by failing to field a viable RISC processor. But, thankfully, its RISC competitors failed to capitalize, and Intel now has a second opportunity to avoid that mistake.

Intel also recognizes that technology flows downhill. Now that Intel has usurped most of the good ideas from the mainframes, it must create its own technology. Bleeding-edge technology is best developed in the context of a high-end business, because the loftier prices make it easier to fund, and the lower volumes make it less challenging to deploy.

Intel's long-range plan is to capture the high-end markets with its new IA-64 architecture. But the first implementation, Merced, has now been delayed to mid-2000 (see MPR 6/22/98, p. 1). Neither Pentium Pro nor Slot 1 Pentium IIs are strong enough to capture a significant share of these markets from the RISCs. Xeon represents Intel's interim attack on these markets.

## Intel Segments the Market by Branding

Significant obstacles block Intel's attack on both the low- and high-end markets. On the low end, Intel's competitors have successfully built competitive processors on Socket 7. To avoid further legitimizing Socket 7, Intel has been forced to push the market to Slot 1 and abandon Pentium.

Without Pentium, Intel is left with one processor, Deschutes, to serve the entire market spectrum. This is not just a technical challenge, it's a business challenge: how is it possible to charge $100 on the low end yet justify $3,500 on the high end—for exactly the same microprocessor? To charge less on the high end would just be leaving money on the table—not something Intel likes to do.

The company has two solutions: repackaging and branding. First, Intel will repackage Deschutes to target specific market segments. For example, Celerons today are packaged with no L2 and limited to 300 MHz for the low end, whereas Xeons have large, fast L2s and run at 400 MHz for the high end. Although the manufacturing-cost or performance differences between them do not justify a huge price delta, the fact that

neither is suitable for the other's market allows Intel to ask completely different prices for them.

Second, Intel is attempting to stratify the market segments in the minds of customers with a strong branding campaign. Branding, Intel hopes, will obscure the underlying microprocessor, thus freeing its products to follow the pricing models appropriate for their respective segments. A secondary benefit of branding is that it gives Intel cover for migrating from one generation of processor to the next. Customers, other than geeks who pry off lids or read this newsletter, may never know what microprocessor is in a future Celeron or Xeon.

Considering the competitive workstation and server landscape, Intel seems justified in demanding high prices for Xeon. Competitor Sun, for example, charges $3,961 for its 360-MHz Ultra-2 module (see MPR 5/11/98, p. 5). Also, the market appears capable of supporting such prices. After all, if a $1,500 PC can use a $300 processor, then surely a $15,000 workstation or server can support a $3,000 processor. Viewed in this way, the Xeon prices do not seem excessive.

## Pentium II Patched for Server Duty

Within the workstation and server markets, Intel differentiates between midrange-to-high-end systems and "volume" systems. The volume systems, generally those below about $5,000, will continue to be served by Pentium II. Xeon is aimed at systems above $5,000. In fact, Intel has taken specific marketing and product actions to prevent crossover.

To encourage companies to use (pay for) Xeon instead of Pentium II in the high-end markets, Intel left a few key features out of Pentium II, as Table 1 shows. Chief among them is its inability to go above two-way SMP.

While Deschutes implements the full cache-coherence protocol for *n*-way SMP operation, the 242-pin Slot 1 interface has signal-integrity limitations that prevent Intel from certifying its operation above two-way. High-end servers, however, require four-way SMP capability. The additional 58 power and ground pins on Xeon's 330-pin Slot 2 interface, and its use of AGTL+ (assisted Gunning transistor logic) signaling, provide the additional signal integrity necessary for 100-MHz four-way SMP systems.

AGTL+ is similar to Pentium II's GTL+, but AGTL+ signals are driven to $V_{CC}$ for one clock cycle after a low-to-high transition to improve rise times and reduce noise. As with GTL+, the input receivers are differential for good noise immunity, but with Xeon the $V_{REF}$ generator was moved onto the processor. Since AGTL+ is electrically compatible with GTL+, it is possible that all Pentium IIs could go to AGTL+. If not, the option is probably electrically or fuse programmable, so the same processor die could be used for either.

Another deficiency of Pentium II for large servers is its 4G memory address limit. Database servers can make good use of memories that are larger than 4G. Transaction processing performance can improve by 10% or more by moving from 4G to 8G of memory. In addition, extra physical address

| Feature | Pentium Pro | Pentium II | Xeon |
|---|---|---|---|
| CPU Frequency | 200 MHz | 400 MHz | 400 MHz |
| System Bus | 66 MHz | 66 or 100 MHz | 100 MHz |
| Interface | Socket 8 | Slot 1 | Slot 2 |
| Package | 387 PGA | 242 SECC | 330 SECC |
| Max L2 Size | 1M | 512K | 1M (2M 4Q98) |
| L2 Speed | 200 MHz | 200 MHz | 400 MHz |
| Chip Set | 450GX | 440BX | 440GX, 450NX |
| Max Memory | 64G | 4G | 64G |
| SMP | 4-way | 2-way | 4-way |
| Process | 0.35 μm | 0.25 μm | 0.25 μm |
| Voltage | 3.3 V | 2.0 V | 2.0 V |
| Power (max) | 44.0 W (1M) | 27.9 W (512K) | 38.1 W (1M) |

**Table 1.** Xeon combines the best attributes of Pentium Pro and Pentium II for the high-end workstation and server markets. SECC stands for single-edge contact cartridge. (Source: Intel)

bits are convenient for partitioning NUMA (nonuniform memory access) MP clusters. The Deschutes CPU implements a 64G physical address space, and the Slot 1 connector brings out the necessary address bits, but Pentium II's L2 cache tag chip limits cachable memory to 4G.

Pentium Pro provides a 36-bit (64G) physical-address space extension called PAE36. Microsoft, however, refused to support this extension, so the feature was rarely used. Deschutes introduced a new 36-bit address extension that Microsoft will support in Windows NT 5.0 and in NT 4.0 with an Intel-supplied driver. The new extension, originally called PSE36, has recently been relabeled Intel extended server memory (IESM). Deschutes also implements PAE36 for compatibility.

Like PAE36, IESM is enabled through a global mode bit. Unlike PAE36, which required the page-directory and page-table formats to be widened to eight bytes, IESM uses reserved bits in the normal four-byte page-table entries. This makes support for IESM more palatable to Microsoft.

The drawback of IESM is that memory above 4G can be accessed using 4M pages only. Such large pages can have deleterious effects on demand paging and can cause fragmentation of physical memory and the virtual address space.

The solution used by NT is to map the first 4G of memory with 4K pages as usual and utilize memory beyond 4G only as a software-managed RAMdisk. Extended memory is never mapped into the application's virtual address space. Applications wanting to take advantage of memory above 4G must be rewritten to use a new API that provides a protocol for explicitly managing buffers in extended memory.

## Xeon Goes for Bandwidth

Pentium Pro, besides being old (0.35-micron BiCMOS) and decrepit (200 MHz), is limited to a 66-MHz system bus. This gives it a puny 528 Mbytes/s of bandwidth to share among processors and I/O. The full-speed backside L2 cache, however, gives it an additional 1.6 Gbytes/s of bandwidth. As Table 2 shows, this provides a total aggregate bandwidth of 2.1 Gbytes/s in a uniprocessor system. This is plenty, considering the 200-MHz processor's meager appetite for data.

Pentium II, which can consume data at a much higher rate, has a far worse problem. Because its L2 cache bus runs at only half the CPU speed, a 300/66-MHz Pentium II has 19% less aggregate bandwidth available to feed the 50%-faster processor. The situation is worse for a 400/100-MHz
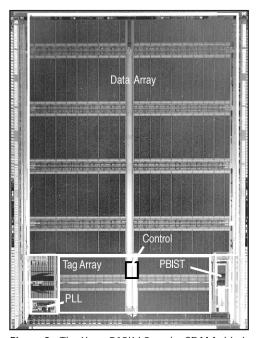


**Figure 2**. The Xeon 512K L2-cache SRAM chip is four-way set-associative and operates at up to 450 MHz. The 0.35-micron four-layer-metal part measures 12.9 × 17.2 mm (222 mm$^2$).

Pentium II, which is twice as fast as Pentium Pro but has essentially the same aggregate memory bandwidth. With this bandwidth limitation, Pentium II's performance just cannot scale above two processors.

Pentium Pro's L2 cache operates at the full processor speed, thanks to its custom SRAM. Intel switched to commodity SRAMs for Pentium II because they cost less, and because putting the fab capacity in place to build SRAMs in Pentium II's volumes would have been prohibitive, even for Intel. Since Xeon will serve a much smaller market, however, Intel can return to the custom SRAM approach for more bandwidth.

## SRAM Delivers 3.6 Gbytes/s

Pentium II's L2 cache comprises two or four commodity BSRAMs and a custom tag chip. This cache operates at half the CPU frequency with a latency of 10 processor clocks (25 ns) and a peak bandwidth of 1.6 Mbytes/s (at 400 MHz). Intel limits its size to 512K, presumably to prevent overlap with Xeon.

The Xeon cache chip, called C6C, puts all the data, tags, and tag-match logic on a single chip, as Figure 2 shows. C6C provides 512K of cache organized as four-way set-associative. The SRAM operates at up to 450 MHz, with a latency of five clocks (11 ns) and a peak bandwidth of 3.6 Gbytes/s.

The SRAM implements 32K tags, each consisting of 19 address bits (A35 to A17) plus 4 bits to track the MESI-coherence state of each line (only 2 bits are really needed, but it is stored redundantly so soft errors can be detected).

The C6C is built in Intel's standard 0.35-micron P854 CMOS process and uses an 18.5-μm$^2$ six-transistor cell. The SRAM die measures 12.9 × 17.2 mm (222 mm$^2$), which is 70% larger than the 131-mm$^2$ Deschutes die. According to the MDR Cost Model, the manufacturing cost of the SRAM is about $90. The SRAM is provided in a 496-ball BGA.

The SRAM's I/O voltage is 2.0 V. Maximum power dissipation of the 2.5-V core, which occurs while it is perform-

| Processor | Frequency | | Sys Bus | Cache Bus | Aggregate B/W | | |
|---|---|---|---|---|---|---|---|
| | Core | Bus | | | Uni | Dual | Quad |
| Pentium II | 300 | 66 | 0.5 | 1.2 | 1.7 | 2.9 | n/a |
| | 400 | 100 | 0.8 | 1.6 | 2.4 | 4.0 | n/a |
| Pentium Pro | 200 | 66 | 0.5 | 1.6 | 2.1 | 3.7 | 6.9 |
| Xeon | 400 | 100 | 0.8 | 3.2 | 4.0 | 7.2 | 13.6 |

**Table 2**. Xeon's theoretical maximum aggregate memory bandwidth (in Gbytes/s) is nearly twice that of any other Intel processor, thanks to its full-speed backside cache. Systems with more than four processors require a cluster bridge that adds an extra load on the bus and limits its speed to 90 MHz (720 Mbytes/s).

ing continuous back-to-back reads, is 4.5 watts at 450 MHz. Maximum power dissipation for the complete device, including I/O drivers, is 7.3 watts.

The Deschutes CPU can support up to four SRAMs. To make space for these SRAMs and to dissipate their extra heat load, Xeon's Slot 2 module was increased in size to nearly twice that of Pentium II's Slot 1 module, as Figure 3 shows.
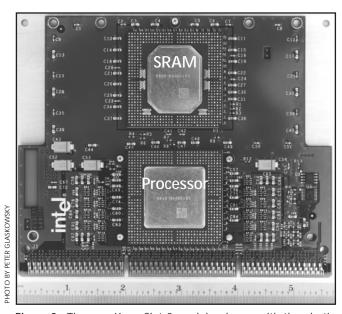


**Figure 3.** The new Xeon Slot 2 module, shown with the plastic cover and heat sink removed, measures 5.5" × 4.6". The module contains the Deschutes processor and enough room for four cache chips (two on each side). Only one SRAM is shown in this photo.

## Cache Expands to 2M

To support various cache sizes, the SRAM operates in one of three modes. In "512K mode," the SRAM implements 4,096 sets with a 32-byte line size and it returns four sequential cycles of 72 bits each (eight data bytes and one ECC byte).

For one-megabyte caches, two SRAMs are paired, with each driving half the data bus; both SRAMs are placed in "1M mode." In this mode, each SRAM implements 8,192 sets with a 16-byte line size, and each delivers to the processor four sequential cycles of 36 bits each. This scheme requires that each SRAM implement enough tags for the full 1M of cache, even though half are unused in 512K mode.

The original design called for cascading four SRAMs to provide a 2M cache. As a result, the cache chips have a "2M mode," in which pairs of SRAMs are cascaded and the next-most-significant address bit (A18) is used as a bank-select signal. In this mode, the chips ignore A18 in the tag-match, making the cache appear to have 16K sets.

Although the design calls for cascading these parts, the extra address- and data-bus loading could pose a problem at high-speed. Built-in self-tests indicate that C6C's internal array will operate at up to 750 MHz, and Intel seems confident it will function alone in systems at up to 500 MHz. But the additional electrical loads of a 2M cache could make 450 MHz operation a challenge. Deep sorting and high prices, however, should enable Intel to deliver adequate quantities.

Heat removal could also be a problem for a 2M Xeon. By our calculations, at 450 MHz a four-SRAM Xeon would dissipate a blistering 59 watts! Not a cooling job for the thermodynamically challenged.

These problems may become moot, however, since our sources indicate that Intel may have its one-megabyte SRAM (called CK1) ready in time for the 450-MHz/2M Xeon later this year. We suspect CK1 is being implemented in Intel's 0.25-micron P856 CMOS process, the same process used for the Deschutes processor. The shrink alone should improve speed by at least 25%, pushing it above 600 MHz without additional cycles of latency. Other double-data-rate (DDR) SRAM vendors are also talking about such speeds, so we do not see SRAM speeds becoming a limiting factor for Xeon. In fact, we expect that 0.18-micron technology could support full-speed Merced caches up to a gigahertz.

## Source-Synchronous Design Enables 400 MHz

The Xeon SRAM uses a source-synchronous DDR clocking scheme, as Figure 4 shows. At the beginning of a cache access, the processor delivers an address precisely centered on the falling edge of a half-speed (200-MHz) clock. The SRAM uses this clock edge to latch the address.

A phase-locked loop (PLL) synchronizes the SRAM's internal clocks to the half-speed external clock and doubles its frequency to 400 MHz for use internally. To save power and minimize access time, the operation of the SRAM array is self-timed (see MPR 10/6/97, p. 18). This technique uses fewer latches than a conventional pipelined array, with a

consequent reduction in clock loading and clock current. Since setup and hold requirements are eliminated at pipeline stages, overall access time is also reduced.

Once latched, the address is used to index the tag and data arrays in parallel. The tag arrays are faster than the data arrays, allowing way-selection to complete before the data is ready. Once the data is available, the correct way is multiplexed onto the output latch. If the access misses the cache entirely, a miss indication is returned to the CPU.

The data is then synchronized to the internal 400-MHz clock and transmitted back to the processor along with a half-speed version of this clock, called the data strobe. The strobe edges are precisely centered on the data periods. The processor uses both edges of the strobe to clock data in. This method minimizes clock skew and maximizes timing margins, allowing the interface to operate at full CPU speed.

An output-driver impedance-control mechanism eliminates noise from electrical reflections on the interface. The driver strength is programmed with four external resistors, one each for the address and data pull-up and pull-down devices. A digital feedback control loop matches the driver impedance to the appropriate resistor. Driver strengths are controlled to 5-bit resolution over a range of 25 to 70 ohms.

As Figure 4 shows, data accesses and data transfers are fully overlapped. The processor can deliver a new address while data is being transferred from the previous access, thus sustaining a continuous 3.6-Gbytes/s read transfer rate at 450 MHz. Surprisingly, the SRAM does not take advantage of late-write timing, so dead cycles must be inserted between reads and writes; two dead cycles are inserted between a read and a write, and six cycles between a write and a read.

Data is sequenced from the SRAM to the processor beginning with the quadword containing the missed datum. This critical-quadword-first delivery technique allows the processor to resume processing at the earliest possible moment, thus minimizing stall time. The burst order is 0123, 1032, 2301, or 3210, depending on the missed quadword.

## Thermal Sensor Enhances Reliability

Server customers demand system-management features that extend all the way into the CPU. For example, in a large MP system it is sometimes necessary to interrogate a CPU when it is not running. To answer this need, Xeon provides a separate system management bus, called SMBus strangely enough, through the Slot 2 connector. This bus gives OEMs access to information in the processor and in a processor-information ROM on the module. This serial flash ROM provides processor stepping and other information that is burned in at the factory, as well as a writable area for general use by the OEM. The ROM is accessed via the I²C protocol.

High reliability is mandatory for any server worth its salt. The biggest source of reliability problems for a microprocessor, other than software from the Pacific Northwest, is elevated silicon-junction temperature. Operation above the manufacturer's temperature specification, even momentarily,
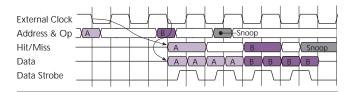


**Figure 4**. The Xeon SRAM uses a source-synchronous double-data-rate clocking scheme. The input clock runs at half the CPU frequency, and data is returned on both edges of a 200-MHz data strobe. The latency of the SRAM is five clocks.

can cause the electrical timing to drift, inducing transient failures. In addition, extended operation of the junctions much above 110°C shortens the silicon's expected lifetime.

Keeping junction temperatures low can be challenging with a hot CPU like Xeon. As a way to monitor the situation, Deschutes provides an on-die thermal diode whose forward-biased voltage drop is proportional to the processor's internal junction temperatures. The voltage is measured with an on-chip analog-to-digital converter and stored in a register accessible via the SMBus.

The diode's voltage-temperature relationship is sensitive to minor process variations, so it must be calibrated for each processor. The calibration data is determined during the manufacturing test flow and is stored in the processor-information ROM, where it can be accessed from the SMBus. Thermal trip points can be programmed that generate an SMB alert when they are crossed.

## Xeon Trounces RISCs on Server Benchmarks

Intel claims that Xeon will be the highest-performing processor—including all the RISCs—on TPC-C, as Figure 5 shows. Results so far seem to bear this out. Admittedly, the Xeon results are Intel's estimates, but we expect that, if anything, Intel is being conservative. This picture may change when Alpha 21264 systems emerge this fall, but the Xeon results are nonetheless quite impressive.

Results on SPEC benchmarks are less remarkable. While Xeon handily outperforms Pentium II on both SPECint95 and SPECfp95, and holds its own against all but the
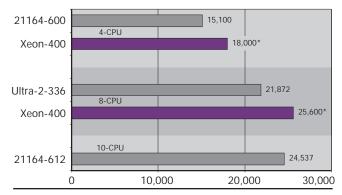


**Figure 5**. Intel's estimates of Xeon's TPC-C performance place it above all similar RISC-based systems. The eight-processor Xeon even beats out the 10-processor 612-MHz 21164 Alpha system. (Source: *www.tpc.org* except *Intel)

Alpha 21164 on SPECint95, it still trails all the leading RISCs on SPECfp95 by at least 65%, as Figure 6 shows.

Xeon's SPECint95 score is only 7% better than Pentium II's (with L2 ECC on). Intel is quick to point out, however, that scores on a few of the individual benchmarks are much higher, indicating that Xeon has a more robust performance profile than Pentium II; i.e., it performs better under a heavy load. Individual scores range from 0% to 14% higher. SPECfp95 results show Xeon with a somewhat larger 10% advantage over Pentium II. Individual SPECfp95 benchmarks range from 5% to 19% better. In general, with Pentium II's L2 ECC turned off, these differences in SPEC scores drop by four to five percentage points.

Intel's workstation performance brief—wonder of wonders—shows Xeon outpacing all the RISCs on several application level benchmarks. For example, it shows Xeon in a $7,500 system edging out a $20,000 767-MHz KryoTech Alpha system (see MPR 7/13/98, p. 4) while running Pro/Engineer's Bench98.

Intel defends these application-level benchmark results as more representative than SPECfp95 by pointing out that most workstation applications, even FP-based ones like Bench98, have fewer than 10% floating-point instructions by
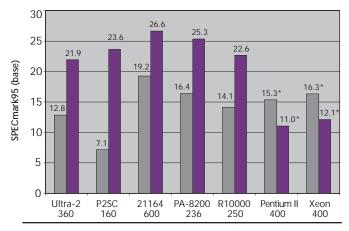
dynamic count. While this may be true, there are other important workstation applications, like Spice2G6, and important FP libraries, like LAPACK, that are heavily dominated by floating-point and will benefit significantly from a faster floating-point engine.

### Xeon Breaks Away From PC Pricing Model

We estimate the total manufacturing cost of a 1M Xeon with a Deschutes CPU and two custom SRAMs to be about $240. A price of $2,800 provides a manufacturing profit margin of over 90%—not bad compared with the 76% overall margin that we compute Intel ran last year.

Clearly, from a cost perspective, Xeon's $2,800 price is indefensible. But Xeon is not priced to cost, it is priced to what the market will bear. There is some evidence that such prices are viable in Xeon's target markets: the price for the 1M Pentium Pro—a dinosaur if ever there was one—has held steady for several quarters at $2,675.

The most amazing thing about the price structure of these markets is the enormous price companies will pay for extra cache. The 1M Xeon costs $1,700 more than the 512K version—all for one measly cache chip! Don't forget, however, that these processors will go into expensive systems, so the performance of the processor is highly leveraged. If the additional cache improves the performance of a $50,000 four-processor server by 10%, the cost per user will actually go down, making the extra $1,700 per processor a bargain.

Intel has done quite an effective job of tuning the Deschutes processor for servers. Xeon has outstanding performance in a server environment and satisfies all the minimum MP and memory-size criteria. Its value proposition in the workstation market, however, is less clear. Vanishingly few workstations are sold with more than two processors, where Xeon shines. The full-speed backside cache helps in the workstation environment, but not as much as in big MP servers, and not nearly enough to make Xeon competitive with the RISCs on floating-point.

Despite this fact, we expect many Xeon workstation systems to appear, primarily on the basis of the strength of the Xeon brand as Intel's premium processor. From a value perspective, however, Pentium II seems like a better choice for the vast majority of NT workstations. For workstation customers looking for the ultimate in performance, the high-end RISC-based workstations still reign supreme.

### Xeon to Get Several Upgrades Before Merced

We expect that Pentium Pro will be retired almost overnight as Xeon takes over completely this quarter. Intel will follow with a 2M 450-MHz Xeon later this year. The rumored 500-MHz Tanner (see MPR 3/9/98, p. 4) is apparently on schedule to be introduced into the Xeon line during 1Q99. Tanner reportedly uses the new 0.25-micron Katmai core and will plug into Slot 2, providing a simple upgrade to the Xeon line.



**Figure 6.** Xeon does well on SPECint95 (gray bars) against all comers, but can't hold a candle to the RISCs on SPECfp95 (purple bars). Pentium II results are shown with ECC on; with ECC off, results are 15.8/11.4. (Source: *www.spec.org* except *Intel)

Sources tell us that a new processor is being planned for the Xeon line in 2H99. This part is likely to be the 0.18-micron shrink of Katmai, which could reach speeds of 600 MHz or more. The 0.18-micron Willamette core could be employed in 2000 to boost Xeon into the 800-MHz range. Whether Xeon gets extended much beyond this will depend on how smoothly the transition to Merced goes. We expect Intel will try to move the market rapidly to Merced and may let Xeon fade away in 2001—although these things have a way of dying slowly.

Intel probably wishes the x86-to-Merced transition could have been avoided in these high-end markets. We believe that Intel originally wanted to ship Merced as early as 1998. Had it done so, Xeon would not have been necessary at all. But when Merced was announced as a mid-1999 product, something like Xeon became necessary to prevent the RISCs from becoming so entrenched that even the mighty Merced would have had difficulty penetrating the market. Now, with Merced delayed further—until mid-2000—Xeon takes on more significance and must have a longer lifetime. The large body of x86 software that will be built up around Xeon cannot help but complicate Intel's transition to Merced. A Slot M Xeon, however, will provide a nice hardware evolution.

In any case, Xeon's existence is clear evidence that Intel is intent on capturing a large portion of the high-end workstation and server markets. As a strategy to increase revenue, raise its ASP, and protect its flank, this is a necessary move. Intel's branding campaign, its product differentiation tactics, and the allure of the x86 software base should allow the company to get away with charging extraordinary prices for the same CPU silicon it sells to the sub-$1,000 PC market. If Intel pulls this off, it will be a remarkable, albeit obscene, feat.

Although the x86 software base—which is the sole reason behind Pentium II's dominance of the PC market—is less compelling in the server and high-end workstation markets, we expect that its halo effect will overcome any marginal advantages the RISC competitors are able to muster. Other than raw floating-point performance, Xeon leaves few advantages the RISCs can still point to.

Therefore, we expect Xeon to be quite successful at penetrating its target markets, especially the server market, where it is somewhat more attractive than in the workstation market. We project that Xeon could grow to 4% of Intel's microprocessor revenues in 1999, accounting for over a billion dollars in sales. Considering the low manufacturing costs of Xeon relative to its selling price, such a success would have a positive impact on Intel's overall margins—and its stock price. Ⓜ