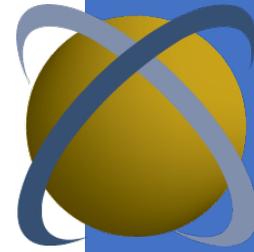# Intersect360 Research White Paper:
# SPEEDING UP HPC INNOVATION:
# THE AMD EPYC 7371 PROCESSOR

## EXECUTIVE SUMMARY

What does it mean for a computer to be faster, or more powerful? "Price/performance" is more complex than it may seem. "Performance" has nuanced components, and "price" is more than the acquisition cost of a system. Power consumption and management costs are both important and influenced by the choice of system. And across the HPC industry, the biggest spending category outside of the computational system is software.

Engineering-driven use cases tend to rely heavily on applications provided by independent software vendors (ISVs). In such cases, software licensing costs may overwhelm server acquisition costs. Another important spending consideration is personnel. Ultimately, "price/performance" may be best thought of in an inverted sense of performance-over-price: how much productivity a company can get from its investment in its skilled people and the tools it provides them.

In some cases, applications perform better with relatively less forced parallelism. Furthermore, many ISV applications are licensed according to the number of cores they are run on. In these cases, there is a direct cost of additional cores, and running faster on fewer cores has an immediate influence on the price/performance equation.

In November 2018, AMD unveiled the newest member of the AMD EPYC™ processor line. Branded as the AMD EPYC™ 7371 processor, the new processor is aimed directly at high-performance markets, with 16 cores and a base frequency of 3.1GHz.

Within the spectrum of offerings, the EPYC 7371 targets the segments of HPC that benefit from higher clock rates. Like most of its brethren, the 7371 is available in either single-socket or dual-socket configurations, with eight DDR4 memory channels, 128 PCIe Gen3 lanes, and 64MB of L3 cache. The higher clock frequency serves applications that exhibit less parallelism or that have higher licensing costs for more cores.

AMD is taking a consultative approach to the market, entrenched in domain-specific expertise. The benchmarks it highlights in its solution brief demonstrate that the company understands HPC workloads, and more importantly, that it is willing to invest in comparative testing to help guide buyers to optimal solutions. With its high-frequency options, AMD is betting that a fast processor can get off to a fast start.

## MARKET DYNAMICS

### Different Ways to Go Fast

Despite appearances, High Performance Computing (HPC) isn't really about the race for faster computers. It is driven by the relentless pursuit of innovation, across science, engineering, and other fields of discovery. Once one problem is solved, it immediately raises more questions. Variables are added; fidelity is increased; degrees of freedom are expanded. The refined models better reflect reality, and new achievements are unlocked. The next insight is the goal.

Each new level of innovation is harder to reach than the last, building on everything learned so far. We need to extend our capabilities just that much further, over the next horizon. It is this link between achievement and capability that pushes the boundaries of computing. It's not the tool; it's what you can build with it.

What does it mean for a computer to be faster, or more powerful? We can read a data sheet to count computational cores or to measure bandwidth. The test that matters is how much innovation is offered. Definitions of performance vary from one industry to the next, one customer to the next. Careful evaluation is required to determine which computational resource will best suit a particular workflow or combination of workflows.

### Cost-Efficient Performance

Naturally, the performance of a computer can't be the only consideration, because there are costs to consider. "Price/performance" is a commonly cited metric in HPC to describe computational efficiency, but this seemingly simple discussion is more complex than it may seem. For one, "performance" itself has nuanced components, as described above. For another, "price" is more than the acquisition cost of a system. Power consumption and management costs are both important and influenced by the choice of system, but it doesn't stop there.

Across the HPC industry, the biggest spending category outside of the computational system is software. Out of $36 billion in worldwide HPC spending, users spent $8.5 billion on software in 2018.[1] The proportion of hardware-to-software spending varies widely by industry. Engineering-driven use cases, such as in manufacturing, chemical engineering, or electronic design automation (EDA), tend to rely heavily on applications provided by independent software vendors (ISVs). In such cases, software licensing costs may overwhelm server acquisition costs.

Another important spending consideration is personnel. After all, it is the end users of the system whose effectiveness is being measured. In an engineering-driven company, HPC is meant to make the engineers more productive. Ultimately, "price/performance" may be best

*In engineering-driven use cases, software licensing costs may overwhelm server acquisition costs.*

---

[1] Intersect360 Research, "Worldwide High Performance Computing 2018 Total Market Model and 2019–2023 Forecast: Products and Services," May 2019.

thought of in an inverted sense of performance-over-price: how much productivity a company can get from its investment in its skilled people and the tools it provides them.

## Cores versus Frequency

An HPC application walks into a bar and orders a beer. It finishes the beer, but never being satisfied, it wants another beer. Only now it orders more beer, faster. Soon its friends arrive, and they want beer too. The more beer they consume, the more beer they demand, until the bartender can't keep up with the queue of beer orders. What's a tavern owner to do?

One solution is to add more beer taps. This allows multiple bartenders to fill mugs simultaneously. But it doesn't do much for the single barman filling a large pitcher, short of rigging a system by which multiple taps feed the same vessel. Moreover, all those beer taps take up room, and it takes time to traverse the length of the bar from one tap to another. Another solution is to increase the flow of each tap. Now each mug, stein, or pitcher gets filled faster. As the bar scales, it likely finds that some combination of "more" and "faster" is desired.

This metaphor illustrates the tradeoffs in microprocessor development since the turn of the century. Previously, the way to make a processor faster was to increase its clock frequency and thereby the number of calculations it could do each second. This is analogous to the higher-volume beer tap. But as process limitations made it impractical to continue the clock race, processor vendors turned their attention to putting multiple computing cores into each processor. This is analogous to adding more beer taps. Competing processors had two cores, then four, then eight, and today it is rare to see one with less than 12 in an HPC environment.

Which approach is better? This depends on workflow context, but as with our bar example, most users will find some combination of these approaches provides the best compromise between more cores and faster cores. Having more cores tends to be better for workloads with many, smaller jobs that can conveniently be run in parallel. Having faster cores is better when it is important to speed up one single job (filling the large pitcher), or when there is an additional cost to maintaining the additional cores.

Many HPC domains lean toward this latter category. In some cases, applications perform better with relatively less forced parallelism—that is, when the cost of decomposing an application into subtasks isn't worth the speedup provided by the parallel computation. Furthermore, many ISV applications are licensed according to the number of cores they are run on. (There is logic to this, since the additional cores put the burden of parallelization on the ISV.) In these cases, there is a direct cost of additional cores, and running faster on fewer cores has an immediate influence on the price/performance equation.

# AMD EPYC 7371 PROCESSORS: HPC WITH MORE FREQUENCY

## Targeting HPC

In November 2018, AMD unveiled details of its second-generation EPYC processor line, codenamed "Rome," the first data center microprocessors based on a 7-nanometer (7nm) process technology, following onto the first-generation EPYC "Naples" processors introduced in early 2017. The new processors aimed directly at high-performance markets, with up to 64 cores based on AMD's Zen 2 microarchitecture.

With expected increased instructions per cycle and core count, the new EPYC processors are projected to offer double the raw performance per socket of the previous generation, with up to four times the performance on floating-point arithmetic. Talking of performance at the "Next Horizons" launch event, AMD CEO Lisa Su declared, "This is our space. This is where we lead."

In advance of the second-generation AMD EPYC processors reaching the market, AMD has introduced a product line extension to the first-generation EPYC 7xx1 series. The EPYC 7xx1 series is characterized by 128 PCIe lanes and up to 341 GB/sec of memory bandwidth. It is complemented by a 7nm GPU, the AMD Radeon Instinct™ MI60, which began shipping soon after the Next Horizons launch. The CPU-to-CPU and GPU-to-GPU connections are over Infinity Fabric™, AMD's technology for coherent component connections.

Now in general availability, the AMD EPYC 7xx1 line comes in a range of configurations, both single-socket and dual-socket, with eight to 32 cores, all targeting high-performance workloads. But at SC18, the International Conference for High Performance Computing, Networking, Storage and Analysis, one week after the Next Horizons launch in November 2018, much of the HPC talk focused in on one particular offering: the AMD EPYC 7371.

## High-frequency AMD EPYC 7371

The entire EPYC 7xx1 line is focused on high-performance workloads. Within the spectrum of offerings, the EPYC 7371 targets the segments of HPC that benefit from higher clock rates. While the other 7xx1 series processors have base clock rates from 2.00GHz to 2.50GHz (with boost frequencies ranging from 2.55GHz to 3.20GHz), the 7371 starts at a base frequency of 3.10GHz, with an "all cores" boost to 3.60GHz, or a burst on half the cores to 3.80GHz.

Like most of its brethren, the 7371 is available in either single-socket or dual-socket configurations, with eight DDR4 memory channels, 128 PCIe Gen3 lanes, and 64MB of L3 cache. The higher clock frequency serves applications that exhibit less parallelism or that have higher licensing costs for more cores.
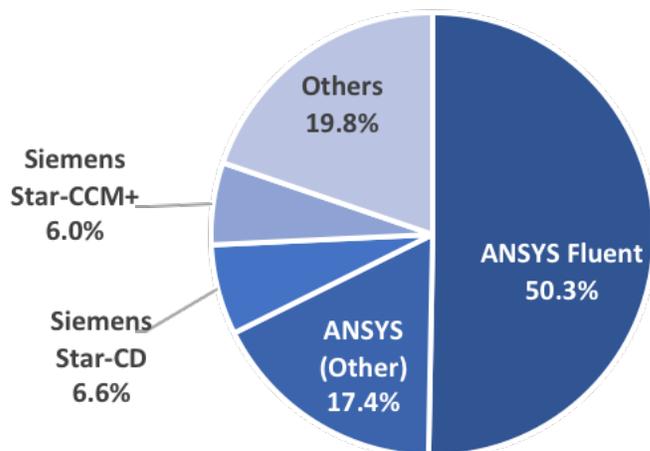
An AMD solution brief[2] compares the 7371 to other AMD EPYC 7xx1 processors on some of the most commonly used HPC applications in computational fluid dynamics (CFD),

---

[2] https://www.amd.com/system/files/documents/amd_epyc_7371_processors_the_right_choice_for_hpc.pdf.

engineering analysis, molecular dynamics, and weather simulation. For ANSYS Fluent, the top CFD package by usage (see figure below), the paper also compares EPYC 7371 performance to that of Intel Xeon Gold 6148 (20 core, 2.4GHz), citing "up to 5% average lead in overall system performance across all benchmarks," with "up to 49%" and "an average of up to 32%" gain in per-core performance, which is "critically important to maximize software license investments."

## Top Named ISV-Provided CFD Software at HPC Sites, by Mentions
### (Open-source and in-house applications excluded)
Intersect360 Research, 2019



The paper continues, giving benchmark comparisons of the EPYC 7371 to its two closest AMD counterparts, the EPYC™ 7351 (16 cores, 2.4GHz) and the EPYC™ 7451 (24 cores, 2.3GHz) processors, showing the effects of adding cores versus adding frequency. These comparisons span commonly used HPC applications: WRF (weather forecasting), NAMD (molecular dynamics), LSTC LS-DYNA and EMI PAM-CRASH (finite element analysis), Altair Radioss (structural analysis), and STAR-CCM+ (CFD and engineering simulation), all among the top-used applications in their categories.

## INTERSECT360 RESEARCH ANALYSIS

In 2006, AMD was the number-one processor vendor to the HPC community,[3] riding the success of its Opteron™ processor line, the first processors with 64-bit arithmetic extensions. But HPC users are fickle and will change vendors quickly if they can find a better-performing alternative elsewhere. AMD soon lost its lead, and Intel has dominated in market share for the past decade.
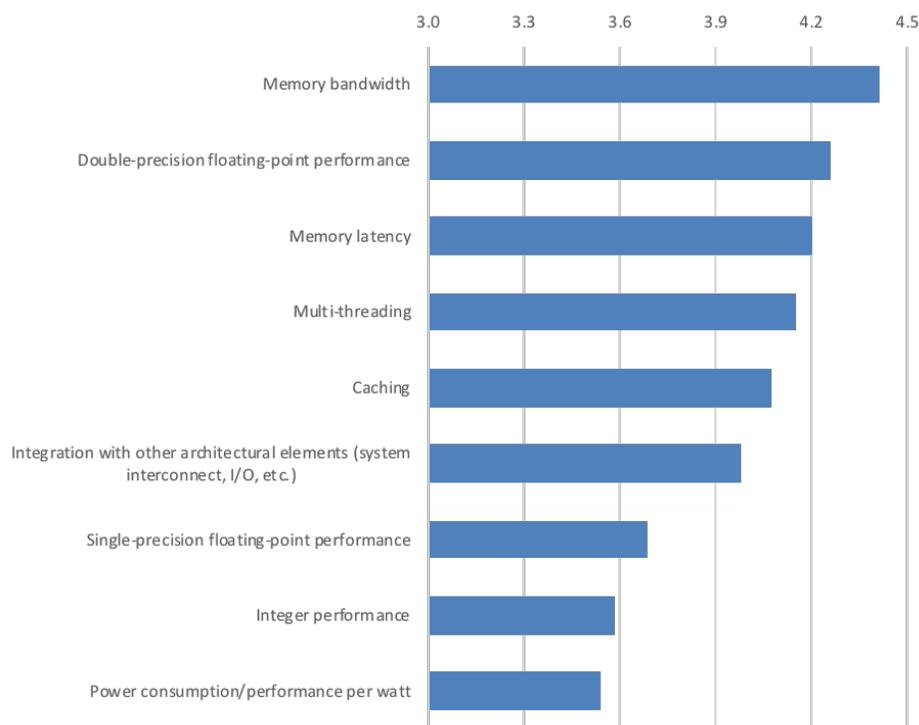
---

[3] Intersect360 Research (as Tabor Research) HPC User Site Census survey data, 2007.

That fickleness could benefit AMD now. In a 2016 study,[4] HPC users identified memory bandwidth as the most important technical criterion in evaluating new processor architectures (see figure below). Moreover, as buyers explore alternatives to Intel processors, AMD has an advantage in that it also provides an x86 architecture, alleviating concerns in application porting and support. In that same study, users cited "availability of compatible software stack (operating system, libraries, compilers, etc.)" as the most important software consideration.

## Importance of Technical Features in Evaluating HPC Processors

Average survey score shown. 1 = not important; 5 = extremely important.
Intersect360 Research special study, 2016



Naturally, many HPC users are also leveraging GPU computing as part of their HPC environments. Here NVIDIA has been dominant in establishing and leading the GPU computing category, but again AMD may have an ace in the hole. Intel and NVIDIA have become natural competitors, and while this hasn't prevented server deployments with Intel CPUs and NVIDIA GPUs, it does stand in the way of cooperation.

With both EPYC and Radeon™ based on the 7nm process and CPU-to-CPU and GPU-to-GPU connections supported over Infinity Fabric, a natural next step will be for AMD to connect CPU-to-GPU using the high-bandwidth, low-latency coherent connection. Notably, this

---

[4] Intersect360 Research special study, "Processing Elements for HPC," 2016.

configuration of EPYC-plus-Radeon-over-Infinity Fabric was selected by the U.S. Department of Energy (DOE) for the upcoming Exascale supercomputer at Oak Ridge National Laboratory (ORNL), to be installed as early as 2021 as one of the first Exascale systems in the U.S., and in the world. Named "Frontier," this Cray Shasta supercomputer is expected to top 1.5 Exaflops[5] of peak computing performance.[6]

One supercomputer, no matter how powerful, isn't enough to make a market, but there is an undeniable echo effect of these large-scale national lab systems to the industry at large, as companies evaluate their next round of HPC purchases. Over half of all HPC systems by revenue are consumed by industry, many of which are in the engineering-driven fields targeted by the AMD EPYC 7371. The manufacturing sector alone spent $4.8 billion on HPC products and services in 2018, 13.2% of the overall market. AMD is going after a substantial piece of that pie.

AMD is taking a consultative approach to the market, entrenched in domain-specific expertise. The benchmarks it highlights in its solution brief demonstrate that the company understands HPC workloads, and more importantly, that it is willing to invest in comparative testing to help guide buyers to optimal solutions. This has historically been a comparative weakness of Intel's, which has tended to let its products compete with each other naturally, with fewer resources to recommend one product over another.

The majority of HPC system vendors have already adopted AMD EPYC into their server lines, and EPYC (including the 7371) is now generally available to HPC users. Ultimately AMD has to prove its mettle not only by hitting its roadmap targets (which it has done consistently since EPYC's introduction), but by delivering in the market with real-world HPC applications. With its high-frequency options, AMD is betting that a fast processor can get off to a fast start.

---

[5] One Exaflop = $10^{18}$ floating point calculations per second. $10^{18}$ = one quintillion = one billion billion = 1,000,000,000,000,000,000.

[6] https://www.ornl.gov/news/us-department-energy-and-cray-deliver-record-setting-frontier-supercomputer-ornl.