

INFORMATION BRIEF

# Cloud Native Computing

Why data center operators should care about 128-core processors



## Introduction

The increased demand for data centers continues as the world races towards a more digital future. This is putting pressure on the finite supplies of land and power in localities where data centers are sited. A 2021 report from the International Energy Agency stated that data centers were responsible for 1% of the world's electricity demand in 2020. The power crunch has prompted data center operators to look at new approaches in an effort to achieve more sustainable data centers. Compute is an area of focus as servers consume 43% of the total data center power.

The new approaches must meet the twin objectives of reducing power consumption per unit of compute capacity while executing workloads at high speed with ultra-low latency. This fresh thinking has been sharpened by the nature of the computing workloads which are migrating rapidly to a cloud native model. According to Outsystems, analyst firms Gartner and IDC both forecast that “90-95% of apps will be cloud native by 2025”. These dynamics have created the perfect conditions for a paradigm shift to a cloud native computing platform which combines high performance and power efficiency, and which is attuned to the requirements of cloud native applications.

Ampere was founded with a mission to lead this paradigm shift. Since Ampere's founding in 2017, it has shipped two Cloud Native Processor products, Ampere®Altra® and Ampere®Altra®Max, that have won wide support from leading cloud service providers and large digital enterprises. The innovative architecture on which these products are built meets the challenge which data centers pose. Unique features of this architecture enable Ampere to implement processor designs which feature more, and more power-efficient, cores. These processors' ultra-high core density executes cloud native applications fast and efficiently and maximizes the utilization of compute resources at rack- and data center-level.

This architecture is the reason that Ampere's Cloud Native Processors deliver significantly higher processing performance while reducing power consumption, compared with the legacy x86 architecture. It is exemplified by the Ampere® Altra® Max, a 128-core processor which offers up to 3.6x the rack-level performance at the same power consumption compared to the best x86-based device. *(Refer to Figure 4)*

This paper explains why a cloud native computing architecture which puts more high-efficiency cores to work on cloud native applications is helping the data center industry to meet growth in demand while limiting power consumption.

## Cloud Native Processing: an architecture for the data centers of today and tomorrow

Ampere introduced the Cloud Native Processing architecture, which produces remarkable performance and power efficiency. Ampere, founded in 2017, is already shipping two products, the 80-core Ampere® Altra® and 128-core Ampere Altra Max, to many of the world's leading cloud service providers.

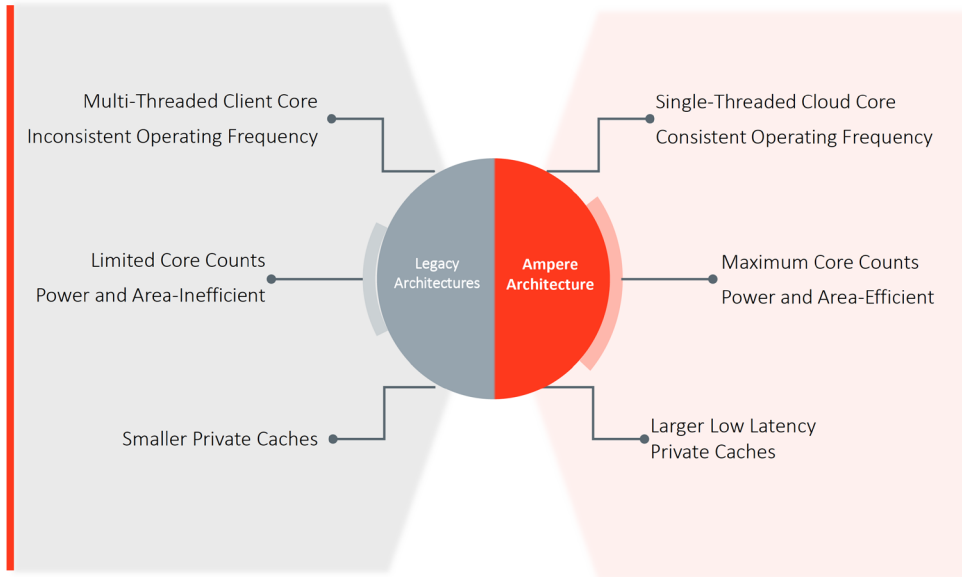


Figure 1

The Ampere Cloud Native Processor architecture has three key features:

**Single-threaded** execution within each of the 128 cores allows for a consistent performance and operating frequency over time and across the processor. By contrast, x86-based processors with fewer cores try to match the compute capacity by implementing multi-threading, where resources within the core are shared by multiple processes. This creates conflict between different workloads as they fight for the same resources, resulting in unpredictable performance. This forces cloud service providers to implement difficult work-arounds or to tolerate inefficient use of compute resources in order to maintain quality of service to all tenants.

**Maximum number of cores** that are power- and area-efficient for a cloud server environment. Ampere Altra Max features 128 high-performance cores in a chip footprint comparable to that of the latest x86-based server processors, while achieving dramatic reductions in system power consumption<sup>1</sup>.

**Larger, high-speed private and low latency caches** when compared to the x86 architecture<sup>2</sup>. The larger private caches accelerate the performance of each core's individual workloads without creating conflict between various users for the same resources. This is appropriate for cloud native computing, where nearly all processes are executed privately in each core. In the x86 architecture, the opposite is the case: a large higher latency shared cache supports the many shared functions in a client computing environment, and there is a correspondingly small provision of private cache.

## Core density: the key to high performance in cloud native computing

In a multi-tenant cloud computing service environment, each workload or microservice is best executed as a single thread in its own core. This ensures predictable high performance and avoids conflicts over allocation and prioritization when compute resources are shared.

Hardware isolation of different workloads also reinforces the security protection embedded in a cloud service provider's software systems.

It is also more power efficient to deliver maximum throughput via balanced system level performance across many cores and many users than to allow certain users to consume an unfair share of resources in a power hungry manner while throttling others.

This means that high core density offers multiple benefits in cloud native computing. The more cores a processor contains:

- The more concurrent workloads it can execute
- The more predictable its performance
- The lower its power consumption across all workloads

In the cloud, more cores means more compute capacity. And high core density is an inherent feature of the Cloud Native Processing architecture developed by Ampere. A comparison of the performance of x86 processors with Ampere products makes this clear.

The Ampere Altra family includes the industry's first 128-core cloud native processor – or any general purpose data center CPU for that matter. As of September 2022, no available processor based on the x86 architecture offered more than 64 cores (AMD EPYC).

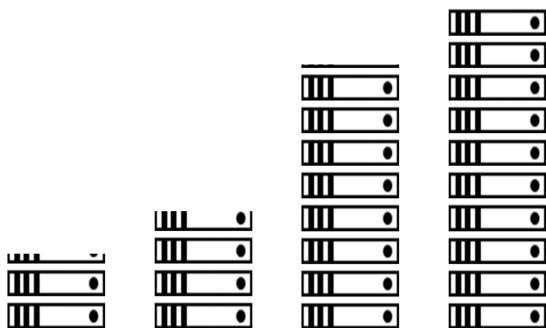
The 128-core Altra Max processor can be compared to the latest x86 server processors from leading manufacturers Intel and AMD in terms of the key parameters which affect data center performance: core density, core efficiency, and core utilization.

The Altra Max offers close to three times higher core density than any x86-based competitor – here, typical aggregate core counts per rack are shown for x86 and for Ampere Altra processors:

## Ampere is the Rack-Level Core Density Leader

*Ampere's Architectural Approach Delivers Highest Capacity per Rack and Data Center*

Xeon 8380	EPYC 7763	Altra Q64-22	Altra Max M128-26
760 Cores	1216 Cores	2688 Cores	3328 Cores



Based on a 12kW Rack

Figure 2

Thanks to cores which are optimized for efficiency, the Ampere Altra family processors also cut average power consumption in half.

And because of the Cloud Native Processing architecture’s ability to maintain a consistent high operating frequency, utilization across multiple heterogeneous workloads is far higher than in the x86-based alternative:

## Ampere is the Utilization Leader

*Ampere’s Architectural Approach Delivers Most Predictable and Scalable Performance*

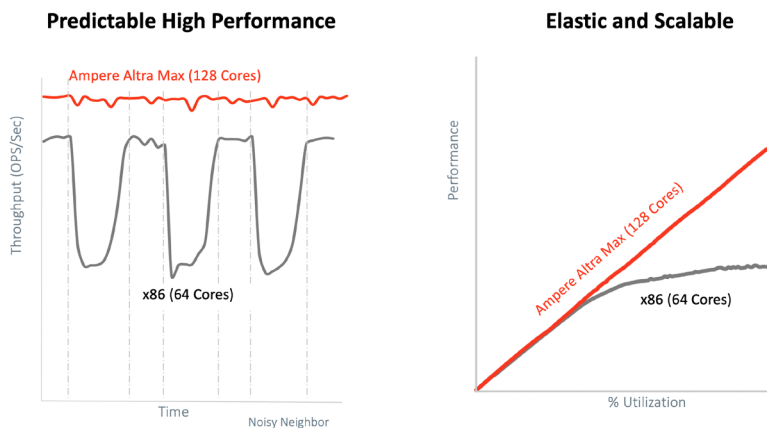
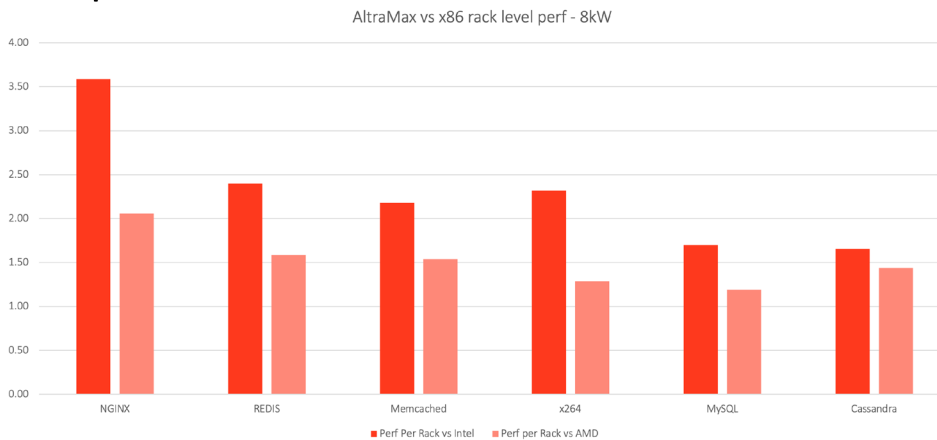


Figure 3

The result of this superiority in individual parameters of processor performance is markedly superior performance at the rack level:

## Leadership in Rack-Level-Perf Across Cloud Native Workloads



1.2X to 3.6X rack-level performance vs x86 across a broad set of cloud native workloads

Figure 4

Across a typical set of workload types, the 128-core Ampere Altra Max offers up to 3.6x the rack-level performance at the same power consumption compared to the best x86-based device.

## Conclusion

Core density is a crucial factor in the superior performance of Ampere Altra family processors in the cloud computing environment. This is clear in benchmark comparisons of rack-level performance when executing typical cloud computing workloads.

Ampere Computing, which has been shipping the 128-core Ampere Altra Max product to customers since 2021, is the clear leader in core density.

Thanks to its development of Cloud Native Processing, Ampere provides a proven, effective implementation of 128-core processing, offering superior performance on the parameters which matter to data center operators: compute capacity per rack, power consumption per rack, and compute utilization per rack.

## End Notes

1: Ampere Altra family processor typically run 30-40% lower than x86 competitors under common loads at high utilization. While running the industry standard benchmark Est. SPECrate<sup>®</sup>2017\_int\_base, Altra Max runs ~30% lower on average over the course of the benchmark run than rated TDP, while competitors run at their TDP or at times higher than their rated TDP.

2: Typical Altra family L2 (private) cache size is 1MB per thread (1/core) as compared to x86 systems which range from 256kB to ~512kB per thread (2/core), based on publicly available specification data. Multi-threaded caches used in x86 systems are shared between the threads and therefore not private.

*Figure 2 (Rack Density)*

Cores count based on 12kW rack with system usage power measured under SIR2017 load, based on publicly available specification data.

- Intel Ice Lake – 40 Cores
- AMD Milan – 64 Cores
- Ampere Altra – 80 Cores
- Ampere Altra Max - 128 Cores

*Figure 3 (Scalability and Predictability)*

- Predictability
  - Performance of Redis combined with StressNG workload run on nonutilized cores/ threads intermittently
  - Test run: Ampere Altra Max (128 cores) vs AMD Milan (64 core/128 threads)
  - Single socket tested
- Scalability
  - Performance of x.264 run on single instance per thread or core
  - Ampere Altra Max – 128 Cores (single socket)
  - AMD Milan – 64 Cores and 128 threads (single socket)

*Figure 4 (Rack Level Performance Configurations)*

System configurations, components, software versions and testing environments that differ from those used in Ampere's tests may result in different measurements from those obtained by Ampere<sup>®</sup>. The system configurations and components used in our testing were performed on bare-metal servers with one CPU socket with equivalent memory, storage, and networking options for all platforms referenced and then scaled linearly to the rack level using an 8kW rack as the limit.

The processors used were AMD EPYC 7763 ("Milan"), Intel Xeon 8380 ("Ice Lake"), and Ampere<sup>®</sup> Altra<sup>®</sup> Max M128-30

Specific test configurations are noted below:

Operating System: CentOS 8.0.1905 (kernel 4.18.0, 64k pages on Ampere<sup>®</sup> Altra<sup>®</sup> Max, 4k pages on Intel and AMD processors)

Memory: 8x 64 GB DIMMs, DDR4-3200

Networking: Mellanox ConnectX-5 100 Gb NIC

Storage: 1-4 NVMe drives depending on the workload across all three platforms

SMT: Enabled on the Intel and AMD platforms, not available on the Ampere<sup>®</sup> Altra<sup>®</sup> Max platform.

## Disclaimer

All data and information contained in or disclosed by this document are for informational purposes only and are subject to change.

This document may contain technical inaccuracies, omissions and typographical errors, and Ampere® Computing LLC, and its affiliates (“Ampere®”), is under no obligation to update or otherwise correct this information. Ampere® makes no representations or warranties of any kind, including express or implied guarantees of noninfringement, merchantability or fitness for a particular purpose, regarding the information contained in this document and assumes no liability of any kind. Ampere® is not responsible for any errors or omissions in this information or for the results obtained from the use of this information. All information in this presentation is provided “as is”, with no guarantee of completeness, accuracy, or timeliness.

This document is not an offer or a binding commitment by Ampere®. Use of the products and services contemplated herein requires the subsequent negotiation and execution of a definitive agreement or is subject to Ampere’s Terms and Conditions for the Sale of Goods.

This document is not to be used, copied, or reproduced in its entirety, or presented to others without the express written permission of Ampere®.

The technical data contained herein may be subject to U.S. and international export, re-export, or transfer laws, including “deemed export” laws. Use of these materials contrary to U.S. and international law is strictly prohibited.

Intel is a trademark of Intel Corporation. AMD is a trademark of Advanced Micro Devices, Inc. All other trademarks and copyrights are property of their respective owners and are only mentioned for informative purposes.

© 2022 Ampere® Computing LLC. All rights reserved. Ampere®, Ampere® Computing, Altra and the Ampere® logo are all trademarks of Ampere® Computing LLC or its affiliates. Other product names used in this publication are for identification purposes only and may be trademarks of their respective companies.

### Ampere Computing

4655 Great American Parkway  
Suite 601 Santa Clara, CA 95054

### Contact Us

Tel: +1-699-770-3700  
Info@amperecomputing.com

### Contact Sales

<https://amperecomputing.com/company/contact-sales.html>