

---

# THE ORACLE SPARC T5 16-CORE PROCESSOR SCALES TO EIGHT SOCKETS

---

THE ORACLE SPARC T5 PROCESSOR MORE THAN DOUBLES THE THROUGHPUT OF THE SPARC T4 PROCESSOR, WHILE INCREASING PER-THREAD PERFORMANCE, SCALABILITY, POWER EFFICIENCY, AND I/O BANDWIDTH. THE AUTHORS DETAIL THE IMPROVEMENTS AND NEW FEATURES LEADING TO THIS LATEST ORACLE SPARC PROCESSOR.

..... The Oracle Sparc T5 processor provides a high level of system integration; excellent throughput, per-thread floating-point, and cryptographic performance; and efficient execution of heterogeneous threads. This makes Sparc T5 an ideal upgrade for customers employing Sparc systems for a range of applications, including Web servers, database and application servers, secure networking, campus backbones, file servers, and high-performance computing. In this article, we detail the Sparc T5, building on our presentation at Hot Chips 24.<sup>1</sup>

## Sparc T5 processor

The Sparc T5 processor (see Figure 1) provides increased performance, scalability, power efficiency, and I/O bandwidth over prior CPU generations. Figure 2 details its major components and interconnects.

The 28-nm Sparc T5 uses the S3 core ported from Sparc T4's 40-nm process. With 16 cores, Sparc T5 has twice as many cores as Sparc T4,<sup>2</sup> and by increasing frequency by 20 percent, it boasts more than twice the throughput. The shared 8-Mbyte, 16-way set-associative, level-three (L3) cache is twice the size of Sparc T4's L3 cache.

This maintains the cache-to-core ratio important for predictably improved performance when migrating customer applications from Sparc T4 to T5 systems.

The highly scalable interprocessor coherence architecture features a distributed, precise duplicate-tag directory, eliminating unnecessary interprocessor messaging. Each socket provides seven coherence ports to scale up to eight sockets in a fully connected single-hop configuration. The interprocessor coherence protocol includes speculative memory reads performed concurrently with directory access, cache-to-cache transfers between processors, and dynamic data routing around congested links. These features increase performance by decreasing memory latency.

The four coherence units service L3 miss requests and, in a multiprocessor system, implement an L3 directory for a subset of addresses over all the system's processors. The double data rate three (DDR3) memory subsystem supports close to 3× the bandwidth of Sparc T4 through the combination of a new memory controller, 12.8 gigabits per second (Gbps) per lane memory link, and a companion buffer-on-board chip.

John Fehrer

Sumti Jairath

Paul Loewenstein

Ram Sivaramakrishnan

David Smentek

Sebastian Turullols

Ali Vahidsafa

Oracle

Each memory controller communicates with up to two custom buffer-on-board chips, each of which drives two DDR3 1066 channels, for a total memory capacity per processor of 512 Gbytes with 32-Gbyte DIMMs.

The I/O subsystem supports  $2 \times 8$  PCI Express Generation 3 ports, doubling Sparc T4's I/O bandwidth. It also adds PCI Express atomic transactions, cache-directed direct memory access (DMA) writes, transaction identification on message-signaled interrupts, and robust direct I/O virtualization.

Power management includes features for dynamically scaling link power with workload, processor core and cache dynamic voltage and frequency scaling (DVFS), and fine-grained core-pair cycle skipping. The number of powered-on coherence links can vary by a factor of 4 as a function of link utilization. The I/O subsystem supports PCI Express power management.

Sparc T5 implements enterprise-level reliability, availability, and service (RAS) throughout, including single-chip correct memory error-correcting code (ECC) with pin sparing, single-error correction and double-error detection (SECDED ECC) in writeback caches, and parity protection in (clean) level-one (L1) caches.

### S3 core port

Sparc T5 ports Sparc T4's out-of-order, dual-issue core from 40 nm to 28 nm.<sup>3</sup> The core is clocked at 3.6 GHz, a 20-percent increase relative to Sparc T4. It has a 16-stage integer pipeline and is dynamically threaded from one to eight strands. Random number generation and encryption acceleration for 16 encryption algorithms are built in as instruction-set architecture (ISA) extensions.

### Core caches

The L1 caches and the L2 cache are private to a core. The L1 data and instruction caches are each 16-Kbyte, four-way set associative and have a line size of 32 bytes. The L1 data cache is write-through. The load-to-use latency of an L1 hit is five cycles.

The L2 cache is a unified MESI (modified, exclusive, shared, invalid) instruction and data cache. It's an eight-way, set-associative

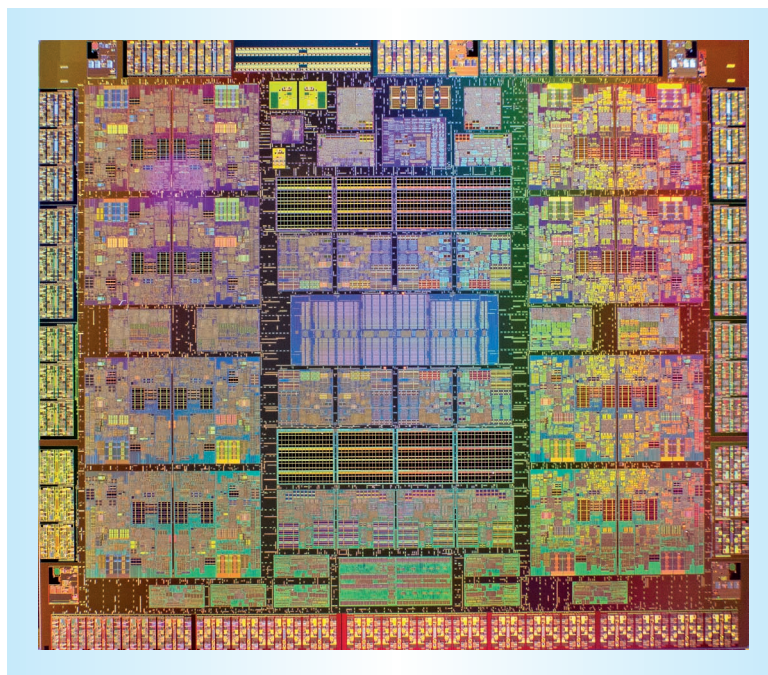


Figure 1. Sparc T5 die micrograph. The crossbar is in the center with the L3 cache banks above and below it. On each side of the crossbar are the memory control units with the cores above and below them. Below the lower L3 cache banks are the coherence units, and above the upper L3 cache banks is the I/O subsystem.

inclusive writeback cache with a capacity of 128 Kbytes and a line size of 32 bytes. The L1 directory in the L2 cache tracks the L1 tags and maintains L1 cache coherence and inclusivity. The load-to-use latency of an L2 hit is 18 cycles.

### L3 cache

The 8-Mbyte MOESI (modified, owned, exclusive, shared, invalid) L3 cache is the last-level cache and comprises eight address-interleaved 1-Mbyte banks to support the core request bandwidth. It is inclusive, 16-way set associative, and has a line size of 64 bytes. The L3 cache is shared by the 16 cores on Sparc T5 through the crossbar interconnect (see Figure 2).

An L3 cache bank is identified by address bits [8:6]. Each L3 cache bank precisely tracks the L2 cache lines that map to it using an L2 directory, which records each L2 cache line in one of three states: private (P), shared (S), or invalid (I). A P state line in the directory corresponds to either an E (exclusive) or M (modified) state line in the

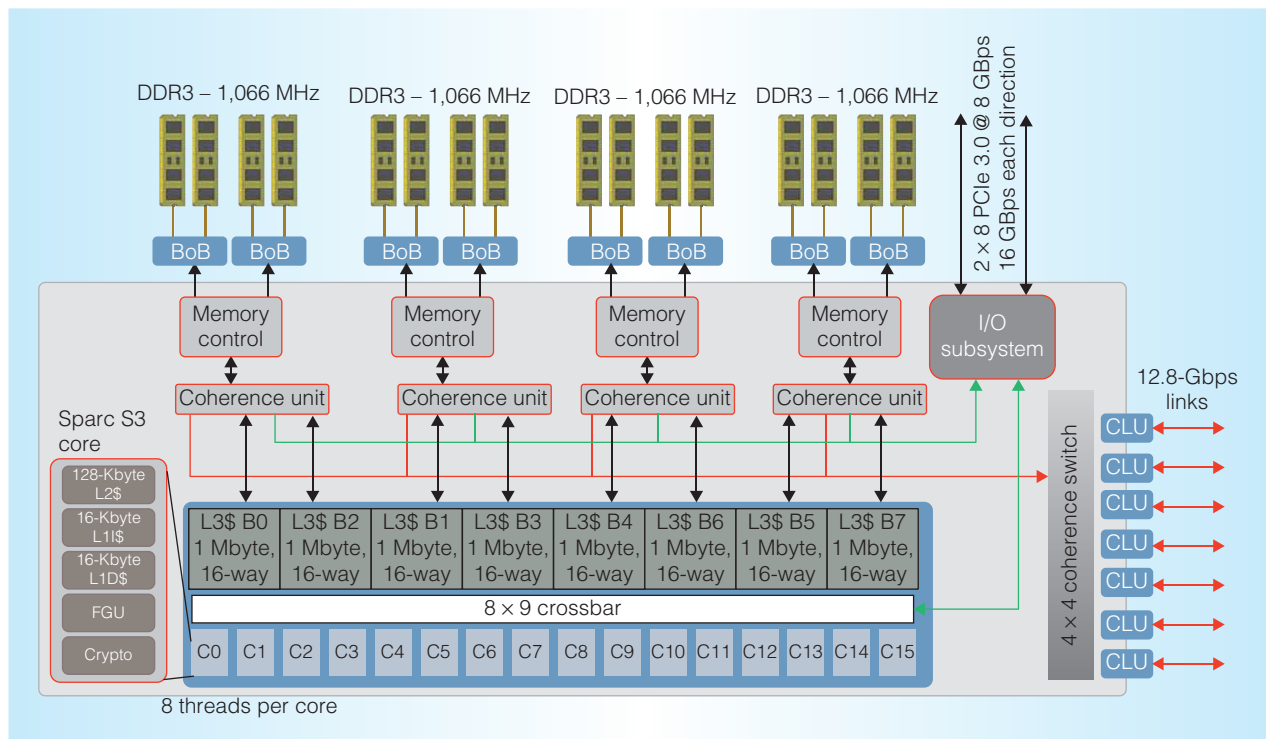


Figure 2. Sparc T5 detailed block diagram. The cores, each with its own L1 caches and L2 cache, communicate through the crossbar to 16 shared L3 cache banks and to the I/O subsystem. The coherence units handle the L3 cache misses and I/O subsystem DMA. (BoB: buffer-on-board chips; CLU: coherence link unit; DDR3: double data rate three; PCIe: PCI Express.)

L2 cache. Inclusion allows the L2 directory to track a line using its L3 location rather than its full physical tag.

Figure 3 shows the L3 cache bank's main components, which include a tag array, a cache state array (CSA), and a data array. The data array in each L3 bank is divided into eight address-interleaved subbanks, only one of which is activated during a read or fill. Two accesses to a data array subbank must be spaced two cycles apart because a data subbank returns data over two cycles. There is no such restriction on accesses to different subbanks.

When an L2 request enters the L3 pipeline, it accesses the tag array and CSA in parallel. The results of the tag compare are used to activate the appropriate data subarray. Tag misses don't result in data array activation.

A "hit" or "miss" indication is sent to the requesting L2's core using the crossbar's control network. The core uses a hit to wake up any dependent instruction. The L3 cache returns data to the L2 cache in two 16-byte beats on the crossbar's data network.

A miss on a load is conditionally used to flush the missing thread from the scheduling window, letting other threads better use the core's resources. The load-to-use latency of an L3 hit is 51 cycles.

For handling L3 misses, each L3 bank has a 60-entry miss buffer, which allows same-address requests from different cores to be chained together. If the oldest request in a chain receives an exclusive copy of a cache line, all the requests in the chain are processed before an interprocessor snoop is serviced to the same address. A shared line's L3 miss latency can thereby be amortized over multiple local requests before the line is relinquished.

The L3 allocates a cache line for a DMA write request when directed to do so by the I/O subsystem. This feature can be used efficiently to consume messages in a message-passing protocol.

## Crossbar

The high-bandwidth,  $8 \times 9$  crossbar network connects the 16 S3 cores with L2 caches to the shared L3 cache.

The crossbar is architecturally similar to the Sparc T4 processor's crossbar, but it has twice the ports on both the core and cache side of the interconnect, providing twice the bandwidth to supply data to the 16 S3 processor cores.

The S3 processor cores are grouped into pairs, with each pair sharing a crossbar port; on the other side of the crossbar, each of the eight L3 banks has a dedicated port. The aggregate bandwidth between the 16 S3 processor cores and the eight L3 cache banks is more than 800 gigabytes per second (GBps). A ninth port on the L3 side of the crossbar connects to the noncacheable unit, which interfaces to the I/O subsystem.

For packets flowing toward the L3 cache, the crossbar implements two independently arbitrated virtual channels—one for requests (for example, memory reads) and one for responses (such as snoop responses and line evictions).

For packets going toward the cores, there's a single virtual channel, carrying both L2 fill data and cache snoops targeting the individual cores. The packets in this direction propagate with a fixed latency, letting the L3 send a hit response to restart the core prior to completing the SRAM lookup in the L3 data array.

## Internode coherence for eight-socket scaling

Sparc T5 scales up to eight sockets using a glueless one-hop interconnect. Each socket has seven 14-lane, 12.8-Gbps serial links to connect to each of the other processors in an eight-socket system (see Figure 4). To allow sufficient intersocket bandwidth, the links are trunked in pairs in four-socket systems, and two-socket systems use four links.

Multiprocessor coherence is maintained using a directory, which precisely tracks all L3 cache lines in the system. This directory is partitioned across all available sockets on the basis of physical address bits [10:8]. For example, in a two-socket system, all L3 cache lines corresponding to PA[8] = 0 are tracked by socket 0, whereas socket 1 tracks all the L3 cache lines with PA[8] = 1. The directory doesn't distinguish between the different valid L3 cache states.

A socket maintains its directory partition in a flexibly organized array. When the

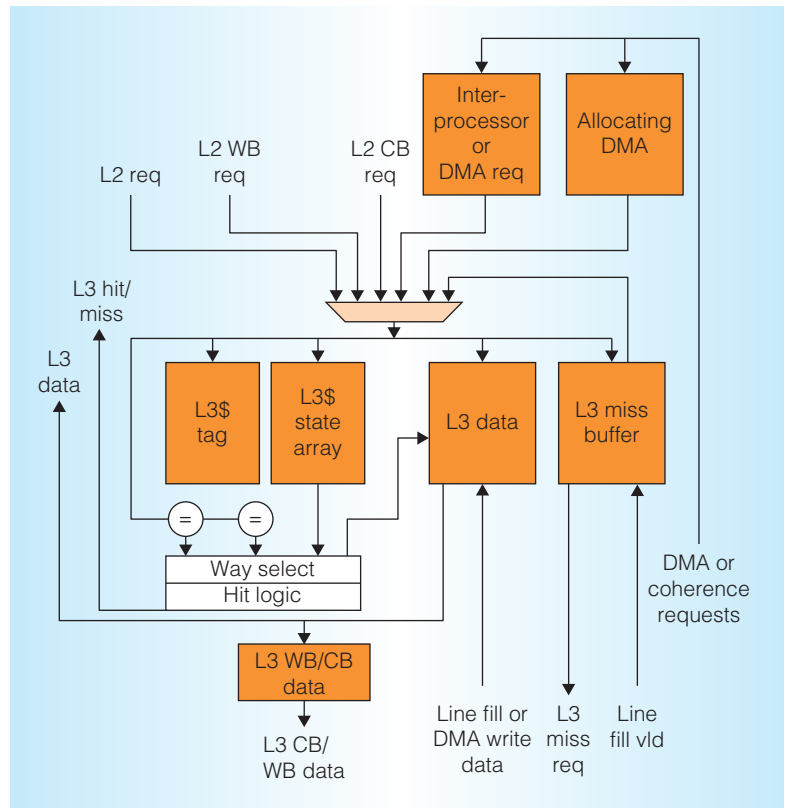


Figure 3. L3 cache. The cache bank's main components include a tag array, a cache state array (CSA), and a data array. (CB: copyback; DMA: direct memory access; WB: writeback.)

system socket count is a power of 2, the directory is organized as follows:

- Two-socket: 8,192 indices  $\times$  32-way set associative.
- Four-socket: 4,096 indices  $\times$  64-way set associative.
- Eight-socket: 2,048 indices  $\times$  128-way set associative.

The directory also supports three- and six-socket configurations. In these modes, a lookup table is used to map an L3 cache line to a directory partition.

The socket ID of a line in the directory is implied by its position in the directory; L3 cache lines from socket  $n$  are maintained in ways  $16n$  to  $16n + 15$ .

## Protocol transaction flow

We explain the directory-based protocol here for an L3 miss request for an exclusive copy of a line. Upon an L3 miss, a request

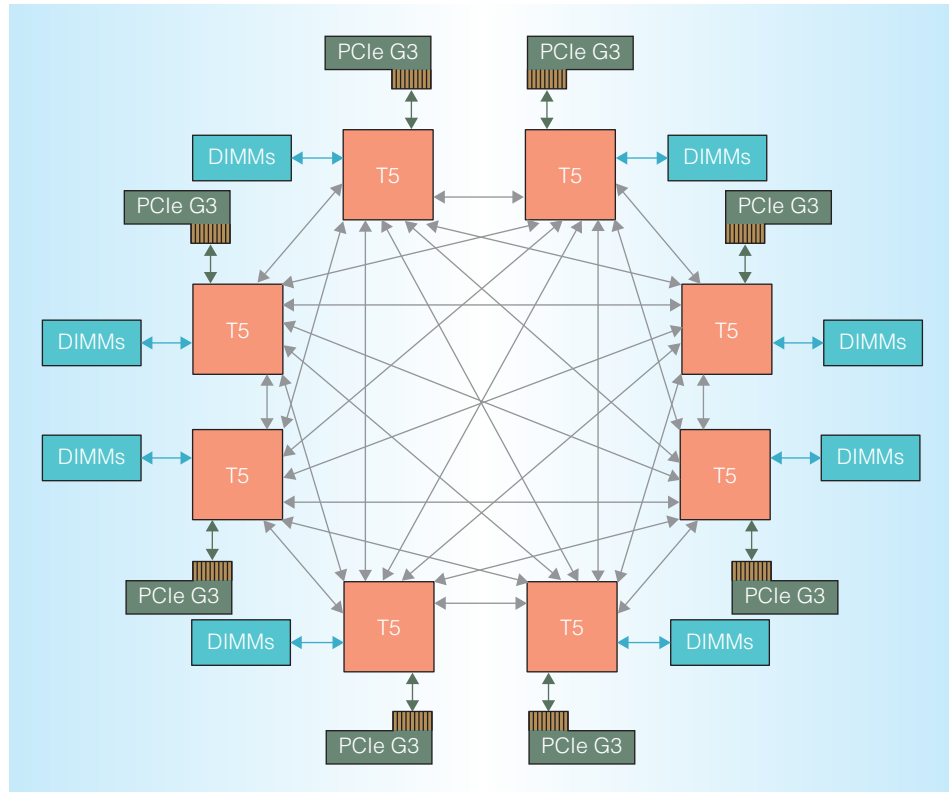


Figure 4. Eight-way interprocessor connections. A fully connected network provides 1-hop communication between any two processors.

is sent to the appropriate directory partition, from which the sharing processors are determined. The directory is updated to indicate the requester as the sole sharer. In the absence of sharers, a memory read request is sent to the memory home; otherwise, a copyback invalidating request is sent to one of the sharers, and an invalidating request is sent to each of the other sharers.

The requester can send a speculative request to memory in parallel with the coherence directory lookup. The speculative request is used to read memory into a tagged prefetch buffer at the memory controller. If the memory read request from the directory hits in the prefetch buffer, the memory controller services it from there instead of initiating a new memory access.

When a socket receives a copyback or invalidation request, it forwards the request to the appropriate L3 bank. The L3 bank responds with data or a snoop

acknowledgment, which is forwarded to the requester.

## Power management

The Sparc T5 processor in conjunction with software provides fine-grained power and performance management.

Power management software selects the target DVFS performance-state (P-state) and cycle skip ratio. The user selects the power management policy using the integrated lights out manager (ILOM). In the *performance* policy, the system runs at peak frequency while throttling to lower P-states only in response to over-temperature or over-current situations. In the *elastic* policy, the system trades off between performance and power. The operating system specifies each core's frequency requirement. This request is translated into the target DVFS P-state and clock cycle skipping ratio by power management software.

Each entry in the hardware P-state table contains a *voltage ID* (VID), an identifier



to set the digital voltage regulator module (VRM) and the clock source divisors. The chip switches dynamically between P-states in response to throttle and resume signals from a T5 mezzanine (T5M) component for processor temperature and current capping (see Figure 5). The DVFS FSM sequences a frequency change to a clock source and a VID transmission to the VRM via the T5M monitoring the VID chip output. In addition to thermal capping, the system relies on current capping to prevent oversubscribing of the VRM. If the VRM current output is greater than a threshold for a specific time period, it asserts a throttling signal to the T5M. The throttling is removed as soon as the current undershoots the threshold by a hysteresis value. The Sparc T5 system also monitors overall power consumption (I/O, DRAM, CPU, fans, and so on) and imposes power constraints on specific sockets by inducing the T5M to assert throttle or resume to the processor.

During production, Sparc T5 processors are characterized at three P-states:  $P_0$ ,  $P_1$ , and  $P_{last}$ .  $P_{last}$  corresponds to the maximum frequency achieved at the minimum operating voltage. The voltage and frequency pairs determined for these points are stored in the electrical fuse ROM (EFUSE). With minimal test overhead, measuring three states improves yield compared with merely extrapolating from a single characterization point. To minimize static voltage guardband, EFUSE stores a power supply calibration (PSC) value corresponding to each of these three voltages and frequency pairs. Firmware reads data from the EFUSE at boot time and populates the P-state table using an interpolation algorithm for entries that don't have a fused value.

Cycle-skipping implemented on the Sparc T5 processor enables software to control the frequency of core-pairs. The cycle-skipping mechanism supports a granularity of one-eighth; target core-pairs operate at  $n/8$  times the core clock frequency, where  $1 \leq n \leq 8$ . The hardware cycle-skipping control gradually changes toward the requested target ratio, ensuring that only one core-pair changes its ratio at a time, in order to prevent  $L(didd)$  voltage droop.

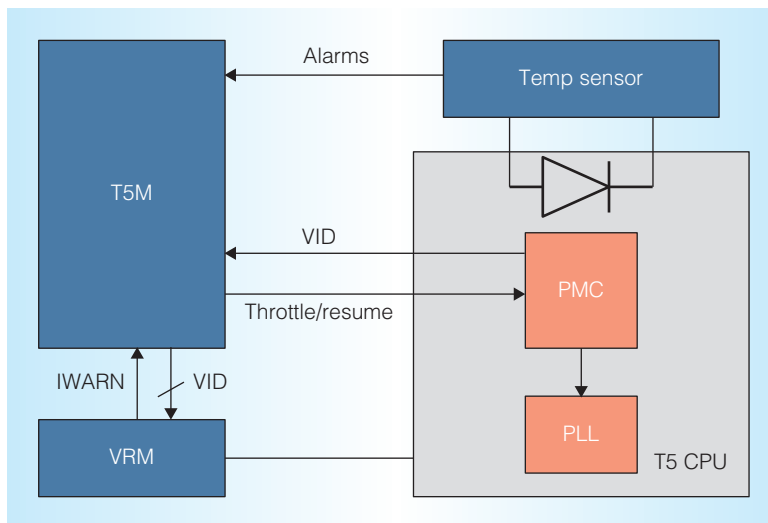


Figure 5. Dynamic voltage and frequency scaling (DVFS) feedback. The T5 mezzanine chip (T5M) mediates the control of throttling and voltage. (PLL: phase-locked loop; PMC: power management controller; VRM: voltage regulator module.)

The T5 crossbar runs at the full processor frequency, while tracking each core port's clock skipping state. The crossbar dynamically adapts the signaling on each core port. The L3 cache banks are completely isolated from the clock-skipping behavior of the processor cores. Because the clock edges are synchronous, the latency of a traditional clock domain crossing scheme is avoided.

In smaller system configurations with trunked links, links can be dynamically shut down to modulate serializer/deserializer (SerDes) power on the basis of bandwidth demand. Lab measurements have demonstrated 25-W power savings in a two-socket configuration when scaling down from four active links to one.

## Reliability, availability, and service

The Sparc T5 processor has RAS features typically found in high-end enterprise systems.

Typically, there is end-to-end parity protection for address and control networks, and ECC for architectural state, caches, and data networks. Hardware and in-memory firmware hide most errors from the user code. An error condition typically forks off two parallel error-handling tracks: damage containment and failure analysis.

Correctable errors are corrected and cleared, and the original instruction is invisibly retried. For uncorrectable errors, the damage is contained to the strand accessing the data. In the background, hardware captures extensive information about the error type and location, and activates error analysis firmware. This firmware retains a record of previously seen errors and can deconfigure components or initiate service action.

An error in a cache evicts the cache line, correcting all single-bit errors and quashing all double-bit errors in clean data. The resulting retry uses a cache bypass path, allowing forward progress with stuck faults. A cache line exhibiting a high error rate can be retired. A line can be retired aggressively on the basis of a quick analysis, and then unretired after a more comprehensive analysis. Both retiring and unretiring of cache lines is performed during normal operation, with no need to quiesce traffic or reboot. If a line is permanently retired, the information is retained in off-chip nonvolatile storage and applied during subsequent boots.

All information carried over each link is protected by a cyclic redundancy code (CRC). The number of check bits and the corresponding polynomials have been optimized so that each link has complete lane and unit interval (UI) coverage. To minimize performance impact when a CRC error is detected, the recovery process starts with a simple replay and escalates through progressively more rigorous forms of training. One faulty lane per link can be deconfigured.

In the memory controller, a scrubber engine touches all memory periodically to prevent accumulation of multiple soft errors in a single line. There's also a characterization engine, which upon detecting a memory error, reads and writes the failure location to determine whether the error is transient or persistent.

Beyond the standard single-chip-correct, double-chip-detect capabilities, memory ECC supports a spare pin per DIMM. When a memory error is detected, the error-analysis firmware looks for trends that would imply a fault localized to a wordline or bitline. Wordline faults trigger page retire, a software sequence that lets a cooperative operating system take a physical page offline.

Bitline faults, which include faults associated with device pins, trigger activation of spare pin mode. The hardware transparently relocates data within the same DIMM, so that the faulty bitline or pin is no longer used. Additional enterprise-level RAS features include a strand identification scheme, which allows for deconfiguring one processor without affecting other processors, and a dynamic interrupt routing scheme, which allows a processor to handle another processor's errors.

## PCI Express Gen3 I/O

The I/O subsystem in Figure 6 implements a root complex compliant with PCI Express Base Specification Revision 3.0, which provides an 8-Gbps serial signaling rate per lane. The root complex presents an eight-lane root port, yielding 8 GBps of raw bandwidth in each direction. Software sees two distinct PCI Express fabrics rooted by the respective I/O subsystems.

The I/O subsystem is a bridge between the host (CPU cores and memory) and the PCI Express fabric. The host connection is through the I/O switch (IOX) interconnect connecting I/O subsystems with coherence units and the noncacheable unit (NCU). The I/O subsystem comprises three units: the data management unit (DMU), the PCI Express unit (PEU), and the PCI Express SerDes (PSR). The PEU implements the layers of the PCI Express hardware protocol stack. Its physical coding sublayer (PCS) at the bottom of the physical layer interfaces with the PSR.

The DMU contains the address-translation unit (ATU), which converts from I/O virtual addresses to physical addresses; the interrupt-management unit (IMU), which processes and records PCI Express message signaled interrupts (MSI) and extensions (MSI-X) interrupts and messages; and the noncacheable processing unit (NPU), which routes programmed I/O (PIO) accesses to internal functions or external PCI Express fabric, and handles errors. The PEU contains the PCI Express switch (PEX), which processes PIOs and implements error handling, and the transaction layer unit (TRU), which implements transaction layer functions and manages PCI Express flow control.

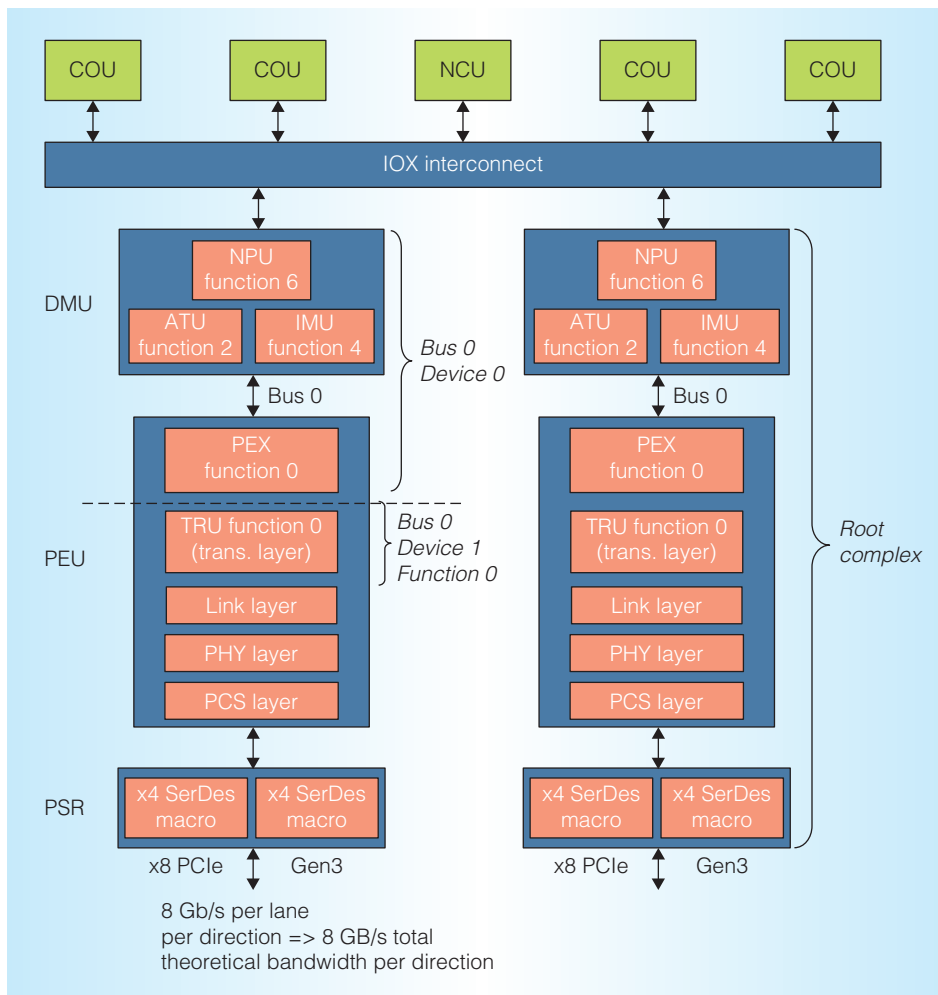


Figure 6. I/O subsystem. The noncacheable unit (NCU) handles programmed I/O and interrupts. DMA is handled by the coherence units (COUs). (ATU: address-translation unit; IMU: interrupt-management unit; IOX: I/O switch; NPU: noncacheable processing unit; PCS: physical coding sublayer; PEX: PCI Express switch; PHY: physical protocol layer; TRU: transaction layer unit.)

The I/O subsystem programming model divides noncacheable space into three main regions: PCI Express configuration space, legacy I/O port space, and PCI Express memory space. The latter is subdivided using PCI base address registers (BARs) in each internal (I/O subsystem) and external PCI Express function. Software can allocate different amounts of PCI Express memory space to different functions, creating a high degree of adaptability to platform needs.

The I/O subsystem supports I/O virtualization (IOV) in its ATU and IMU so that multiple guest operating-system instances can share the I/O subsystem and the

underlying PCI Express fabric. PCI Express endpoint devices that contain multiple virtual functions (VFs) following the single-root IOV (SR-IOV) scheme facilitate this sharing. The hypervisor can allocate different guest operating-system instances to different VFs on the physical PCI Express network or disk controller device, bringing better efficiency and reduced cost in the I/O subsystem.

The ATU includes multiple translation storage buffers (TSBs) maintained in memory and cached in the ATU. TSBs are indexed by the PCI Express requester's ID, thus meeting the protection, performance,



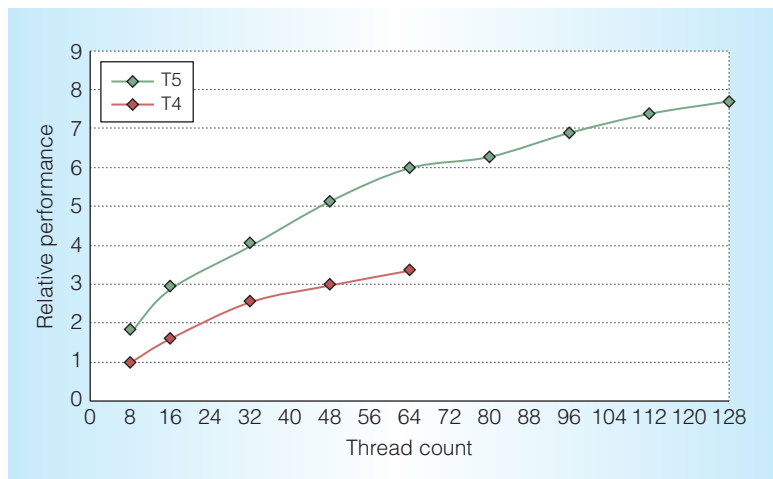


Figure 7. Sparc T4 per-thread Java comparison. Sparc T5 achieves 1.75 $\times$  the per-thread of Sparc T4 while providing twice as many threads.

and RAS requirements for virtualization. If one guest operating system panics or suffers throughput degradation, it doesn't disrupt another guest operating system sharing the I/O subsystem or PCI Express devices. The ATU's bounds checking prohibits one guest from accessing physical memory allocated to another guest.

The IMU supports a total of 1,024 MSI and MSI-X vectors, which are allocated by the hypervisor to endpoint devices during initialization and bus probe. A PCI Express message or MSI-X vector is mapped to one of 256 in-memory event queues (EQs). The IMU records PCI Express messages and MSI-X vectors to their associated EQs in memory, signaling an interrupt if an EQ becomes nonempty. Each EQ is mapped to a specific processor, core, and strand. The hypervisor distributes interrupt handling load and can steer interrupts from specific VFs to the same cores and strands managing the DMA from those VFs. The distribution can be adjusted dynamically as workloads change to optimize system performance.

To implement PCI Express ordering rules, an acknowledgment from the coherence unit signals to the I/O subsystem that a DMA write has been ordered with respect to other memory operations. Vital to achieving a high DMA write line rate on large systems, the I/O subsystem supports the optional relaxed-ordering PCI Express feature to increase the performance of DMA

writes, allowing writes to pass each other in the memory system.

Supported optional PCI Express-architected RAS and security features include advanced error reporting (AER), end-to-end CRC (ECRC) checking and generation, and access-control services (ACS). Other optional PCI Express features are more complete support for atomic transactions and TLP processing hints, allowing the targeting of DMA writes at an L3 cache instead of memory.

It's possible to reset an I/O subsystem without disrupting the other I/O subsystems or any other units attached to the coherence switch (see Figure 2). This contains the effects of a fatal I/O subsystem error, such as a parity error detected on an internal header. When the hypervisor is informed of the error, it can reset and reinitialize the I/O subsystem and reprobe the PCI Express fabric attached to it, with no interruption of service to PCI Express devices connected to the other I/O subsystem.

Each I/O subsystem unit contains performance counters that have proved quite useful in hardware and software validation and in benchmarking for the Sparc T5 platform.

Sparc T5 has 16 S3 cores running at 3.6 GHz. With twice the number of cores as Sparc T4 and each running 20 percent faster, Sparc T5's instruction execution rate is 2.4 $\times$  that of its predecessor. Sparc T5 also has double the L3 cache size and has 2.5 $\times$  the memory bandwidth of Sparc T4.

Sparc T5 achieves 2.25 $\times$  the peak-throughput Java performance of Sparc T4. Faster cores and twice the cache size improve per-thread performance by 1.75 $\times$  for low-thread-count applications. As Figure 7 shows, more cores and higher-memory bandwidth helps maintain that improvement over the full range of thread counts.

Sparc T5 also doubles the cryptographic performance delivered by Sparc T4. Sparc T5 extends its single-socket performance gains to multisocket systems. Sparc T5 delivers excellent multiprocessor performance, scaling across a range of workloads (see Figure 8). Online transaction processing (OLTP) scaling is an important indicator of multisocket scaling, because the

dataset is spread across all processors with low affinity to processing threads. A directory-based intersocket coherence protocol and 840-Gbps intersocket bisectional bandwidth enable Sparc T5 to achieve greater than  $7\times$  OLTP scaling on eight-processor systems. Sparc T5 supports shared-memory systems of up to eight processors, delivering  $4.5\times$  the peak throughput of the largest Sparc T4 system.

MICRO

## References

1. S. Turullols and R. Sivaramakrishnan, "Sparc T5: 16-Core CMT Processor with Glueless 1-Hop Scaling to 8-Sockets," *Hot Chips* 24, 2012.
2. M. Shah et al., "Sparc T4: A Dynamically Threaded Server-on-a-Chip," *IEEE Micro*, vol. 32, no. 2, 2012, pp. 8-19.
3. J.L. Shin et al., "The Next-Generation 64b Sparc Core in a T4 SoC Processor," *Proc. IEEE Int'l Solid-State Circuits Conf.*, IEEE Press, 2012, pp. 60-62.

**John Feehrer** is a senior staff engineer at Oracle. His research interests include I/O subsystem architecture, distributed computing, and optical interconnects and switching. Feehrer has a PhD in electrical engineering from the University of Colorado Boulder.

**Sumti Jairath** is a senior hardware architect at Oracle. His research interests include hardware-software interaction, software-specific hardware design, and coherence scaling. Jairath has a BTech in electronics and communication engineering from National Institute of Technology, Kurukshetra, India.

**Paul Loewenstein** is a senior principal hardware engineer at Oracle. His research interests include multiprocessor cache coherence protocols, formal memory models, deadlock and livelock avoidance, and error detection and correction. Loewenstein has a PhD in computer science from the University of Cambridge.

**Ram Sivaramakrishnan** is the director of hardware at Oracle. His research interests include low-power cache design for high throughput multithreaded architectures

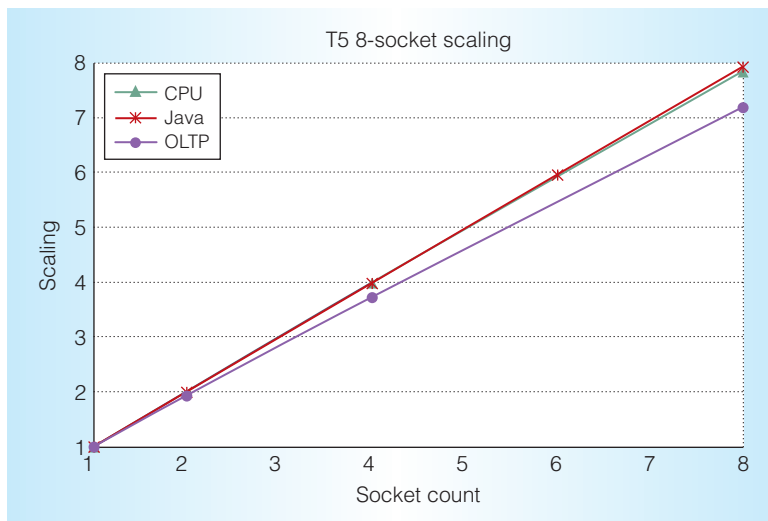


Figure 8. 8-socket scaling performance. Even on low-affinity applications such as online transaction processing (OLTP), Sparc T5 achieves greater than  $7\times$  scaling on 8-processor systems.

and low-overhead multiprocessor coherence protocols. Sivaramakrishnan has an MS in electrical engineering from the University of Nevada, Las Vegas.

**David Smentek** is a senior hardware architect at Oracle. His research interests include computer architecture, computer arithmetic, performance modeling, and on-chip networks for multicore architectures. Smentek has an MS in electrical engineering from Stanford University.

**Sebastian Turullols** is a hardware director at Oracle. His research interests include high-speed interchip interconnect, power efficiency, and multiprocessor performance scaling. Turullols has an MS in electrical engineering from Stanford University.

**Ali Vahidsafa** is a senior principal hardware engineer at Oracle. His research interests include characterization and testability of analog-digital interfaces and improving the efficiency of dynamic voltage scaling. Vahidsafa has an MS in electrical engineering from San Diego State University.

Direct questions and comments about this article to Sebastian Turullols, Oracle, Floor 3, 4120 Network Circle, Santa Clara, CA 95054; sebastian.turullols@oracle.com.