

## IT@INTEL

# Hyperscale High-Performance Computing For Silicon Design

- Five generations of HPC have successfully enabled Intel® silicon tape-out, reducing tape-out time from 25 to less than 10 days.<sup>1</sup>
- >90x growth in HPC in the last 10 years since inception with >64x increase in stability.
- Intel has taped out several silicon products with HPC-1 alone, delivering a return on investment of USD 44.72 million.<sup>1</sup>

### Shesha Krishnapura

Intel IT Chief Technology Officer and Senior Principal Engineer

### Ty Tang

Senior Principal Engineer, Intel IT

### Vipul Lal

Senior Principal Engineer, Intel IT

### Matt Ammann

HPC Network Architect, Intel IT

### Raju Nallapa

Principal Engineer, Intel IT

### Doug Austin

Storage and Backup Architect, Intel IT

### Shaji Achuthan

Senior Staff Engineer, Intel IT

## Executive Overview

Designing Intel® microprocessors is extremely compute intensive. Tape-out is a final step in silicon design and its computational demand is growing continuously for each generation of silicon process technology. Intel IT adopted high-performance computing (HPC) to address this very large computational scale and realized significant improvements in computing performance, reliability, and cost.

We treated the HPC environment as a holistic computing capability—ensuring all key components were well designed, integrated, and operationally balanced with no bottlenecks. We designed our HPC model to scale to meet future needs, with HPC generations aligned with successive generations of Intel® process technology.

The first-generation HPC environment (HPC-1), supporting 45nm processor tape-out, included innovative approaches and technologies to increase scalability, such as the following:

- A parallel storage system providing 10x scalability compared with our previous system based on traditional file servers, together with high-speed backup.
- Large-memory compute servers based on a unique modular non-uniform memory access (NUMA) design, offering significant cost advantages.
- Batch compute servers based on multi-core Intel® Xeon® processors, offering substantial performance increases.
- Optimization of our license server and job scheduler to handle thousands of simultaneous design jobs.

HPC-1 successfully enabled 45nm processor tape-out, delivering net present value (NPV) of USD 44.72 million to Intel.<sup>1</sup> We subsequently developed four new generations of HPC environments (HPC-2, HPC-3, HPC-4, and HPC-5), with further scalability increases to support the tape-out of 32nm, 22nm, 14nm, and 10nm processors, respectively.

Since deployment, our HPC environment has supported a 90.21x increase in compute demand, with a 64.4x increase in stability. In addition, tape-out

**Contents**

- 1 Executive Overview**
- 2 Business Challenge**
  - Tape-out Challenges
- 4 Solution: High-Performance Computing Strategy**
  - Storage and Backup
  - Compute Servers
  - Network
  - Batch Clustering: Job Scheduler Improvements
  - EDA Application License Servers
  - Enterprise Linux\* OS
  - Application and Platform Tuning
- 16 HPC Benefits**
- 17 Key Learnings and Future Plans**
- 18 Conclusion**

**Acronyms**

<b>DRC</b>	design rule check
<b>EDA</b>	electronic design automation
<b>FSB</b>	front-side bus
<b>GB</b>	gigabyte
<b>Gb/s</b>	gigabits per second
<b>GHz</b>	gigahertz
<b>HPC</b>	high-performance computing
<b>HPC-1</b>	first-generation HPC environment
<b>HPC-2</b>	second-generation HPC environment
<b>HPC-3</b>	third-generation HPC environment
<b>HPC-4</b>	fourth-generation HPC environment
<b>HPC-5</b>	fifth-generation HPC environment
<b>MB</b>	megabyte
<b>Mb/s</b>	megabits per second
<b>MB/s</b>	megabytes per second
<b>NAS</b>	network-attached storage
<b>NPV</b>	net present value
<b>NUMA</b>	non-uniform memory access
<b>OPC</b>	optical proximity correction
<b>RET</b>	resolution enhancement techniques
<b>SSD</b>	solid-state drive
<b>TB</b>	terabyte
<b>TiB</b>	tebibyte

time was reduced from 25 days for the first 65nm process technology-based microprocessor in a non-HPC compute environment to 10 days for the first 45nm process technology-based microprocessor in an HPC-enabled environment.<sup>1</sup> The success of the HPC environment was due to factors such as careful alignment of technology with business needs, informed risk taking, and disciplined execution. We are continuing to develop the future HPC generations to enable tape-out of successive generations of Intel® processors.

We have continuously worked over HPC generations to improve the environmental reliability. We test the HPC environment at 1.15-2x scale before we optimize. To add redundancy, we have multiple physical sites across the WAN to enable sites to be operational to run silicon design workloads, even if any other site is not operational. With HPC-5, we are working to introduce the concept of sub-sites within a site, to address continued operation during local failures. This is currently in progress.

## Business Challenge

Microprocessor design is extraordinarily complex—and as a result, requires huge amounts of computing capacity. About 120,000 of the servers in Intel’s worldwide environment are dedicated to silicon design.

Each new generation of process technology—such as the transition from 65nm to 45nm processors—brings a substantial increase in complexity, requiring a major increase in design compute performance.

Though increased performance is needed across the entire design process, the requirement is particularly acute at the highly compute-intensive tape-out stage.

Tape-out is a process where Intel® chip design meets manufacturing. As shown in Figure 1, it is the last major step in the chain of processes leading to the manufacture of the masks used to make microprocessors.

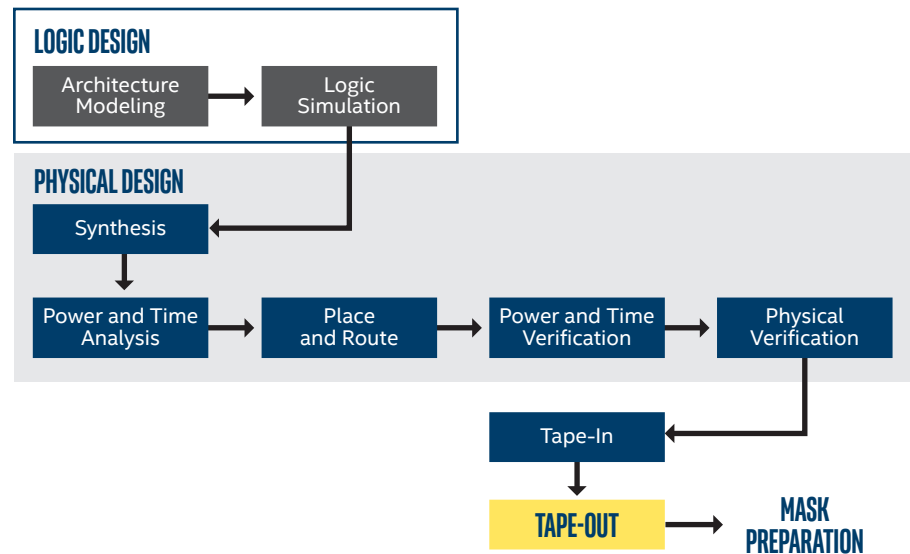


Figure 1. The phases of silicon design.

During tape-in, the stage that immediately precedes tape-out, Intel chip design teams create multi-gigabyte hierarchical layout databases specifying the design to be manufactured. During tape-out, these layout databases are processed using electronic design automation (EDA) tools. These tools apply extremely compute-intensive resolution enhancement techniques (RET) to update layout data for mask manufacturability and verify the data for compliance to mask manufacturing rules.

A key EDA application within the tape-out stage is optical proximity correction (OPC), which makes it possible to create circuitry that contains components far smaller than the wavelength of light directed at the mask. OPC is a complex, compute-bound process. To accelerate the process, OPC applications take advantage of distributed parallel processing; tasks are divided into thousands of smaller jobs that run on large server clusters.

It is critical to complete tape-out as fast as possible—and to minimize errors—because delays at this stage can mean slipped project deadlines and even a missed market window.

## Tape-out Challenges

Up to and including the 65nm process technology generation, tape-out computing was managed as an extension of our general-purpose design computing environment.

However, as we prepared for the transition to the first Intel® 45nm processors, it became apparent that we needed a new approach to create a cost-effective, reliable, and predictable environment capable of supporting the increased demands of 45nm processor tape-out.

Overall, we anticipated that we would need a 10-fold increase in scalability. Key challenges included:

- **Storage.** We anticipated a requirement for a 10x increase in storage system throughput. However, our existing production network-attached storage (NAS) file servers were already experiencing I/O bottlenecks even before the transition to 45nm technology.
- **Compute servers.** The compute servers used to run the largest tape-out jobs could not support the anticipated 4x increase in physical memory requirements.
- **Stability.** Our existing production environment was not designed to support very large-scale tape-out computing. Because of this, it was less reliable than desired, leading to more than 20 tape-out delays per quarter.
- **Cost.** We needed to solve these technical challenges while meeting the requirement to reduce capital expenditure by USD 20 million.

We expected this growth trend to continue in future process generations. This meant we needed an approach that could both support 45nm tape-out and subsequently scale to meet future needs.



## 4 Tape-out Challenges to 10x Scalability

- Storage
- Compute Servers
- Stability
- Cost

To solve these challenges, we set out to develop a high-performance computing (HPC) environment optimized for tape-out processing, using large compute server clusters and disruptive technologies to deliver substantial increases in scalability and performance.

## Solution: High-Performance Computing Strategy

In 2005, we created an HPC strategic program to develop a highly scalable and reliable tape-out compute environment that is capable of delivering optimal results. Developing our HPC environment presented significant challenges because this was the first time HPC was attempted for semiconductor design.

Strategic objectives included the following:

- Leverage industry and internal expertise to architect a leading-edge HPC environment
- Design a solution that is highly customized for tape-out
- Use open standards to develop an agile environment
- Regularly benchmark and adopt best-in-class HPC technology

Our immediate goal was to enable the tape-out of the first Intel 45nm processors to meet our committed deadline to Intel® product groups.

Our longer-term objective was to develop an HPC generational model that could meet future needs, aligned in lockstep with successive generations of Intel® process technology, as shown in Figure 2. Each HPC generation would provide a major increase in capacity to support the demands of the corresponding new processor generation.

For the first generation of the HPC environment (HPC-1), our goal was to achieve an overall 10x increase in scalability.

Our approach was to treat the HPC environment as a holistic computing capability—ensuring that critical components were well-designed, integrated, and operationally balanced with no single bottleneck. These components were the following:

- Storage and backup
- Compute servers
- Network
- Batch clustering and job scheduling
- Application license servers
- Enterprise Linux\* OS
- Application and platform tuning

The solution stack that delivers our HPC environment is shown in Figure 3, on the next page.

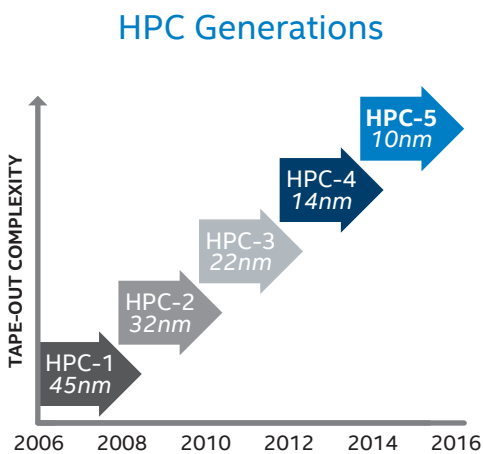


Figure 2. We aligned our high-performance computing (HPC) environment with process technology generations.

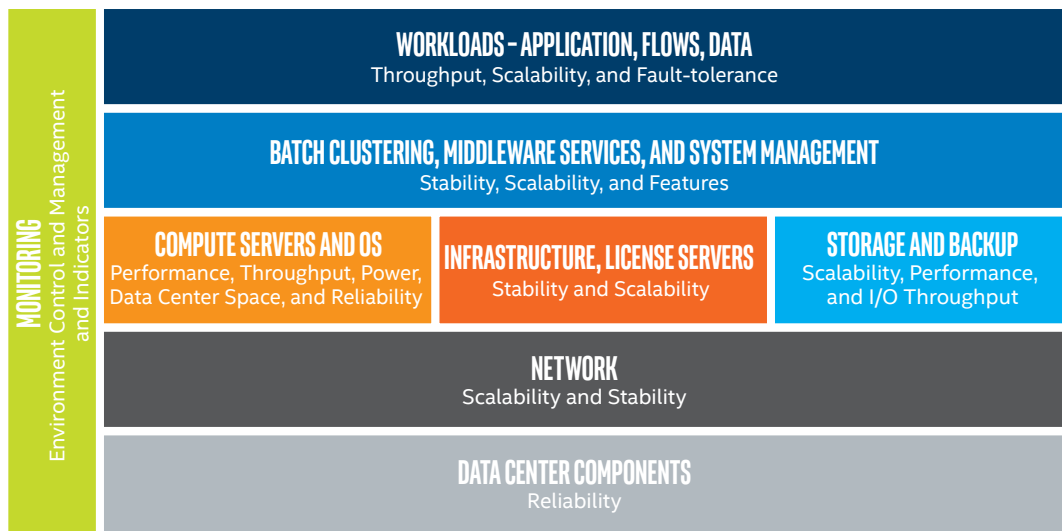


Figure 3. The high-performance computing (HPC) solution stack is well-designed, integrated, and operationally balanced.

We assessed the performance of each component using real tape-out workloads. We identified bottlenecks and the improvements needed in each area. Then members of the HPC program and Intel’s manufacturing group jointly defined the HPC-1 specifications.

We have continued this approach with subsequent HPC generations to achieve the increases in scalability required for successive processor generations.

Beginning in 2007 we designed and implemented the second, third, fourth, and fifth generations of our HPC environment to provide the increased compute resources required to support tape-out of 32nm, 22nm, 14nm, and 10nm processors.

We have made substantial improvements in the key components, outlined in the following sections.

## Storage and Backup

We identified storage performance and scalability as significant bottlenecks. We implemented a parallel storage system to deliver the anticipated 10x increase in required scalability in HPC-1; by HPC-5, compute demand has increased 90.21x. For HPC-5, we enhanced scalability beyond the single datacenter environment. We combined this with a faster backup solution capable of handling the required throughput and much larger disk volumes.

### Storage

For the 65nm processor generation, we had been using traditional NAS file servers, which were able to serve only 400 distributed clients and had a volume size limit of 400 gigabytes (GB).

For the 45nm generation, we needed to support up to at least 4,000 clients—a 10x increase—and volume sizes up to 3 terabytes (TB). To achieve this with the existing solution would have required at least 10 additional storage server racks. This was not an option because of the resulting increases in our data center footprint as well as power and cooling costs. An additional problem was that the need to replicate large design datasets across multiple storage servers to work around scalability limitations affected the productivity of our design engineers.

We therefore decided to research parallel storage solutions that would not only satisfy our current storage needs but also easily scale to future demands. The storage solution needed to deliver higher performance with a significantly lower total cost of ownership (TCO). We considered more than 20 possible solutions and selected one after an extensive evaluation, including onsite testing with real tape-out workloads that consumed more than 1 million CPU hours.

The deployment of our parallel storage solution was a milestone; it was a pioneering use of parallel storage in an IT organization in the semiconductor industry. Specifications of HPC-1 through HPC-5 parallel storage are summarized in Table 1.

To support the >2x computational demand and >4x increase in system count between HPC-4 and HPC-5 and anticipated future increase, we had to revisit our storage selection and our computing infrastructure setup. After intensive testing with synthetic and real compute workloads that consumed several million hours over an extended pilot period, we selected a scale-out NAS solution. The second challenge was to enable us to continue to scale without having a single failure domain which impacts the entire capability. This led to working on having multiple “independent sites” within a single physical site and also supporting multiple physical sites across the WAN enabling us to achieve hyperscaling. The concept of “independent sites” is to eventually be able to have one site continue to support running workloads even if the rest of the sites are unavailable. Our key applications have been updated to run across these independent sites and to also significantly reduce their storage IOPS demand. Key infrastructure including storage, compute, and scheduling have been set up to support the independent sites. We currently have three such sites and expect to go to four sites before the end of 2015.

After intensive testing with synthetic and real compute workloads that consumed several million hours over an extended pilot period, we selected a scale-out NAS solution.

Table 1. Storage System Specifications for High-Performance Computing (HPC) Environments, HPC-1 through HPC-5

COMPONENT	GENERATIONS				
	HPC-1	HPC-2	HPC-3	HPC-4	HPC-5
Storage Server	Intel® Celeron® processor	Intel® Celeron® M processor 370	Intel® Celeron® M processor 370	Intel® Xeon® processor LC3518	Intel® Xeon® processor E5-2658
CPU Specification	1.2 GHz, 256 KB L2 cache	1.5 GHz, 1 MB L2 cache	1.5 GHz, 1 MB L2 cache	1.73 GHz, 256 KB L2, 2 MB L3 cache	2.1 GHz, 20 MB Intel® Smart Cache
Bus	100 MHz FSB	400 MHz FSB	400 MHz FSB	DMI 2.5 GT/s	QPI 8 GT/s
RAM	512 MB	2 GB	4 GB	8 GB	64 GB
RAM Type	PC 100 SDRAM	ECC DDR2-400	ECC DDR2-400	ECC DDR3-800	ECC DDR3-1600
Raw Storage System Capacity	80 TB	100 TB	93 TB	180 TB	192 TB

DDR2 – Double Data Rate 2; DDR3 – Double Data Rate 3; DMI – direct media interface; ECC – error correction code; FSB – front-side bus; GB – gigabyte; GHz – gigahertz; GT/s – giga-transfers per second; KB – kilobyte; MB – megabyte; MHz – megahertz; QPI – Intel® QuickPath Interconnect; RAM – random access memory; SDRAM – Synchronous dynamic random access memory; TB – terabyte

---

We have continuously worked over HPC generations to improve the environmental reliability.

---

### Reliability Improvements over Generations

We have continuously worked over HPC generations to improve the environmental reliability. We follow multiple key enablers stated below to achieve reliability and stability:

- Establish and adhere to standards starting from physical hardware components all the way to the tools for executing the number and mix of jobs running in the environment.
- Ensure technology changes follow a very strict control process with the tape-out teams as the decision maker. Approval for changes that have a proven business values depends on the following:
  - Successful completion of a mix of synthetic test and a representative mix of end-to-end workflows
  - Testing at 1.15x – 2x the anticipated peak scale of running jobs, job ramp, and exit rates
  - A successful pilot covering production usage at scale
  - The deployment schedule to deliver maximum value at minimum risk
- Ensure technology changes are aligned to meet the demands of the process change and are an integral part of the process development and ramp.
- Perform fully automated configuration checking and control of all infrastructure components starting from firmware to tool flows with configuration changes reviewed through a change control process.
- Continuously monitor all components including trending data to foresee upcoming bottlenecks.
- Root causing all issues that impact workflows with follow-up to improve existing environment and mandatory test addition for future changes. In the continuous improvement process based on our learning, tape-out teams updated their workflow to use multiple sites across the WAN to get redundancy and scaling. In the latest generation, we have embarked on enabling sub-sites within a physical site to further limit the impact of a failure.

### Parallel Storage Advantages

The parallel storage system has delivered major advantages over our previous file servers.

- **Scalability.** In HPC-1 we were able to substitute one parallel server for every 10 conventional storage servers. This 10:1 consolidation ratio translated into huge cost savings because of reduced space requirements and energy consumption.



- **Performance.** For specific portions of the workflow, we achieved a 300-percent performance improvement in HPC-1 compared to the previous storage solution.
- **Volume size.** The maximum volume size increased by a factor of 16, from 400 GB to 6.4 TB, easily supporting our requirement for 3 TB-plus volumes in HPC-1.

### Scale-out Storage

We have shifted to a scale-out NAS solution that enables us to scale performance and capacity on demand without any disruption. Each scale-out node is configured with a set of filesystems and they are accessed via a single mount point. This node-based setup helps to ensure that workloads can be isolated to avoid impacting other workloads running on a different node in the cluster. We are able to utilize 4-TB disk-based shelves to deliver the required performance scaling while reducing our cost per usable tebibyte (TiB). Scaling for increased workload is obtained by adding nodes in pairs or by adding disk shelves. HPC has a total capacity of 160,000 cores. The HPC-5 storage environment scaled to support up to 50,000 distributed client CPU cores and volume sizes up to 10 TB to support the 10nm process generation.

### Backup

The HPC-1 requirements greatly exceeded the capabilities of our previous backup infrastructure. HPC-1 included disk volumes larger than 3 TB; to meet our service-level agreement, we needed to complete backup of these volumes within 24 hours. This required a single-stream throughput of at least 35 megabytes per second (MB/s).

At the time, this requirement was challenging because few available tape drives offered this level of performance. However, we identified a product that offered 120 MB/s raw performance per tape drive. After verifying performance, we coupled two of these drives with media servers running Linux\*, which enabled us to more easily use them with the parallel storage system.

When combined with the parallel storage system, this setup delivered aggregate read throughput of more than 200 MB/s. As a result, we were able to support 3 TB volumes without compromising our backup, archive, and restore service levels.

We have been able to continue scaling our HPC backup solution to support volume sizes up to 10 TB today.

## Compute Servers

Our tape-out environment includes thousands of servers that support highly compute-intensive applications. The increased demands of 45nm tape-out and beyond presented significant challenges in the following areas:

### Large-Memory Compute Systems

The largest tape-out jobs, such as design rule check (DRC) workloads, require servers with a very large RAM capacity. We also use these large-memory servers as master servers for distributed OPC applications.

The maximum certified memory capacity of servers in our pre-HPC tape-out environment was 128 GB. However, we knew that the increased complexity of 45nm processors would result in tape-out jobs that required up to 4x this memory capacity.



Moving to a higher-end system based on our existing architecture to support large memory capacity would have increased costs significantly. We therefore set a goal of implementing a system based on a modular architecture that could scale to meet future needs while meeting our aggressive cost objectives.

We identified a unique modular system based on NUMA architecture, capable of accommodating up to 32 Intel® Xeon® processors and 512 GB of RAM.

While this system provided the scalability we needed, the situation also created new challenges. There wasn't a Linux OS optimized for NUMA platforms, and neither the server nor the EDA applications were qualified for use in our environment.

We took a two-step approach: We first focused on deploying a 256 GB configuration to enable tape-out of the first 45nm processor, followed by a larger 512 GB system for tape-out of subsequent high-volume 45nm processors.

**256 GB Solution**

Our initial objective was to create a system based on four nodes, each with four processors and 64 GB of RAM, and compare performance with the previous solution. The architecture is shown in Figure 4.

This required close collaboration with the suppliers of the server hardware and the OS. We formed a joint system enablement team and worked intensively with a pre-release version of the OS to help ensure that it ran effectively on the system. We also worked with the OS supplier to conduct numerous performance and reliability tests.

As a next step, we worked closely with the EDA supplier to certify its memory-intensive DRC application on the new platform. Our efforts to resolve critical functionality, reliability, and performance issues achieved a remarkable result: We deployed the production system on the same day that the OS release was officially launched.

The new system successfully delivered substantial performance improvements and the ability to run bigger workloads. Large workloads ran 79 percent faster, compared with the previous server architecture.

**512 GB Solution**

Our objective was to enable an eight-node system with 32 CPUs and up to 512 GB of RAM, analyze the scalability and stability, and qualify the system in time to support tape-out of high-volume 45nm processors. We connected eight of the nodes illustrated in Figure 4; the interconnectivity is shown in Figure 5.

We evaluated this system when running DRC workloads consuming up to 512 GB of RAM. We tested multiple workloads in a variety of configurations, including single and multiple concurrent workloads using local and network file systems. We found that the system was able to scale to run these workloads with no performance degradation.

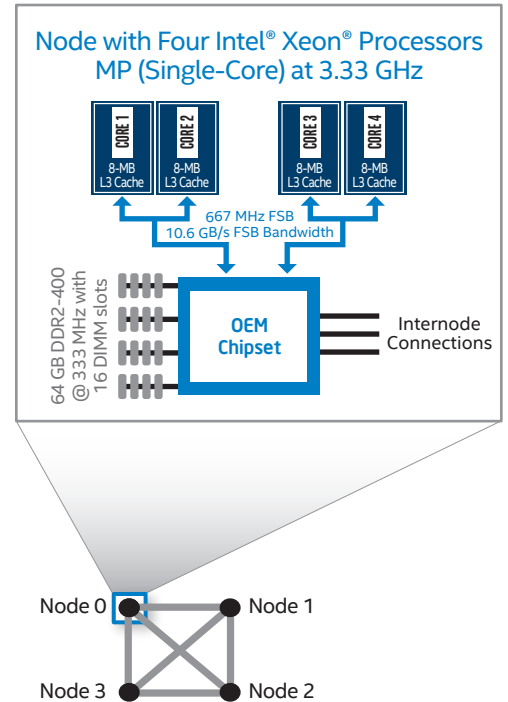


Figure 4. First-generation high-performance computing environment (HPC-1) large-memory system. Top: One node based on four single-core processors with 64 GB of RAM. Bottom: How four nodes interconnect to create a 256-GB system.

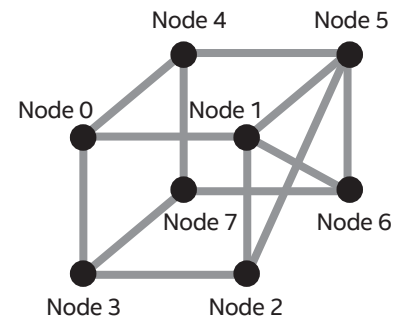


Figure 5. Interconnectivity for first-generation high-performance computing environment (HPC-1) large-memory compute server with eight nodes and 512 GB of RAM.

### HPC-1 Large-Memory Compute Server Refresh

When Intel® Xeon® processor 7100 series was released, with two cores per processor, we adopted these processors as standard. The overall system architecture remained the same, but each individual node now was equipped with additional cores and a larger L3 cache. An individual node is shown in Figure 6.

### HPC-2 Large-Memory Compute Server with 1 TB of RAM

For HPC-2, we took advantage of the introduction of the 45nm Intel® Xeon® processor 7400 series, with six cores per processor, to create a 96-core system with 1 TB of RAM. This consists of a four-node cluster in which each node has 256 GB of RAM and 24 processor cores. The architecture is shown in Figure 7.

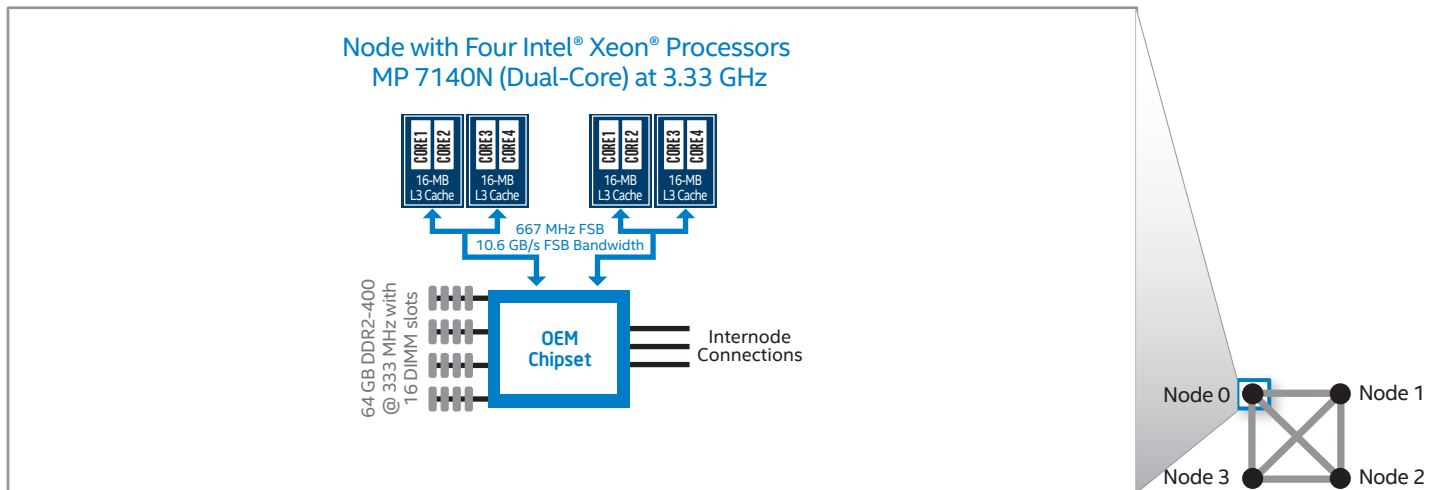


Figure 6. First-generation high-performance computing environment (HPC-1) large-memory refresh server. Left: One node based on four dual-core processors with 64 GB of RAM. Right: How four nodes interconnect to create a 256-GB system.

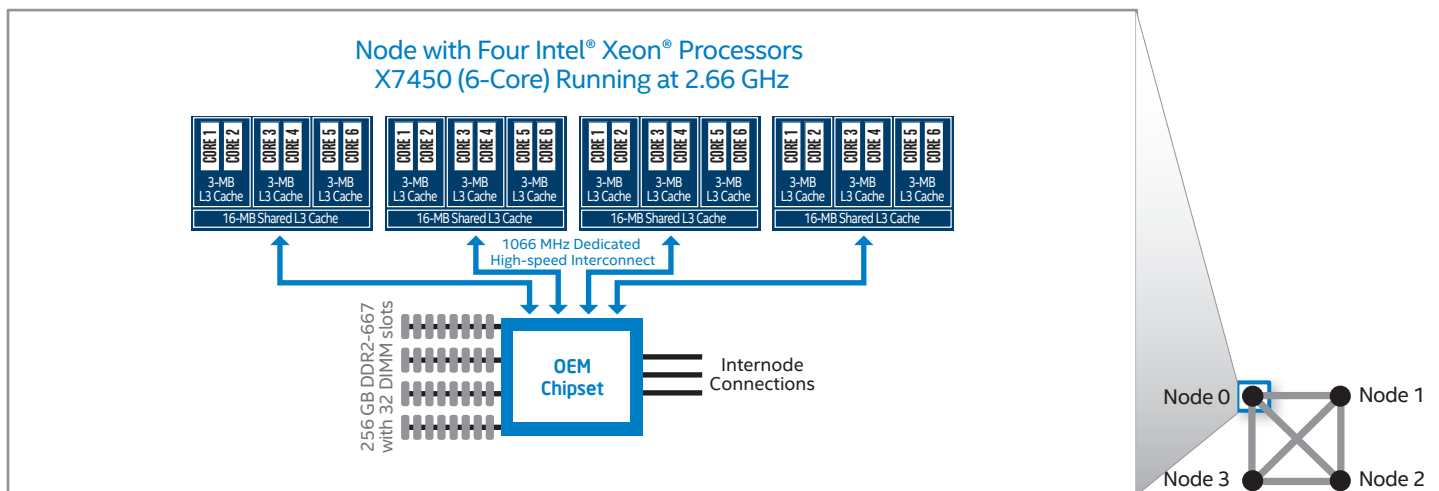


Figure 7. The Intel® Xeon® processor X7450-based node in a second-generation high-performance computing environment (HPC-2) large-memory compute server system. Left: One node based on four 6-core processors with 256 GB of RAM. Right: How four nodes interconnect to create a 1-TB system.

### HPC-3 and HPC-4 Large-Memory Compute Server with 2 TB of RAM

For HPC-3, we started by connecting two nodes of the Intel® Xeon® processor 7500 series with 8 cores per processor and 1 TB of RAM per node to create 64 cores per system with 2 TB of RAM. Later, we took advantage of the newer Intel® Xeon® processor E7-8800 product family 10-core processor platforms to create 80 cores per system with 2 TB of RAM for high-volume HPC-3 and HPC-4. For HPC-4, we adopted Intel® Solid-State Drives to enable fast swap on the Intel Xeon processor E7-8800 product family large-memory compute servers. The two architectures are shown in Figures 8 and 9, respectively.

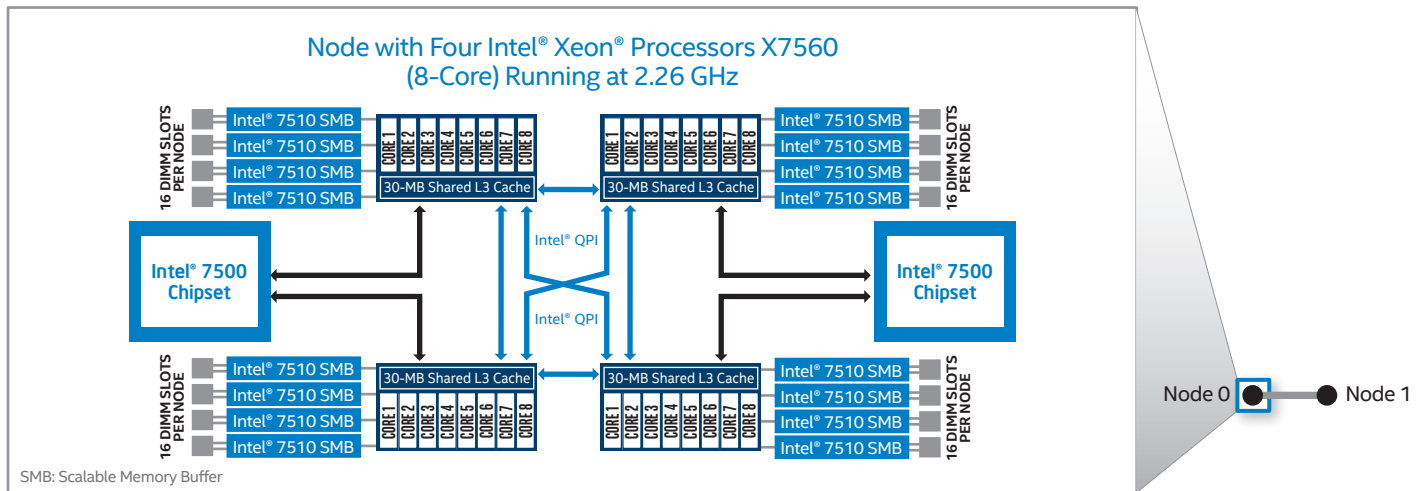


Figure 8. The Intel® Xeon® processor X7560-based node in a third-generation high-performance computing environment (HPC-3) large-memory compute server. Two such nodes are connected to form a single 2-TB system. Left: One node based on four 8-core processors with 1 TB of RAM. Right: Two nodes interconnect to create a 2-TB system.

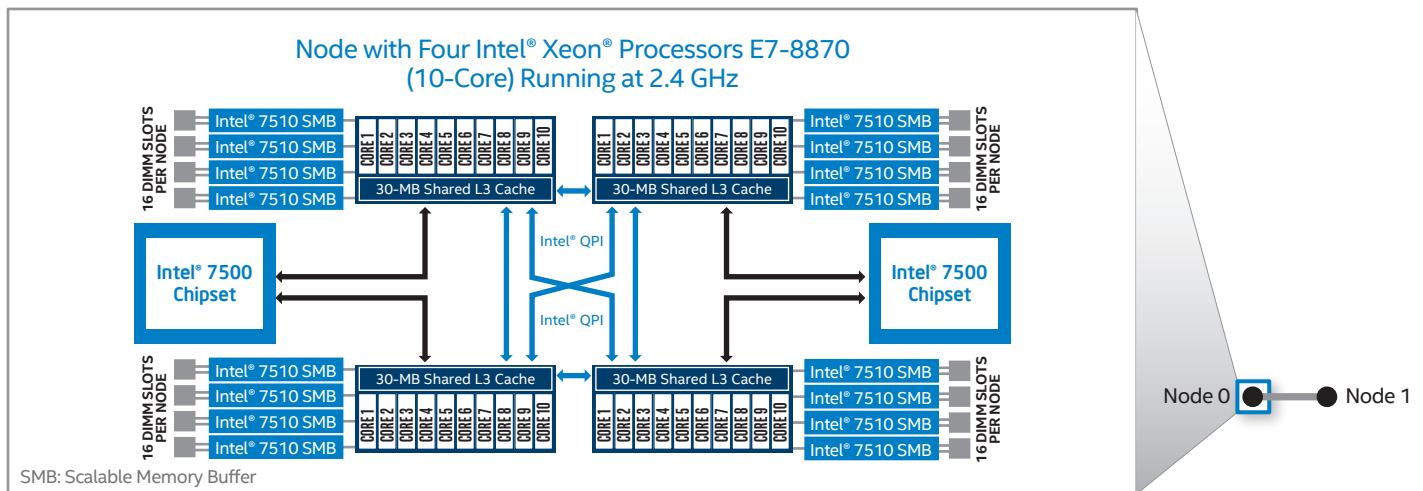


Figure 9. The Intel® Xeon® processor E7-8870-based node in a third-generation and fourth-generation high-performance computing environment (HPC-3 and HPC-4) large-memory compute server. Two such nodes are connected to form a single 2-TB system. Left: One node based on four 10-core processors with 1 TB of RAM. Right: Two nodes interconnect to create a 2-TB system.

### HPC-5 Large-Memory Compute Server with 3 TB and up to 6 TB of RAM

For HPC-5, we took advantage of the newer Intel® Xeon® processor E7-4800 v2 and E7-8800 v3 product families with up to 10-18 cores to create up to 72 cores per system. We increased the RAM capacity to either 3 TB with 32-GB DIMMs or 6 TB with 64-GB DIMMs. Also with the support of high-capacity DIMMs, large-memory footprints were achievable using single node. We continued to use Intel Solid-State Drives to enable fast swap on the Intel Xeon processor E7-4800 v2 and E7-8800 v3 product families' large-memory compute servers. The two architectures are shown in Figures 10 and 11, respectively.

All five generations of the large-memory compute servers are compared in Table 2, on the next page.

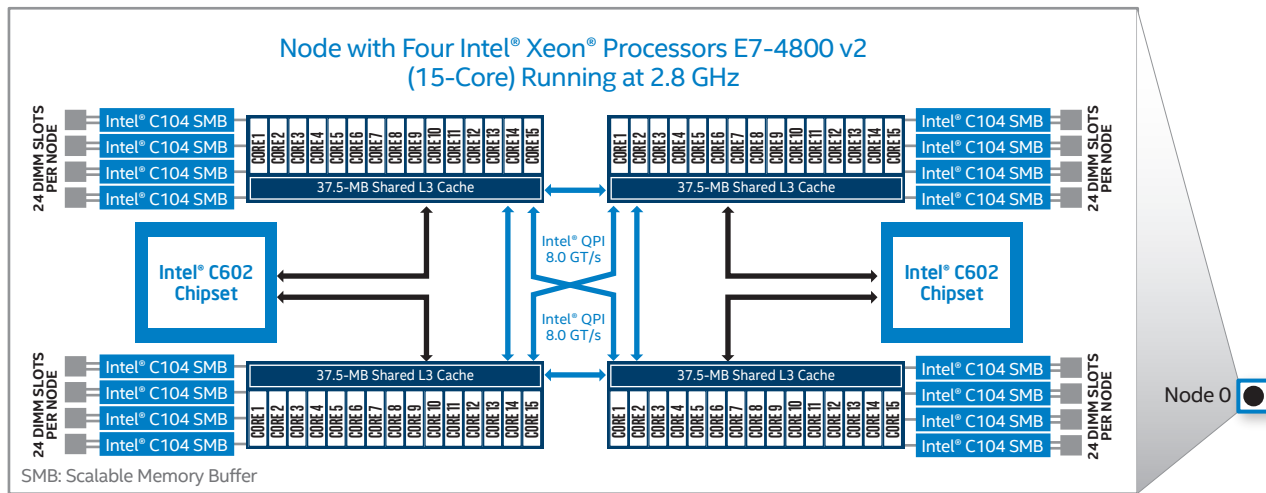


Figure 10. Intel® Xeon® processor E7-4800 v2-based node in a fifth-generation high-performance computing environment (HPC-5) large-memory compute server. HPC-5 used a processor-based, 15-core, single-node server with 3 TB of RAM.

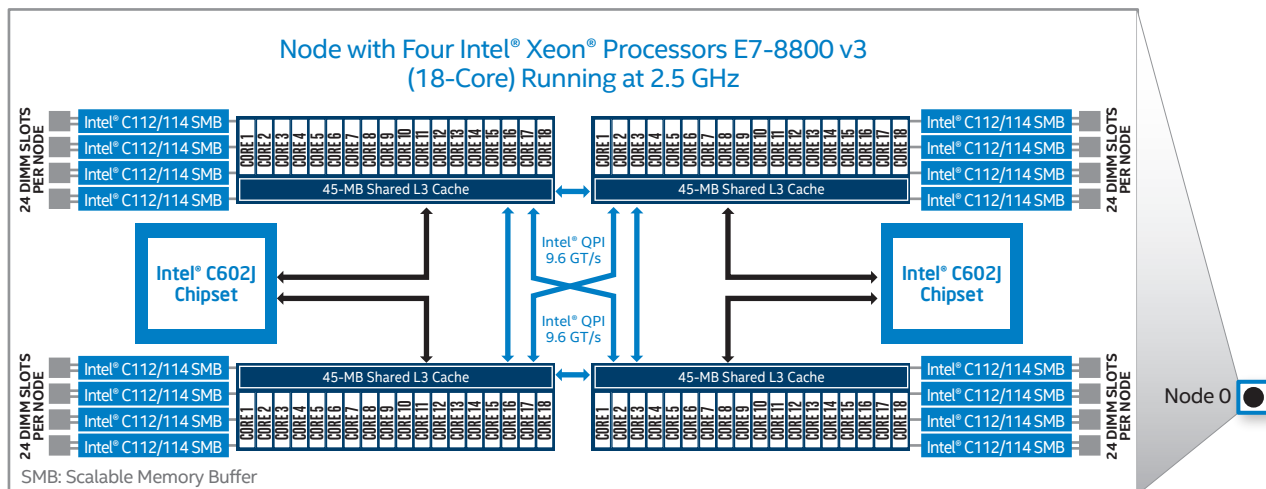


Figure 11. Intel® Xeon® processor E7-8800 v3-based node in a fifth-generation high-performance computing environment (HPC-5) large-memory compute server. HPC-5 used a processor-based, 10-18-core, single-node server with up to 6 TB of RAM.

Table 2. High-Performance Computing (HPC) Environment Comparison of Large-Memory Compute Servers HPC-1 through HPC-5

	GENERATIONS				
	HPC-1	HPC-2	HPC-3	HPC-4	HPC-5
Total CPU Cores	32 or 64	96	64 or 80	80	40-72
Memory Capacity	512 GB	1 TB	2 TB	2 TB	3 TB or 6 TB
Intel® SSD Fast Swap	—	—	—	Yes	Yes
Data Center Space Needed (Rack Units)	24	16	8	8	4
Power Consumed	7.3 kW	3.6 kW	2.52 kW or 2.34 kW	2.34 kW	1.49 kW

GB – gigabyte; kW – kilowatt; SSD – solid-state drive; TB – terabyte

### Batch Compute Servers

Compute-intensive tape-out jobs such as OPC are handled by large clusters of batch compute servers operating in parallel in a master-slave configuration. To illustrate the scale of the challenge, there may be as many as 40,000 OPC jobs executing concurrently on thousands of servers.

We achieved major throughput performance improvements by taking advantage of multi-core Intel Xeon processors as they became available. Our pre-HPC environment relied on single-core processors, but we subsequently moved to dual-core and then quad-core, six-core, eight-core, ten-core, and twelve-core processors.

Our tape-out workload results provided real-world proof of a key theoretical advantage of multi-core processors: that performance scales with the number of cores within an HPC cluster.

Servers based on Intel Xeon processors with four cores showed a consistent ability to run twice as many jobs as servers with prior-generation dual-core processors and delivered faster runtimes with a relative throughput of 4.8x compared to older generation single-core processors. The relative throughput scaling has increased to 37.66x by HPC-5 with twelve-core processors.

The performance benefits achieved with faster Intel Xeon processor-based batch compute servers in HPC-1 translated directly into a reduction in data center space and energy requirements.

As new Intel® server processors are released, we have continued to incorporate servers based on these processors into our environment. This delivers continuing increases in performance for key applications such as OPC and simulation, as shown in Figure 12, on the next page.

Our tape-out also utilized high-frequency Intel Xeon processors with lower core count per system in their batch computing environment. This has shown a >1.35x increase in core-to-core performance, improving time-critical stages of the tape-out computing needs.

We achieved major throughput performance improvements by taking advantage of multi-core Intel Xeon processors as they became available.

### Intel® Architecture Performance Improvement for Optical Proximity Correction (OPC)

Higher is Better

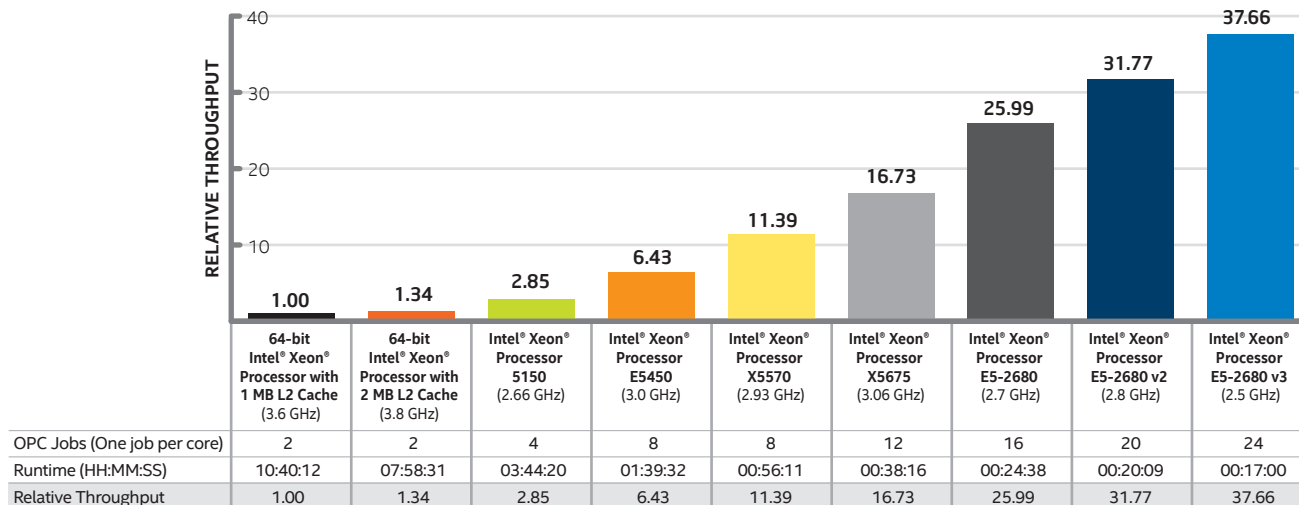


Figure 12. Servers based on successive generations of multicore Intel® Xeon® processors continue to deliver improvements in batch computing performance.

## Network

By carefully characterizing data transfer requirements, we determined the need to increase bandwidth and provide high availability across the tape-out environment. We made the following upgrades:

- For HPC-1: all master and large-memory compute servers to at least 2x 1 gigabit per second (Gb/s) network connection with switch-level failover capabilities; all slave servers to at least 100 megabits per second (Mb/s).
- By HPC-4: all master and large-memory compute servers to 2x 10 Gb/s; all slave servers to 1 Gb/s.

We provided 6x 1 Gb/s uplinks per rack in HPC-1 and 2x 10 Gb/s uplinks per rack by HPC-4. Also, we configured the two uplinks to connect to two different switches and virtual LANs (VLANs) for redundancy in case of link or upstream switch failure.

## Batch Clustering: Job Scheduler Improvements

Tape-out involves scheduling thousands of simultaneous OPC batch jobs as efficiently as possible. Heavy job loading exposed quality issues in the batch job scheduler, resulting in a higher level of job failures and lower server utilization.

We devised a systematic test method based on synthetic jobs that did not generate load on the CPU. This enabled us to analyze and stress test the job scheduler code—scheduling up to 11,000 production machines while the machines were still being used for regular production work. As a result, we were able to execute several million test jobs per day.

This method was key to developing an improved scheduler as well as to detect and fix bugs, because it allowed us to rapidly test combinations of hardware and OS scheduler configurations.

Our improved scheduler cut in half the time required for job submission and scheduling. It also supported three independent job queues, resulting in a 13.5x increase in the total number of jobs supported by our tape-out resources.

## EDA Application License Servers

EDA application license server performance was a factor constraining the growth of our tape-out environment. Random job failures occurred when the license servers were heavily loaded, resulting in an inability to check out more licenses.

As when optimizing the job scheduler, testability was a key challenge. It was impractical to extensively test the license servers using the actual EDA application, because this would have required the dedicated use of more than 5,000 production server CPUs over several days.

We overcame this obstacle by working with suppliers to develop a methodology for testing simultaneous license checkout of 1,000 keys per second from a single machine—while running regular production jobs. This enabled us to stress-test the license servers and validate new software and configuration combinations.

This approach led to the discovery of a fundamental bug in the license server application that limited scalability and enabled suppliers to fix it before it impacted our growing production environment.

We used the same method to demonstrate to our EDA application supplier that license servers based on Intel® architecture were stable and more scalable than the RISC-based used in our pre-HPC production environment. The move to Intel® architecture-based license servers meant that our design and tape-out environment was completely based on Intel architecture.

## Enterprise Linux\* OS

To improve the stability of batch computing, we standardized on the same enterprise Linux OS on all our HPC large-memory and batch computing servers. As we took advantage of new hardware platforms, we worked with the OS supplier to enhance and optimize the OS to support new hardware features. We also worked with the OS supplier to resolve bugs and to help ensure interoperability between new and existing platforms.

## Application and Platform Tuning

To take full advantage of multicore platforms, we have optimized BIOS settings for processors, memory, and hard drive operation modes to achieve a further 20-percent performance improvement. We also periodically performed internal stress tests to help ensure that the efficiency of our HPC cluster is comparable with top-ranked GbE supercomputing clusters in the Top500\*.

---

To take full advantage of multicore platforms, we have optimized BIOS settings for processors, memory, and hard drive operation modes to achieve a further 20-percent performance improvement.

---



# HPC Benefits

The use of HPC-1 to enable tape-out of Intel's breakthrough 45nm processors delivered significant value to Intel. Financial analysis showed that HPC-1 alone delivered net present value (NPV) of USD 44.72 million, of which USD 22.68 million was directly attributable to the first generation of the parallel storage solution and USD 16.64 million to the large-memory compute servers. Batch compute server improvements reduced requirements for data center space, power, and cooling, resulting in USD 5.4 million NPV.

HPC-2, HPC-3, HPC-4, and HPC-5 have continued to deliver substantial increases in scalability and performance, as shown in Table 3.

In addition to providing the major increases in compute capacity required for new processor generations, HPC has dramatically improved the stability of our tape-out environment. The number of issues impacting tape-out declined sharply after the implementation of HPC-1, and this improvement has been sustained even as the environment has supported continuous growth in demand. As shown in Figure 13, since deployment, HPC has supported more than a 90.21x increase in demand, with a 64.4x increase in stability.

## Intel® Tape-Out Computing Metrics

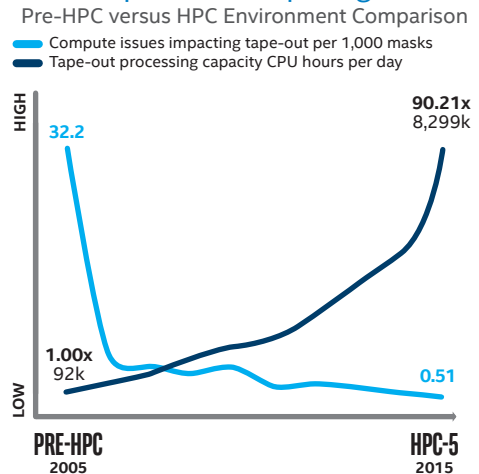


Figure 13. High-performance computing (HPC) has shown increased stability even as demand has increased. Note: Tape-out processing capacity normalized to Intel® Xeon® processor 5150

Table 3. Summary of HPC Generational Performance Improvements Note: Generation scaling improvements shown in parenthesis.

	HPC-1 Optimize for 45nm Support >=65nm	HPC-2 Optimize for 32nm Support >=45nm	HPC-3 Optimize for 22nm Support >=32nm	HPC-4 Optimize for 14nm Support >=22nm	HPC-5 Optimize for 10nm Support >=14nm
<b>BATCH CLUSTERING.</b> Stability, Scalability, Features					
Systems per Pool	8,500 (1.15-2x)	11,000 (1.15-2x)	11,000 (1x)	11,000 (1x)	11,000 (1x)
Virtual Sites/Pools	2	2	2	2	3
Jobs per Pool	20,000+ (1.5x)	30,000+ (1.5x)	120,000+ (4x)	180,000+ (1.5x)	180,000+ (1x)
<b>STORAGE AND BACKUP.</b> Scalability, Performance, I/O Throughput					
Volume Size	3.2 TB (8x)	6.4 TB (2x)	6.4 TB (1x)	11 TB (1.72x)	11 TB (1x)
Single-Stream Performance	70 MB/s (1x)	160 MB/s (2.3x)	240 MB/s (1.5x)	240 MB/s (1x)	TBD
Hardware and Software Capacity	Parallel Storage-Generation 1 100 TB (8x)	Parallel Storage-Generation 2 100 TB (1x)	Parallel Storage-Generation 3 100 TB (1x, SSD)	Parallel Storage-Generation 4 360 TB (3.6x)	Scale out Storage-Generation 1 1,532 TB (4.25x)
<b>NETWORK.</b> Scalability, Stability					
Storage	40 Gb/s (10x)	40 Gb/s (1x)	100 Gb/s (2.5x)	300 Gb/s (3x)	40 Gb/s per scale out node
Master	1 Gb/s (10x)	2x 1 Gb/s (1x, Redundancy)	2x 10 Gb/s (10x, Redundancy)	2x 10 Gb/s (1x, Redundancy)	2x 10 Gb/s (1x, Redundancy)
Slave	100 Mb/s (1x)	100 Mb/s (1x)	100 Mb/s-1 Gb/s (1x-10x)	1 Gb/s (1x)	1 Gb/s (1x)
<b>COMPUTE SERVERS.</b> Optimized for Performance, Throughput, Capacity, Power, and Data Center Space					
Large RAM Server	512 GB (4x)	1 TB (2x)	2 TB (2x)	2 TB (1x, SSD)	3 TB, 6 TB (1.5-3x, SSD)
Batch Node	2-Socket, Dual-Core 16 GB	2-Socket, Quad-Core 32 GB	2-Socket, 4-8 Cores 72-128 GB	2-Socket, 8+ Cores 128-256 GB	1 and 2-Socket, 4+ Cores 32-256 GB
<b>OPERATING SYSTEM.</b> New Hardware Feature Support, Scalability, Stability, Performance					
Enterprise Feature	Stable, Intersystem NUMA Support	Multi-Core Optimized	Power, Performance Optimized	Power, Performance Optimized	Power, Performance Optimized

Gb/s – gigabits per second; GB – gigabyte; Mb/s – megabits per second; MB/s – megabytes per second; NUMA – non-uniform memory access; SSD – solid-state drive; TB – terabyte

## Key Learnings and Future Plans

The success of HPC was based on several key factors.

- **Alignment of technology with business requirements.** In specifying the HPC solution, we carefully aligned business and technical requirements, resulting in a system that delivered the scalability required to support 45nm and beyond processor tape-out. We are continuing to use this model to align successive HPC and process technology generations.
- **Informed risk-taking.** To optimize solutions for HPC, we needed to take risks, such as our pioneering decisions to use our parallel storage system and the modular large-memory compute servers. Implementing these solutions required significant ecosystem development. Our team understood that there was a significant risk, with concerns about supplier maturity and the viability of the solution in production use, yet we strongly believed the system would deliver great rewards to Intel. The fact that these solutions worked and enabled 45nm processor—and newer generations, up to 14nm processor—tape-out demonstrated that the risk level was appropriate.
- **Governance.** We adopted a holistic view of HPC capabilities and created a clear computing roadmap. Disciplined governance then helped ensure that we executed according to this roadmap. Intel IT and business groups acted as a single team with collective responsibility; a joint manufacturing and IT committee reviewed and approved computing recommendations.
- **Hyperscale, agility, and business continuity.** With increasing complexity of chip design, Intel has adopted hyperscale computing to address the exponential growth of data volume and demand from tape-out workloads. This is achieved using agile cost-efficient multi-site concepts within a data center accommodating thousands of servers and operating at the optimal power usage effectiveness.

---

We are currently developing HPC-6, which is planned to support the tape-out of 7nm processors.

---

We are currently developing the sixth-generation HPC (HPC-6), which is planned to support the tape-out of 7nm processors. As with previous generations, we expect to optimize the throughput of 7nm tape-out applications with significant, balanced improvements across all HPC components. We are constantly improving our hyperscale data centers to meet these workload demands. Additionally, this includes major performance improvements in the areas of storage, compute servers, batch clustering, and network bandwidth. These concepts and learnings are also being applied to other phases of silicon design.

## Conclusion

Our pioneering HPC approach to silicon design enabled tape-out of the industry's first 45nm processors and numerous follow-on products.

Delivering this solution required replacing our old computing model with an innovative approach aligned with the requirements of Intel process technology generations. Intel's manufacturing group recognized two components of our environment—the parallel storage solution and large-memory Intel® Xeon® processor E7-based NUMA systems—as pillars supporting the successful completion of the first 45nm processors. Intel has taped out several silicon products with HPC-1 alone, delivering ROI of USD 44.72 million and reducing tape-out time from 25 to less than 10 days.<sup>1</sup> We are continuing to support the increasing design complexity as Intel process technology advances to meet scalability requirements, while keeping tape-out time to less than 10 days.

For more information on Intel IT best practices, visit [www.intel.com/IT](http://www.intel.com/IT).

Receive objective and personalized advice from unbiased professionals at [advisors.intel.com](http://advisors.intel.com). Fill out a simple form and one of our experienced experts will contact you within 5 business days.

### IT@Intel

We connect IT professionals with their IT peers inside Intel. Our IT department solves some of today's most demanding and complex technology issues, and we want to share these lessons directly with our fellow IT professionals in an open peer-to-peer forum.

Our goal is simple: improve efficiency throughout the organization and enhance the business value of IT investments.

Follow us and join the conversation:

- [Twitter](#)
- [#IntelIT](#)
- [LinkedIn](#)
- [IT Center Community](#)

Visit us today at [intel.com/IT](http://intel.com/IT) or contact your local Intel representative if you would like to learn more.



<sup>1</sup> Tape-out time was reduced from 25 days for the first 65nm process technology-based microprocessor in a non-HPC compute environment to less than 10 days through an HPC-enabled environment. Financial analysis showed that HPC-1 alone delivered a net present value (NPV) of USD 44.72 million.

Software and workloads used in performance tests may have been optimized for performance only on Intel® microprocessors. Performance tests, such as SYSmark\* and MobileMark\*, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products.

Configurations: System configurations, SSD configurations and performance tests conducted are discussed in detail within the body of this paper. For more information go to [intel.com/performance](http://intel.com/performance).

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Performance varies depending on system configuration. Check with your system manufacturer or retailer or learn more at [intel.com](http://intel.com).

THE INFORMATION PROVIDED IN THIS PAPER IS INTENDED TO BE GENERAL IN NATURE AND IS NOT SPECIFIC GUIDANCE. RECOMMENDATIONS (INCLUDING POTENTIAL COST SAVINGS) ARE BASED UPON INTEL'S EXPERIENCE AND ARE ESTIMATES ONLY. INTEL DOES NOT GUARANTEE OR WARRANT OTHERS WILL OBTAIN SIMILAR RESULTS.

INFORMATION IN THIS DOCUMENT IS PROVIDED IN CONNECTION WITH INTEL PRODUCTS AND SERVICES. NO LICENSE, EXPRESS OR IMPLIED, BY ESTOPPEL OR OTHERWISE, TO ANY INTELLECTUAL PROPERTY RIGHTS IS GRANTED BY THIS DOCUMENT. EXCEPT AS PROVIDED IN INTEL'S TERMS AND CONDITIONS OF SALE FOR SUCH PRODUCTS, INTEL ASSUMES NO LIABILITY WHATSOEVER AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY, RELATING TO SALE AND/OR USE OF INTEL PRODUCTS AND SERVICES INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT OR OTHER INTELLECTUAL PROPERTY RIGHT.

Intel, the Intel logo, Celeron, and Xeon are trademarks of Intel Corporation in the U.S. and other countries.

\*Other names and brands may be claimed as the property of others.

Copyright © 2015 Intel Corporation. All rights reserved.