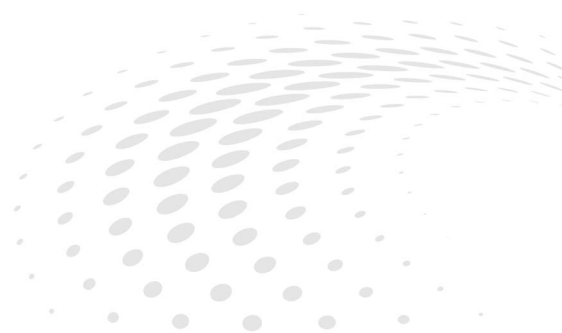


# TUNING GUIDE

## AMD EPYC 8004



# BIOS & Workload

|             |                 |
|-------------|-----------------|
| Publication | 58306           |
| Revision    | 1.0             |
| Issue Date  | September, 2023 |



© 2023 Advanced Micro Devices, Inc. All rights reserved.

The information contained herein is for informational purposes only and is subject to change without notice. While every precaution has been taken in the preparation of this document, it may contain technical inaccuracies, omissions and typographical errors, and AMD is under no obligation to update or otherwise correct this information. Advanced Micro Devices, Inc. makes no representations or warranties with respect to the accuracy or completeness of the contents of this document, and assumes no liability of any kind, including the implied warranties of noninfringement, merchantability or fitness for particular purposes, with respect to the operation or use of AMD hardware, software or other products described herein. No license, including implied or arising by estoppel, to any intellectual property rights is granted by this document. Terms and limitations applicable to the purchase or use of AMD's products are as set forth in a signed agreement between the parties or in AMD's Standard Terms and Conditions of Sale.

## Trademarks

AMD, the AMD Arrow logo, AMD EPYC, and combinations thereof are trademarks of Advanced Micro Devices, Inc. Other product names and links to external sites used in this publication are for identification purposes only and may be trademarks of their respective companies.

\* Links to third party sites are provided for convenience and unless explicitly stated, AMD is not responsible for the contents of such linked sites and no endorsement is implied.

| Date      | Version | Changes                |
|-----------|---------|------------------------|
| Jul, 2023 | 0.1     | Initial NDA release    |
| Sep, 2023 | 1.0     | Initial public version |
|           |         |                        |
|           |         |                        |
|           |         |                        |
|           |         |                        |

## Audience

This tuning guide describes best practices for optimizing performance using AMD default BIOS settings. It is intended for a technical audience such as system architects, production deployment, and performance engineering teams with:

- A background in configuring servers.
- Administrator-level access to both the server management Interface (BMC) and the OS.
- Familiarity with both the BMC and OS-specific configuration, monitoring, and troubleshooting tools.

## Authors

Muhammad Ashfaq, Anil Rajput, and Jesse Rangel

# Table of Contents

|                  |   |           |
|------------------|---|-----------|
| <b>Chapter 1</b> | <b>Introduction</b>                                     | <b>1</b>  |
| <b>Chapter 2</b> | <b>BIOS Defaults Summary</b>                            | <b>3</b>  |
| 2.1              | Processor Core Settings                                 | 4         |
| 2.2              | Power Efficiency Settings                               | 5         |
| 2.3              | NUMA and Memory Settings                                | 6         |
| 2.4              | Infinity Fabric Settings                                | 7         |
| 2.5              | PCIe, I/O, Security, and Virtualization Settings        | 8         |
| 2.6              | Higher-Level Settings                                   | 9         |
| <b>Chapter 3</b> | <b>BIOS Option Details</b>                              | <b>11</b> |
| 3.1              | Processor Core Settings                                 | 11        |
| 3.1.1            | Symmetric Multithreading (SMT) Settings                 | 11        |
| 3.1.2            | Cache Prefetchers                                       | 12        |
| 3.1.3            | Core Performance Boost                                  | 13        |
| 3.1.4            | Global C-States Control                                 | 13        |
| 3.2              | Power Management Settings                               | 14        |
| 3.2.1            | Power Profile Selection                                 | 14        |
| 3.2.2            | Power vs. Performance Determinism Settings              | 14        |
| 3.2.3            | Processor Cooling and Power Dissipation Limit Settings  | 14        |
| 3.2.4            | ACPI–Collaborative Processor Performance Control (CPCC) | 15        |
| 3.3              | NUMA and Memory Settings                                | 15        |
| 3.3.1            | L3 Cache as NUMA Domain                                 | 15        |
| 3.3.2            | NUMA Nodes per Socket (NPS)                             | 16        |
| 3.3.3            | Memory Target Speed                                     | 17        |
| 3.3.4            | Memory Interleaving                                     | 17        |
| 3.4              | Infinity Fabric Settings                                | 17        |
| 3.4.1            | Link Speed  | 17        |
| 3.4.2            | xGMI Link Width Management                              | 18        |
| 3.4.3            | Power States  | 19        |
| 3.4.4            | DF C-States   | 19        |
| 3.5              | PCIe, I/O, Security, and Virtualization Settings        | 20        |
| 3.5.1            | APIC Settings   | 20        |
| 3.5.2            | PCIe Speed PMM Control                                  | 20        |
| 3.5.3            | SR-IOV Settings   | 21        |
| 3.5.4            | PCIe Ten Bit Tag  | 21        |
| 3.5.5            | Input-Output Memory Management Unit (IOMMU) Settings    | 21        |
| 3.5.6            | Transparent Secure Memory Encryption (TSME)             | 22        |
| 3.5.7            | SEV, SEV-ES, and SEV-SNP                                | 22        |
| 3.5.8            | AVIC & x2AVIC   | 22        |

|                  |  |           |
|------------------|--|-----------|
| <b>Chapter 4</b> | <b>Workload-Specific BIOS Settings</b> | <b>23</b> |
| 4.1              | General-Purpose Workloads              | 23        |
| 4.1.1            | Processor Core Settings                | 23        |
| 4.1.2            | Power Management Settings              | 24        |
| 4.1.3            | NUMA and Memory Settings               | 24        |
| 4.1.4            | Infinity Fabric Settings               | 25        |
| 4.1.5            | I/O Settings                           | 25        |
| 4.2              | Memory and I/O Intensive Workloads     | 26        |
| 4.2.1            | Processor Core Settings                | 26        |
| 4.2.2            | Power Management Settings              | 26        |
| 4.2.3            | NUMA and Memory Settings               | 27        |
| 4.2.4            | Infinity Fabric Settings               | 27        |
| 4.2.5            | I/O Settings                           | 27        |
| 4.3              | Virtualization and Containers          | 28        |
| 4.3.1            | Processor Core Settings                | 28        |
| 4.3.2            | Power Management Settings              | 28        |
| 4.3.3            | NUMA and Memory Settings               | 29        |
| 4.3.4            | Infinity Fabric Settings               | 29        |
| 4.3.5            | I/O Settings                           | 29        |
| 4.4              | OS, Analytics, and Telco Settings      | 30        |
| 4.4.1            | Processor Core Settings                | 30        |
| 4.4.2            | Power Management Settings              | 30        |
| 4.4.3            | NUMA and Memory Settings               | 31        |
| 4.4.4            | Infinity Fabric Settings               | 31        |
| 4.4.5            | I/O Settings                           | 31        |
| <b>Chapter 5</b> | <b>Processor Identification</b>        | <b>33</b> |
| 5.1              | CPUID Instruction                      | 33        |
| 5.2              | New Software-Visible Features          | 34        |
| 5.2.1            | AVX-512                                | 34        |
| <b>Chapter 6</b> | <b>Debugging BIOS Setting Changes</b>  | <b>35</b> |

## Chapter

## 1

# Introduction

Default BIOS options generally produce the best overall performance for generic workloads, but these defaults may not be optimal for a specific workload. AMD continually tests various workloads; this tuning guide discusses BIOS options to deliver both maximum performance and performance-per-watt (power efficiency).

- [“BIOS Defaults Summary” on page 3](#) provides a quick overview of default AMD EPYC 8004 BIOS settings.
- [“BIOS Defaults Summary” on page 3](#) provides detailed information about the AMD EPYC 8004 BIOS options and the potential benefit of modifying each one.
- [“Workload-Specific BIOS Settings” on page 23](#) presents sample workloads and recommended BIOS settings. Keep in mind that these BIOS settings are not “one size fits all” because your specific workload(s) are not identical to synthetic benchmarks.

*Note: Not all platforms support all of the BIOS settings described in this Tuning Guide. Please contact your platform vendor if you cannot see one or more needed settings.*

AMD strongly recommends that customers download and install the latest BIOS update for your AMD EPYC 8004 Series Processor-based server from your platform vendor. BIOS updates often help customers by providing new and updated features, bug fixes, enhancements, security features, and other improvements. These improvements can help your system software stability and dependency modules (such as hardware, firmware, drivers, and software) by giving you a more robust environment to run your applications.



*This page intentionally left blank.*

## Chapter

## 2

# BIOS Defaults Summary

This chapter provides high-level lists of the default AMD EPYC 8004 BIOS settings and their default values. Please see Chapter 4 of the *BIOS & Workload Tuning Guide for AMD EPYC™ 8004 Series Processors* (available from [AMD EPYC Tuning Guides](#)) for detailed descriptions. Later chapters in this Tuning Guide discuss the BIOS options as they relate to a specific workload or set of workload(s).

*Note: The default setting name and values described in this chapter are the AMD default names and values that serve as recommendations for OEMs. End users must confirm their OEM BIOS setting availability and options.*

AMD strongly recommends that customers download and install the latest BIOS update for your AMD EPYC 8004 Series Processor-based server from your platform vendor. BIOS updates often help customers by providing new and updated features, bug fixes, enhancements, security features, and other improvements. These improvements can help your system software stability and dependency modules (such as hardware, firmware, drivers, and software) by giving you a more robust environment to run your applications.

## 2.1 Processor Core Settings

| Name                    | Default | Description  |
|-------------------------|---------|--|
| SMT Control             | Auto    | <ul style="list-style-type: none"> <li><b>Enabled/Auto:</b> Two hardware threads per core.</li> <li><b>Disabled:</b> Single hardware thread per core.</li> </ul>   |
| L1 Stream HW Prefetcher | Auto    | <ul style="list-style-type: none"> <li><b>Enabled/Auto:</b> Enables the prefetcher.</li> <li><b>Disabled:</b> Disables the prefetcher.</li> </ul>  |
| L1 Stride Prefetcher    | Auto    | <ul style="list-style-type: none"> <li><b>Enabled/Auto:</b> Enables the prefetcher.</li> <li><b>Disabled:</b> Disables the prefetcher.</li> </ul>  |
| L1 Region Prefetcher    | Auto    | <ul style="list-style-type: none"> <li><b>Enabled/Auto:</b> Enables the prefetcher.</li> <li><b>Disabled:</b> Disables the prefetcher.</li> </ul>  |
| L1 Burst Prefetch Mode  | Auto    | <ul style="list-style-type: none"> <li><b>Enabled/Auto:</b> Enables the prefetcher.</li> <li><b>Disabled:</b> Disables the prefetcher.</li> </ul>  |
| L2 Stream HW Prefetcher | Auto    | <ul style="list-style-type: none"> <li><b>Enabled/Auto:</b> Enables the prefetcher.</li> <li><b>Disabled:</b> Disables the prefetcher.</li> </ul>  |
| L2 Up/Down Prefetcher   | Auto    | <ul style="list-style-type: none"> <li><b>Enabled/Auto:</b> Enables the prefetcher.</li> <li><b>Disabled:</b> Disables the prefetcher.</li> </ul>  |
| ACPI CST C2 Latency     | Auto    | <ul style="list-style-type: none"> <li><b>Default/Auto:</b> 800 us (microseconds)</li> <li><b>Range Min Value:</b> 18 us Max Value: 1000 us</li> </ul>   |
| Core Performance Boost  | Auto    | <ul style="list-style-type: none"> <li><b>Enabled/Auto:</b> Enables Core Performance Boost.</li> <li><b>Disabled:</b> Disables Core Performance Boost.</li> </ul>  |
| BoostFmaxEn             | Auto    | <ul style="list-style-type: none"> <li><b>Auto:</b> Use the default Fmax</li> <li><b>Manual:</b> User can set the boost Fmax</li> </ul>  |
| BoostFmax               | Auto    | Specify the boost Fmax frequency limit to apply to all cores (MHz in decimal)  |
| Global C-State Control  | Auto    | <ul style="list-style-type: none"> <li><b>Enabled/Auto:</b> Controls IO based C-state generation and DF C-states, including core processor C-States</li> <li><b>Disabled:</b> AMD strongly recommends not disabling this option because this also disables core processor C-States.</li> </ul> |

Table 2-1: Processor core BIOS settings



## 2.2 Power Efficiency Settings

| Name                    | Default | Description  |
|-------------------------|---------|--|
| Power Profile Selection | Auto    | <ul style="list-style-type: none"> <li><b>Auto/0:</b> High-performance mode</li> <li><b>1:</b> Efficiency mode</li> <li><b>2:</b> Maximum I/O performance mode</li> </ul>  |
| Determinism Control     | Auto    | <ul style="list-style-type: none"> <li><b>Auto:</b> Use default performance determinism settings.</li> <li><b>Manual:</b> Specify custom performance determinism settings.</li> </ul>  |
| Determinism Enable      | Auto    | <ul style="list-style-type: none"> <li><b>Auto:</b> Performance.</li> <li><b>1:</b> Power.</li> </ul>  |
| TDP Control             | Auto    | <ul style="list-style-type: none"> <li><b>Auto:</b> Use platform- and OPN-default TDP.</li> <li><b>Manual:</b> Set custom configurable TDP.</li> </ul>   |
| TDP                     | OPN Max | This option appears once the user sets the <b>TDP Control</b> to <b>Manual</b> . <ul style="list-style-type: none"> <li><b>Values 70-225:</b> Set configurable TDP, in watts.</li> </ul>   |
| PPT Control             | Auto    | Enables or disables the <b>PPT</b> control. <ul style="list-style-type: none"> <li><b>Auto:</b> Automatically set PPL in watts.</li> <li><b>Manual:</b> Specify a custom PPL.</li> </ul>   |
| PPT                     | OPN Max | This option appears once the user sets the <b>PPT Control</b> to <b>Manual</b> . <ul style="list-style-type: none"> <li><b>Values 70-225:</b> Set configurable PPT, in watts.</li> </ul>   |
| CPPC                    | Auto    | <ul style="list-style-type: none"> <li><b>Enabled/Auto:</b> Allows the OS to make performance/power optimization requests using ACPI CPPC.</li> <li><b>Disabled:</b> Prevents the OS from making performance/power optimization requests using ACPI CPPC.</li> </ul> |

Table 2-2: Power efficiency BIOS settings

## 2.3 NUMA and Memory Settings

| Name  | Default  | Description  |
|---|----------|--|
| LLC as NUMA Domain<br>(ACPI SRAT L3 Cache as NUMA Domain) | Disabled | <ul style="list-style-type: none"> <li><b>Disabled (recommended):</b> Both NUMA nodes (<code>cpubind</code>) and memory interleaving (<code>membind</code>) are determined by the NPS setting.</li> <li><b>Enabled:</b> Overrides the NPS setting for # of NUMA nodes by mapping each LLC as a NUMA node. This does not impact the memory interleaving</li> </ul>  |
| Nodes Per Socket (NPS)                                    | 1        | <p><b>Memory Interleaving:</b> The <b>NPS</b> setting always determines the memory interleaving regardless of whether <b>LLC as NUMA</b> is <b>Enabled</b> or <b>Disabled</b>.</p> <p># of NUMA nodes (if <b>LLC as NUMA Domain</b> is <b>Disabled</b>):</p> <ul style="list-style-type: none"> <li><b>NPS1:</b> One NUMA node per socket (Most cloud providers use this as it provides consistent average memory latency to all the accesses within a socket).</li> <li><b>NPS2:</b> Two NUMA nodes per socket. This option is not supported by all AMD EPYC 8004 Series Processors.</li> <li><b>NPS4:</b> Four NUMA nodes per socket. This option is not supported by all AMD EPYC 8004 Series Processors and is only applicable for certain use cases.</li> </ul> <p>AMD recommends NPS1, NPS2, or NPS4 depending on your use case.</p> |
| Memory Target Speed                                       | Auto     | <ul style="list-style-type: none"> <li><b>Auto:</b> Determine the maximum memory speed based on SPD information from populated DIMMs and platform memory speed support.</li> </ul> <p>Alternatively, you can select:</p> <ul style="list-style-type: none"> <li><b>Values 3200–5600 MT/s:</b> Run the DRAM memory target clock speed at the specified speed. The DRAM memory target is the DDR rate.</li> </ul> <p>Your OEM system default value may vary.</p>   |
| Memory Interleaving                                       | Auto     | <ul style="list-style-type: none"> <li><b>Auto/Enable:</b> Enables memory interleaving.</li> <li><b>Disable:</b> Allows for disabling memory interleaving. The <b>NUMA Nodes per Socket</b> setting will be honored regardless of this setting. AMD strongly recommends not disabling this setting because most production deployments benefit from memory interleaving.</li> </ul>  |

Table 2-3: NUMA and memory BIOS settings

## 2.4 Infinity Fabric Settings

| Name  | Default | Description  |
|---|---------|--|
| 4 xGMI Link Max Speed*                            | Auto    | <ul style="list-style-type: none"> <li>20 Gbps</li> <li>25 Gbps/Auto</li> <li>32 Gbps</li> </ul> <p>Your OEM system default value may vary.</p>  |
| xGMI Link Width Control*                          | Auto    | <ul style="list-style-type: none"> <li><b>Auto:</b> Use the default xGMI link width controller settings.</li> <li><b>Manual:</b> Specify a custom xGMI link width controller setting.</li> </ul>   |
| xGMI Force Link Width* Control                    | Auto    | <ul style="list-style-type: none"> <li><b>Unforce:</b> Do not force the xGMI to a fixed width.</li> <li><b>Force:</b> Use the xGMI link to the user-specified width.</li> </ul>  |
| xGMI Force Link Width*                            | Auto    | <ul style="list-style-type: none"> <li><b>0:</b> Force xGMI link width to x4.</li> <li><b>1:</b> Force xGMI link width to x8.</li> <li><b>2:</b> Force xGMI link width to x16.</li> </ul>  |
| xGMI Max Link Width Control*                      | Auto    | <ul style="list-style-type: none"> <li><b>Auto:</b> Use the default xGMI link width controller settings.</li> <li><b>Manual:</b> Specify a custom xGMI link with controller setting.</li> </ul>  |
| xGMI Max Link Width*                              | Auto    | <ul style="list-style-type: none"> <li><b>0:</b> Set max xGMI link width to x4.</li> <li><b>1:</b> Set max xGMI link width to x8.</li> <li><b>2:</b> Set max xGMI link width to x16.</li> </ul>  |
| APBDIS  | Auto    | <ul style="list-style-type: none"> <li><b>0/Auto:</b> Dynamically switch the Infinity Fabric P-state based on link usage.</li> <li><b>1:</b> Enabled fixed Infinity Fabric P-state control.</li> </ul>   |
| DfPstate Range Support                            | Auto    | <ul style="list-style-type: none"> <li><b>Auto:</b> If this feature is enabled, the range value setting should follow the rule that MaxDfPstate ≤ MinDfPstate. Otherwise, it will not work.</li> <li><b>Enable:</b> Add the values MaxDfPstate &amp; MinDfPstate.</li> <li><b>Disable:</b> No MaxDfPstate &amp; MinDfPstate option.</li> </ul> |
| *This option only applies to dual-socket systems. |         |  |
| DF C-States                                       | Auto    | <p>Controls DF C-states.</p> <ul style="list-style-type: none"> <li><b>Disabled:</b> Prevents the AMD Infinity Fabric from entering a low-power state.</li> <li><b>Enabled/Auto:</b> Allows the AMD Infinity Fabric to enter a low-power state.</li> </ul>   |

Table 2-4: Infinity Fabric BIOS settings

## 2.5 PCIe, I/O, Security, and Virtualization Settings

| Name                     | Default    | Description   |
|--------------------------|------------|---|
| Local APIC Mode          | Auto(0x02) | <ul style="list-style-type: none"> <li>• <b>xAPIC:</b> Use xAPIC, supports up to 255 cores.</li> <li>• <b>x2APIC:</b> Supports more than 255 cores.</li> <li>• <b>Auto:</b> The system will choose the mode that best fits the number of active cores in the system.</li> <li>• <b>Compatibility:</b> Threads below 255 run in xAPIC with xAPIC ACPI structures, and threads 255 and above run in x2 mode with x2 ACPI structures.</li> <li>• <b>XApicMode (0x01):</b> Forces legacy xAPIC mode.</li> <li>• <b>X2ApicMode (0x02):</b> Forces x2APIC mode independent of thread count.</li> </ul>  |
| PCIe Speed PMM Control   | Auto       | <ul style="list-style-type: none"> <li>• <b>0:</b> Dynamic link speed determined by power management functionality.</li> <li>• <b>1:</b> Static Target Link Speed (Gen4); sets the maximum idle link speed to 16 GT/s.</li> <li>• <b>Auto/2:</b> Static Target Link Speed (Gen5); sets the maximum idle link speed to 32 GT/s, thereby disabling the feature).</li> </ul>   |
| PCIe ARI Support (SRIOV) | Auto       | <ul style="list-style-type: none"> <li>• <b>Enabled/Auto:</b> Enables Alternative Routing ID interpretation.</li> <li>• <b>Disabled:</b> Disables Alternative Routing ID interpretation.</li> </ul>   |
| PCIe Ten Bit Tag Support | Auto       | <ul style="list-style-type: none"> <li>• <b>Enabled/Auto:</b> Enables PCIe 10-bit tags for supported devices.</li> <li>• <b>Disabled:</b> Disables PCIe 10-bit tags for all devices.</li> </ul>   |
| IOMMU                    | Auto       | <ul style="list-style-type: none"> <li>• <b>Enabled/Auto:</b> Enables IOMMU. AMD recommends setting this to <code>pt:pass-through</code> in the Linux kernel settings.</li> <li>• <b>Disabled:</b> Disables IOMMU.</li> </ul>   |
| AVIC                     | Enabled    | <p>Advanced Virtual Interrupt Controller. AMD recommends enabling this option for virtualized deployments.</p> <ul style="list-style-type: none"> <li>• <b>Disabled:</b> Disables AVIC.</li> <li>• <b>Enabled:</b> Enables AVIC.</li> </ul>   |
| x2AVIC                   | Enabled    | <p>x2AVIC is an extension of the advanced virtual interrupt controller. AMD recommends enabling this option for virtualized deployments where the total number of vCPUs &gt;255.</p> <p>x2AVIC support varies on platforms powered by AMD EPYC processors based on the AMD EPYC generation and OS (Linux, Windows, etc.). Please refer to respective OS support pages for help enabling x2AVIC. Linux Kernel 5.20 and onward include upstream x2AVIC support, and AMD strongly recommends upstream Linux Kernel 6.5 because this version contains several optimizations.</p> <ul style="list-style-type: none"> <li>• <b>Disabled:</b> Disables x2AVIC.</li> <li>• <b>Enabled:</b> Enables x2AVIC.</li> </ul> |

Table 2-5: PCIe, I/O, security, and virtualization BIOS settings

|         |          |   |
|---------|----------|---|
| TSME    | Auto     | <ul style="list-style-type: none"> <li>• <b>Auto/Disabled:</b> Disables transparent secure memory encryption.</li> <li>• <b>Enabled:</b> Enables transparent secure memory encryption.</li> </ul>   |
| SEV     | Disabled | <p>In a multi-tenant environment (such as a cloud), Secure Encrypted Virtualization (SEV) mode isolates virtual machines from each other and from the hypervisor.</p> <ul style="list-style-type: none"> <li>• <b>Disabled:</b> SEV is disabled.</li> <li>• <b>Enabled:</b> SEV is enabled.</li> </ul>  |
| SEV-ES  | Disabled | <p>Secure Encrypted Virtualization-Encrypted State (SEV-ES) mode extends SEV protection to the contents of the CPU registers by encrypting them when a virtual machine stops running. Combining SEV and SEV-ES can reduce the attack surface of a VM by helping protect the confidentiality of data in memory.</p> <ul style="list-style-type: none"> <li>• <b>Disabled:</b> SEV-ES is disabled.</li> <li>• <b>Enabled:</b> SEV-ES is enabled.</li> </ul>   |
| SEV-SNP | Disabled | <p>Secure Encrypted Virtualization-Secure Nested Paging (SEV-SNP) mode builds on SEV and SEV-ES by adding strong memory integrity protection to create an isolated execution environment that helps prevent malicious hypervisor-based attacks such as data replay and memory re-mapping. SEV-SNP also introduces several additional optional security enhancements that support additional VM use models, offer stronger protection around interrupt behavior, and increase protection against recently-disclosed side channel attacks.</p> <ul style="list-style-type: none"> <li>• <b>Disabled:</b> SEV-SNP is disabled.</li> <li>• <b>Enabled:</b> SEV-SNP is enabled.</li> </ul> |

Table 2-5: PCIe, I/O, security, and virtualization BIOS settings (Continued)

## 2.6 Higher-Level Settings

The system powers on to an initial state, after which succeeding software layers may affect system settings:

1. System firmware validates basic hardware functionality and allows users to change various settings via the BIOS Setup menus.
2. UEFI provides a shell environment that allows users to further interact with the system.
3. The operating system or hypervisor is the next software layer that provides control over system hardware.
4. Lastly, certain applications can also affect underlying hardware.

Each of the lines above may alter settings made by prior line, and some user changes require a reboot to take effect.

Please refer to your OEM documentation and/or applicable AMD Tuning Guide(s) for further guidance.



*This page intentionally left blank.*

**Chapter****3**

# BIOS Option Details

## 3.1 Processor Core Settings

### 3.1.1 Symmetric Multithreading (SMT) Settings

Enabling SMT causes neutral to negative performance impacts on some workloads, especially HPC. Also, some application licenses count the number of hardware threads enabled instead of the physical core count. It may therefore be best to disable SMT on your AMD EPYC 8004 Series Processor.

| Setting     | Options   |
|-------------|---|
| SMT Control | <ul style="list-style-type: none"><li>• <b>Enable/Auto:</b> Two hardware threads per core.</li><li>• <b>Disable:</b> Single hardware thread per core.</li></ul> |

*Table 3-1: SMT settings*

### 3.1.2 Cache Prefetchers

Most workloads and production deployments benefit from the L1 & L2 Stream Hardware prefetchers gathering data and keeping the core pipeline busy, but some workloads that stress the memory bandwidth to its maximum capacity may perform better when some or all prefetchers are disabled. All prefetchers are enabled by default. Be sure to evaluate the prefetchers for your deployments.

| Setting                 | Options   |
|-------------------------|---|
| L1 Stream HW Prefetcher | <p>This prefetcher uses the history of L1 cache memory access patterns to fetch additional sequential lines in ascending or descending order.</p> <ul style="list-style-type: none"> <li>• <b>Disable:</b> Disable prefetcher.</li> <li>• <b>Enable:</b> Enable prefetcher.</li> </ul>  |
| L1 Stride Prefetcher    | <p>The prefetcher uses the L1 cache memory access history of individual instructions to fetch additional lines when each access is a constant distance from the previous.</p> <ul style="list-style-type: none"> <li>• <b>Disable:</b> Disable prefetcher.</li> <li>• <b>Enable:</b> Enable prefetcher.</li> </ul>  |
| L1 Region Prefetcher    | <p>This prefetcher uses the L1 cache memory access history to fetch additional lines when the data access for a given instruction that tends to be followed by a consistent pattern of other accesses within a localized region.</p> <ul style="list-style-type: none"> <li>• <b>Disable:</b> Disable prefetcher.</li> <li>• <b>Enable:</b> Enable prefetcher.</li> </ul> |
| L1 Burst Prefetch Mode  | <p>This prefetcher uses the L1 cache memory access history to fetch additional lines when the data access for a given instruction that tends to be followed by a consistent pattern of other accesses within a localized region.</p> <ul style="list-style-type: none"> <li>• <b>Disable:</b> Disable prefetcher.</li> <li>• <b>Enable:</b> Enable prefetcher.</li> </ul> |
| L2 Stream HW Prefetcher | <p>This prefetcher uses the history of L2 cache memory access patterns to fetch additional sequential lines in ascending or descending order.</p> <ul style="list-style-type: none"> <li>• <b>Disable:</b> Disable prefetcher.</li> <li>• <b>Enable:</b> Enable prefetcher.</li> </ul>  |
| L2 Up/Down Prefetcher   | <p>Uses the L2 cache memory access history to determine whether to fetch the next or previous line for all memory accesses.</p> <ul style="list-style-type: none"> <li>• <b>Disable:</b> Disable prefetcher.</li> <li>• <b>Enable:</b> Enable prefetcher.</li> </ul>  |

Table 3-2: Cache prefetcher settings



### 3.1.3 Core Performance Boost

**Core Performance Boost** can be enabled or disabled. Enabling this setting allows the processor to opportunistically increase a set of CPU cores to higher than the CPU's rated base clock speeds based on the number of active cores, power, and thermal headroom in a system.

Some workloads don't need maximum core frequency to achieve acceptable performance. Limiting the maximum core boost frequency can reduce power consumption. The **BoostFmax** setting limits the maximum boost frequency but does not set a fixed frequency. The SoC will not exceed the maximum algorithm-allowable frequency if **BoostFmax** is set too high. Actual boost performance depends on many factors, including the other settings discussed in this tuning guide.

| Setting                | Options   |
|------------------------|---|
| Core Performance Boost | <ul style="list-style-type: none"> <li>• <b>Enable/Auto:</b> Enables Core Performance Boost.</li> <li>• <b>Disable:</b> Disables Core Performance Boost.</li> </ul> |
| BoostFmaxEn            | <ul style="list-style-type: none"> <li>• <b>Manual:</b> Use specified BoostFmax setting.</li> <li>• <b>Auto:</b> Use default BoostFmax setting.</li> </ul>          |
| BoostFmax              | Specify the boost Fmax frequency limit to apply to all cores (MHz in decimal).  |

Table 3-3: Core boost settings

### 3.1.4 Global C-States Control

The **Global C-States Control** enables and disables C-states on the server across all cores. Disabling this feature means that the CPU cores can only be in C0 (active) or C1 state because the C1 state cannot be disabled. A CPU core will be in C1 state if the core is halted by the OS. IO based C-state generation and DF C-states include core processor C-States. If you have a low latency or extremely low jitter use case, then consider disabling DF C-states as described in this Tuning Guide. AMD strongly recommends not disabling Global C-states except for debugging.

| Setting                | Options  |
|------------------------|--|
| Global C-State Control | <ul style="list-style-type: none"> <li>• <b>Enabled/Auto:</b> Controls IO based C-state generation and DF C-states, including core processor C-States</li> <li>• <b>Disabled:</b> AMD strongly recommends not disabling this option because this also disables core processor C-States.</li> </ul> |

## 3.2 Power Management Settings

### 3.2.1 Power Profile Selection

| Setting                         | Options   |
|---------------------------------|---|
| Power Profile Selection Control | <ul style="list-style-type: none"> <li><b>0:</b> High Performance mode (Default)</li> <li><b>1:</b> Efficiency Mode</li> <li><b>2:</b> Maximum IO Performance Mode</li> </ul> |

Table 3-4: Power profile selection

### 3.2.2 Power vs. Performance Determinism Settings

The **Determinism Enable** selects between:

- Performance (default for most OPNs):** Uniform performance across identically configured systems in a datacenter. Set TDP and PPL to the same value, as described in [“Processor Cooling and Power Dissipation Limit Settings” on page 14](#).
- Power:** Maximum performance of any individual system by leveraging the capabilities of a given CPU to the maximum, resulting in a varying performance range across the datacenter or larger deployments.

| Setting             | Options  |
|---------------------|--|
| Determinism Control | <ul style="list-style-type: none"> <li><b>Auto:</b> Use default performance determinism settings.</li> <li><b>Manual:</b> Specify custom power/performance determinism.</li> </ul>   |
| Determinism Enable  | <ul style="list-style-type: none"> <li><b>Auto:</b> This setting may be either <b>Power</b> or <b>Performance</b> based on OEM Platform and OPN selection.</li> <li><b>0:</b> Power (default)</li> <li><b>1:</b> Performance</li> </ul> <p>See <a href="#">“Processor Cooling and Power Dissipation Limit Settings” on page 14</a> for additional information.</p> |

Table 3-5: Power/performance settings

### 3.2.3 Processor Cooling and Power Dissipation Limit Settings

Thermal Design Power (TDP) allows modifying the CPU cooling limit and the Package Power Limit (PPL) allows modifying the CPU Power Dissipation Limit. Many platforms configure TDP to the maximum CPU-supported value. Most platforms also set the PPL to the same value as the TDP.

If you are using **Performance** determinism, then both the TDP and PPT must be set to the same value, as described in [“Power vs. Performance Determinism Settings” on page 14](#). You can set the PPL to a lower value than the TDP to reduce system operating power. If you do this, then the CPU will control the Core Boost to keep the socket power dissipation at or below the PPL value.

If you are using **Power** determinism, then you can obtain maximum performance by setting the TDP and PPL to the maximum TDP value supported by the CPU. Setting the TDP and PPL to **Auto** sets both parameters to the CPU default TDP value for energy-efficient operation.

| Setting     | Options   |
|-------------|---|
| TDP Control | <ul style="list-style-type: none"> <li>• <b>Auto:</b> Use platform- and OPN-default TDP.</li> <li>• <b>Manual:</b> Set custom configurable TDP.</li> </ul>                                    |
| TDP         | This option is available if the user sets the <b>TDP Control</b> to <b>Manual</b> . <ul style="list-style-type: none"> <li>• <b>Values 70–225:</b> Set configurable TDP, in watts.</li> </ul> |
| PPT Control | Enables or disables the <b>PPT</b> control. <ul style="list-style-type: none"> <li>• <b>Auto:</b> Use platform- and OPN-default-PPL.</li> <li>• <b>Manual:</b> Set customized PPL.</li> </ul> |
| PPT         | This option is available if the user sets the <b>PPT Control</b> to <b>Manual</b> . <ul style="list-style-type: none"> <li>• <b>Values 70–225:</b> Set PPT, in watts.</li> </ul>              |

Table 3-6: TDP settings

### 3.2.4 ACPI—Collaborative Processor Performance Control (CPCC)

Enabling CPCC allows the OS to help maintain energy efficiency by controlling when and how much turbo can be applied. ACPI 5.0 introduced this feature. Not all operating systems support CPCC. Microsoft began supporting CPCC with Windows® Server® 2016.

| Setting | Options   |
|---------|---|
| CPCC    | <ul style="list-style-type: none"> <li>• <b>Disabled:</b> Disabled.</li> <li>• <b>Enabled:</b> Allow the OS to make performance/power optimization requests using ACPI CPPC.</li> </ul> |

Table 3-7: CPCC settings

## 3.3 NUMA and Memory Settings

This section describes NUMA- and memory-related BIOS settings.

### 3.3.1 L3 Cache as NUMA Domain

This setting controls automatic or manual generation of distance information in the ACPI System Locality Information Table (SLIT) and NUMA proximity domains in the System Resource Affinity Table (SRAT). Some hypervisors and operating systems do not perform L3-aware scheduling, and some workloads will benefit from having the L3 declared as a NUMA domain. In dual-socket systems, the remote socket distance can affect memory allocation decisions. Setting this to a value of at least 32 (32 recommended) may improve scheduling of lightly-threaded workloads. Setting this to a value less than 32 (22 recommended) may improve scheduling of heavily-threaded workloads. In general:

- If a workload spans two sockets, then set the distance to < 32.
- If the workload can be confined to a socket, then set the distance to 32.

### 3.3.2 NUMA Nodes per Socket (NPS)

This setting enables a trade-off between minimizing local memory latency for NUMA-aware or highly parallelizable workloads vs. maximizing per-core memory bandwidth for non-NUMA-friendly workloads. NPS2 and/or NPS4 may not be an option on certain OPNs or with certain memory populations.

- **NPS1:** Indicates a single NUMA node per socket. This setting configures all memory channels on the processor into a single NUMA domain. All of the processor cores, all attached memory, and all PCIe devices connected to the SoC are in that one NUMA domain. Memory accesses are interleaved across all 24 memory channels into a single address space. The default configuration (one NUMA domain per socket) is recommended for most workloads.
- **NPS2:** 2 NUMA domains per socket, which interleaves the corresponding six memory channels within the same 6 CCD NUMA domain. Half of the cores and half of the memory channels of the SoC are grouped together into one NUMA domain, with the remaining cores and memory channels grouped into the second NUMA domain. Memory is interleaved across the six memory channels of each NUMA domain. Not all AMD EPYC 8004 Series Processors support this option.
- **NPS4:** 4 partitions the processor into four NUMA domains with each logical quadrant configured as its own NUMA domain. Memory is interleaved across the three memory channels within each quadrant. PCIe devices will be local to one of the four NUMA domains depending on the quadrant (of the I/O die) that has the PCIe root complex for that device. Every pair of memory channels is interleaved. AMD suggests NPS1 or NPS2 for most workloads. This option is not supported by all AMD EPYC 8004 Series Processors and is only applicable for certain use cases.

Enabling **ACPI SRAT L3 Cache as NUMA Domain** (another name for **L3 as NUMA**) determines the number of NUMA nodes and overrides the number of NUMA nodes specified by the NPS setting while still using the NPS setting to determine the memory interleaving granularity.

| Setting   | Options   |
|---|---|
| L3 Cache as NUMA Domain (ACPI SRAT L3 Cache as NUMA Domain) | <ul style="list-style-type: none"> <li>• <b>Disabled/Auto:</b> Do not report each L3 cache to the OS as a NUMA domain.</li> <li>• <b>Enabled:</b> Report each L3 cache to the OS as a NUMA domain.</li> </ul>   |
| NUMA Node Per Socket  | <ul style="list-style-type: none"> <li>• <b>NPS1:</b> Interleaves memory accesses across all channels in each socket and report one NUMA node per socket unless <b>L3 Cache as NUMA</b> is enabled.</li> <li>• <b>NPS2:</b> Interleaves memory accesses the channels associated with each half of a socket and reports two NUMA nodes per socket unless <b>L3 Cache as NUMA</b> is enabled. This option is not supported by all AMD EPYC 8004 Series Processors.</li> <li>• <b>NPS4:</b> Interleaves memory accesses across the channels associated with a quadrant of each socket and reports four NUMA nodes per socket unless <b>L3 Cache as NUMA</b> is enabled. This option is not supported by all AMD EPYC 8004 Series Processors and is only applicable for certain use cases.</li> </ul> |

Table 3-8: NPS settings

### 3.3.3 Memory Target Speed

By default, the 4th Gen AMD EPYC processor BIOS runs at the maximum target frequency allowed by the platform and DIMM. This configuration allows maximum memory bandwidth and lowest latency for the processor. Lowering the memory clock speed reduces memory controller power consumption and allows the rest of the SoC to consume more power, thereby potentially boosting performance elsewhere for certain workloads.

| Setting             | Options   |
|---------------------|---|
| Memory Target Speed | <ul style="list-style-type: none"> <li><b>Auto:</b> Determine maximum memory speed based on SPD information from populated DIMMs and platform memory speed support.</li> <li><b>Values 3200–5600 MT/s:</b> Run the DRAM memory target speed at the specified speed (the DRAM memory target is the DDR rate.)</li> </ul> |

Table 3-9: Memory clock settings

### 3.3.4 Memory Interleaving

This setting allows you to enable or disable memory interleaving within a NUMA node. The NUMA Nodes per Socket (NPS) setting will be honored regardless of this setting. This BIOS setting does not impact the number of NUMA nodes or how memory channels are mapped to the NUMA nodes.

*Note: AMD strongly recommends not disabling this setting because most applications and deployments benefit from memory interleaving.*

| Setting             | Options   |
|---------------------|---|
| Memory Interleaving | <ul style="list-style-type: none"> <li><b>Auto/Enable:</b> Enables memory interleaving.</li> <li><b>Disable:</b> Allows for disabling memory interleaving. The NUMA Nodes per Socket (NPS) setting will be honored regardless of this setting.</li> </ul> |

Table 3-10: Memory interleaving settings

## 3.4 Infinity Fabric Settings

This section discusses BIOS settings related to AMD Infinity Fabric technology.

### 3.4.1 Link Speed

Lowering the link speed decreases cross-socket bandwidth and increases cross-socket latency but can also save uncore power (CPU power not consumed by the cores) to either:

- Increase core frequency.
- Reduce overall power consumption.

| Setting               | Options   |
|-----------------------|---|
| 4 Link xGMI Max Speed | <ul style="list-style-type: none"> <li><b>25 Gbps/Auto</b></li> <li>Additional options depend on the OEM platform: 20, 25, 32 Gbps</li> </ul> |

Table 3-11: Link speed settings

### 3.4.2 xGMI Link Width Management

xGMI Dynamic Link Width Management saves power during periods of low socket-to-socket data traffic by reducing the number of active xGMI lanes per link from 16 to 8, or x4 if the xGMI links have limited traffic. Latency may increase in some scenarios involving low-bandwidth, latency-sensitive traffic as the processor transitions from a low-power xGMI state to full-power xGMI state. Setting **xGMI Link Width Control** to **Manual** and specifying a **Force Link Width** eliminates any such latency jitter. Applications that are not sensitive to both socket-to-socket bandwidth and latency can use a forced link width of 8 (or 2 on certain platforms) to save power, which can divert more power to the cores for boost.

| Setting   | Options  |
|---|--|
| xGMI Link Max Speed*                              | <p>NUMA-unaware workloads may need maximum xGMI bandwidth because of extensive cross-socket communications. NUMA-aware workloads may want to minimize xGMI power because they do not have a lot of cross socket traffic and prefer to use the increased CPU boost.</p> <ul style="list-style-type: none"> <li>25 Gbps/Auto</li> <li>32 Gbps</li> </ul> |
| xGMI Link Width Control*                          | <ul style="list-style-type: none"> <li><b>Auto:</b> Hide the <b>Max Link Width</b> and <b>Force Link Width</b> control options.</li> <li><b>Manual:</b> Show <b>Max Link Width</b> and <b>Force Link Width</b> control options.</li> </ul>   |
| xGMI Max Link Width*                              | <ul style="list-style-type: none"> <li><b>0:</b> Max width x4 (not supported on all platforms).</li> <li><b>1:</b> Max width x8.</li> <li><b>2:</b> Max width x16.</li> </ul>  |
| xGMI Max Link Width Control*                      | <ul style="list-style-type: none"> <li><b>Auto:</b> Hide the <b>xGMI Max Link Width</b> control.</li> <li><b>Manual:</b> Show the <b>xGMI Max Link Width</b> control.</li> </ul>   |
| xGMI Force Link Width Control*                    | <ul style="list-style-type: none"> <li><b>Unforce:</b> Use automatic xGMI Link Width selection.</li> <li><b>Force:</b> Use the xGMI Force Link Width link width.</li> </ul>  |
| xGMI Force Link Width*                            | <ul style="list-style-type: none"> <li><b>0:</b> Force xGMI link width to x4 (not supported on all platforms).</li> <li><b>1:</b> Force xGMI link width to x8.</li> <li><b>2:</b> Force xGMI link width to x16.</li> </ul>   |
| *This option only applies to dual-socket systems. |  |

Table 3-12: DLWM settings

### 3.4.3 Power States

Enable or disable Algorithm Performance Boost (APB). By default, the AMD Infinity Fabric selects between a full- and low-power fabric clock and memory clock based on usage. Latency may increase in some scenarios involving low-bandwidth, latency-sensitive traffic as the processor transitions from low to full power. Setting **APBDIS** to 1 (APB disabled) and specifying a fixed Infinity Fabric P-state of 0 forces the AMD Infinity Fabric and memory controllers into full-power mode and significantly reduces latency jitters.

| Setting                | Options   |
|------------------------|---|
| APB Disable (APBDIS)   | <ul style="list-style-type: none"> <li><b>0:</b> Dynamically switch Infinity Fabric P-state based on link usage.</li> <li><b>1:</b> Enable fixed Infinity Fabric P-state control.</li> </ul>  |
| DfPstate               | DfPstate index to set below when APBDIS [1]. <ul style="list-style-type: none"> <li><b>Min Value:</b> 0 (default); highest-performing AMD Infinity Fabric P-state.</li> <li><b>Max Value:</b> 2</li> <li><b>Pn:</b> Next-highest-performing AMD Infinity Fabric P-state.</li> </ul> |
| DfPstate Range Support | DF Pstate selection is overridden by the APB_DIS BIOS option if it is selected. If this feature is enabled, then range value setting should follow the rule that $\text{MaxDfPstate} \leq \text{MinDfPstate}$ . Otherwise, it will not work.  |

Table 3-13: Power state settings

### 3.4.4 DF C-States

Much like CPU cores, the AMD Infinity Fabric can enter lower-power states while idle, but a delay occurs when transitioning back to full-power mode that causes some latency jitter. Disabling this feature for workloads requiring low latency and/or bursty I/O will increase both performance and power consumption.

| Setting     | Options   |
|-------------|---|
| DF C-states | <ul style="list-style-type: none"> <li><b>Auto/Enabled:</b> Allow the AMD Infinity Fabric to enter a low-power state.</li> <li><b>Disabled:</b> Prevent the AMD Infinity Fabric from entering a low-power state.</li> </ul> |

Table 3-14: C-state settings

## 3.5 PCIe, I/O, Security, and Virtualization Settings

### 3.5.1 APIC Settings

Interrupt delivery is generally faster when using x2APIC compared to the legacy xAPIC mode, but not all operating systems include AMD x2APIC support. AMD recommends this mode if your OS supports it, including for configurations with fewer than 256 logical processors.

| Setting         | Options  |
|-----------------|--|
| Local APIC Mode | <ul style="list-style-type: none"> <li>• <b>APIC:</b> Use xAPIC, which supports up to 255 cores.</li> <li>• <b>x2APIC:</b> Supports more than 255 cores.</li> <li>• <b>Auto:</b> The system will choose the mode that best fits the number of active cores in the system.</li> <li>• <b>Compatibility:</b> Threads below 255 run in xAPIC with xAPIC ACPI structures, and threads 255 and above run in x2 mode with x2 ACPI structures.</li> <li>• <b>XApicMode (0x01):</b> Force legacy xApic mode</li> <li>• <b>X2ApicMode (0x02):</b> Force x2Apic mode independent of thread count.</li> </ul> |

Table 3-15: APIC settings

### 3.5.2 PCIe Speed PMM Control

The **PCIe Speed PMM Control** is an activity-based power management feature designed for PCIe Gen5 endpoints. After a device is trained, the controller monitors activity and adjusts the link speed accordingly. An idle PCIe Gen5 link will be turned down to the next highest available device speed and will return to the Gen5 speed when activity increases.

*Note: PCIe Gen5 devices that train using Equalization Bypass do not enable Gen3 or Gen4 speed. When idle, they will operate at Gen2 speed. Gen5 devices that train with full equalization support all speeds and will be turned down to Gen4 speed when idle.*

| Setting                | Options   |
|------------------------|---|
| PCIe PMM Speed Control | <ul style="list-style-type: none"> <li>• <b>0:</b> Dynamic link speed determined by power management functionality.</li> <li>• <b>1:</b> Static Target Link Speed (Gen4); sets the maximum idle link speed to 16 GT/s.</li> <li>• <b>Auto/2:</b> Static Target Link Speed (Gen5); sets the maximum idle link speed to 32 GT/s disabling the feature.</li> </ul> |

Table 3-16: PCIe Speed PMM Control settings



### 3.5.3 SR-IOV Settings

SR-IOV requires enabling PCIe Alternative Routing-ID interpretation (ARI) on both root complexes and endpoints. ARI devices interpret the PCI address as an 8-bit function number instead of a 3-bit function number, and the device number is implied to be 0.

| Setting                  | Options  |
|--------------------------|--|
| PCIe ARI Support [SRIOV] | <ul style="list-style-type: none"> <li>• <b>Disable:</b> Disable Alternative Routing ID interpretation.</li> <li>• <b>Enable:</b> Enable Alternative Routing ID interpretation.</li> </ul> |

Table 3-17: SR-IOV settings

### 3.5.4 PCIe Ten Bit Tag

A PCIe adapter must support 10-bit extended tags to achieve maximum PCIe Gen 5 bandwidth. This boosts adapter performance by allowing a 3x increase over the previous number of non-posted requests. Not all PCIe Gen 5 devices support 10-bit extended tags, which can cause issues during boot. Disabling this feature allows the server to boot if the adapter is having issues.

| Setting                  | Options   |
|--------------------------|---|
| PCIe Ten Bit Tag Support | <ul style="list-style-type: none"> <li>• <b>Auto/Disable:</b> Enable PCIe 10-bit tags for supported devices.</li> <li>• <b>Enable:</b> Disable PCIe 10-bit tags for all devices.</li> </ul> |

Table 3-18: PCIe 10-bit settings

### 3.5.5 Input-Output Memory Management Unit (IOMMU) Settings

Enabling the IOMMU allows devices such as the AMD EPYC processor-integrated SATA controller to present separate IRQs for each attached device instead of one IRQ for the subsystem. The IOMMU also allows operating systems to provide additional protection for DMA capable I/O devices. If you believe the IOMMU is limiting performance, then leave it enabled in BIOS and disable it via OS options (e.g., `iommu=pt` on the Linux® kernel command line). Enabling IOMMU is required when using x2APIC.

| Setting | Options  |
|---------|--|
| IOMMU   | <ul style="list-style-type: none"> <li>• <b>Enabled:</b> Enable IOMMU.</li> <li>• <b>Disabled:</b> Disable IOMMU.</li> </ul> |

Table 3-19: IOMMU settings

### 3.5.6 Transparent Secure Memory Encryption (TSME)

This feature provides hardware memory encryption of all data stored on system DIMMs that is invisible to the OS and slightly increases memory latency.

| Setting | Options   |
|---------|---|
| TSME    | <ul style="list-style-type: none"> <li>• <b>Auto/Disabled:</b> Disable transparent secure memory encryption.</li> <li>• <b>Enabled:</b> Enable transparent secure memory encryption.</li> </ul> |

Table 3-20: TSME settings

### 3.5.7 SEV, SEV-ES, and SEV-SNP

Please see the *AMD EPYC™ 8004 Cloud Infrastructure and Datacenter Design & Configuration Guide* (available from [AMD EPYC Tuning Guides](#)) for information about AMD EPYC security features.

### 3.5.8 AVIC & x2AVIC

AMD recommends enabling AVIC for virtualized deployments using platforms powered by 4th Gen AMD EPYC processors. Inter-processor interrupt (IPI) delivery without hardware acceleration is an expensive process because the hypervisor must take execution control from the guest, take a snapshot of the virtual machine control block (VMCB), inject the interrupt, and return control to the guest. AVIC provides a virtual APIC (vAPIC) backing page mapped into system memory and can be accessed directly by the guest for performance-sensitive operations such as IPI initiation and completion.

x2AVIC is an extension of the advanced virtual interrupt controller (AVIC) that supports more than 255 virtual CPUs and offers better performance than AVIC. It can also be thought of as hardware virtualization of x2APIC. With x2AVIC, the guest's local APIC hardware-assisted virtualization extends to 512 virtual CPUs. Before AMD Socket SP5 Processors, the VM needed to disable x2APIC capabilities because using them would bypass hardware AVIC and use software-emulated x2APIC. However, with x2AVIC, the guest will be able to leverage x2APIC performance advantages.

| Name   | Default | Description   |
|--------|---------|---|
| AVIC   | Enabled | <p>Advanced Virtual Interrupt Controller. AMD recommends enabling this option for virtualized deployments.</p> <ul style="list-style-type: none"> <li>• <b>Disabled:</b> Disables AVIC.</li> <li>• <b>Enabled:</b> Enables AVIC.</li> </ul>   |
| x2AVIC | Enabled | <p>x2AVIC is an extension of the advanced virtual interrupt controller. AMD recommends enabling this option for virtualized deployments where the total number of vCPUs &gt;255.</p> <p>x2AVIC support varies on platforms powered by AMD EPYC processors based on the AMD EPYC generation and OS (Linux, Windows, etc.). Please refer to respective OS support pages for help enabling x2AVIC. Linux Kernel 5.20 and onward include upstream x2AVIC support, and AMD strongly recommends upstream Linux Kernel 6.5 because this version contains several optimizations.</p> <ul style="list-style-type: none"> <li>• <b>Disabled:</b> Disables x2AVIC.</li> <li>• <b>Enabled:</b> Enables x2AVIC.</li> </ul> |

Table 3-21: AVIC and x2AVIC settings

## Chapter

## 4

# Workload-Specific BIOS Settings

Use these guidelines for general-purpose workloads. Some cases list the benchmarks used in order to better describe the workloads used to obtain the recommended settings. Default settings are used when labeled default.

## 4.1 General-Purpose Workloads

### 4.1.1 Processor Core Settings

| Setting                 | CPU Intensive  | Java Throughput | Java Latency   | Power Efficiency |
|-------------------------|----------------|-----------------|----------------|------------------|
| SMT Control             | <i>default</i> | <i>default</i>  | <i>default</i> | <i>default</i>   |
| L1 Stream HW Prefetcher | <i>default</i> | <i>default</i>  | <i>default</i> | <i>default</i>   |
| L1 Stride Prefetcher    | <i>default</i> | <i>default</i>  | <i>default</i> | <i>default</i>   |
| L1 Region Prefetcher    | <i>default</i> | <i>default</i>  | <i>default</i> | <i>default</i>   |
| L1 Burst Prefetch Mode  | <i>default</i> | <i>default</i>  | <i>default</i> | <i>default</i>   |
| L2 Stream HW Prefetcher | <i>default</i> | <i>default</i>  | <i>default</i> | <i>default</i>   |
| L2 Up/Down Prefetcher   | <i>default</i> | <i>default</i>  | <i>default</i> | <i>default</i>   |
| BoostFmaxEn             | <i>default</i> | <i>default</i>  | <i>default</i> | <i>default</i>   |
| BoostFmax               | <i>default</i> | <i>default</i>  | <i>default</i> | <i>default</i>   |

Table 4-1: Processor core settings

## 4.1.2 Power Management Settings

| Setting                 | CPU Intensive    | Java Throughput  | Java Latency     | Power Efficiency |
|-------------------------|------------------|------------------|------------------|------------------|
| Power Profile Selection | High Performance | High Performance | High Performance | Efficiency Mode  |
| Determinism Control     | Manual           | Manual           | Manual           | Manual           |
| Determinism Enable      | Power            | Power            | Power            | Power            |
| TDP Control             | Manual           | Manual           | Manual           | Manual           |
| TDP                     | OPN Max          | OPN Max          | OPN Max          | OPN Max          |
| PPT Control             | Manual           | Manual           | Manual           | Manual           |
| PPT                     | OPN Max          | OPN Max          | OPN Max          | OPN Max          |
| CPPC                    | <i>default</i>   | <i>default</i>   | <i>default</i>   | <i>default</i>   |

Table 4-2: Power efficiency settings

## 4.1.3 NUMA and Memory Settings

| Setting                           | CPU Intensive  | Java Throughput  | Java Latency   | Power Efficiency |
|-----------------------------------|----------------|------------------|----------------|------------------|
| ACPI SRAT L3 Cache as NUMA Domain | <i>default</i> | <i>default</i>   | <i>default</i> | Enabled          |
| NUMA Nodes per Socket (NPS)       | <i>default</i> | 2 (if available) | <i>default</i> | 2 (if available) |
| Memory Target Speed               | <i>default</i> | <i>default</i>   | <i>default</i> | <i>default</i>   |
| Memory Interleaving               | <i>default</i> | <i>default</i>   | <i>default</i> | <i>default</i>   |

Table 4-3: NUMA and memory settings

#### 4.1.4 Infinity Fabric Settings

| Setting                        | CPU Intensive  | Java Throughput | Java Latency   | Power Efficiency |
|--------------------------------|----------------|-----------------|----------------|------------------|
| xGMI Link Max Speed*           | <i>default</i> | <i>default</i>  | <i>default</i> | <i>default</i>   |
| xGMI Link Width Control*       | <i>default</i> | <i>default</i>  | <i>default</i> | <i>default</i>   |
| xGMI Max Link Width*           | <i>default</i> | <i>default</i>  | <i>default</i> | <i>default</i>   |
| xGMI Max Link Width Control*   | <i>default</i> | <i>default</i>  | <i>default</i> | <i>default</i>   |
| xGMI Force Link Width Control* | <i>default</i> | <i>default</i>  | <i>default</i> | <i>default</i>   |
| xGMI Force Link Width*         | <i>default</i> | <i>default</i>  | <i>default</i> | <i>default</i>   |
| APBDIS                         | <i>default</i> | <i>default</i>  | <i>default</i> | 1                |
| DF C-States                    | <i>default</i> | <i>default</i>  | <i>default</i> | <i>default</i>   |

Table 4-4: Infinity Fabric DP settings

\*This option only applies to dual-socket systems.

#### 4.1.5 I/O Settings

| Setting                  | CPU Intensive  | Java Throughput | Java Latency   | Power Efficiency |
|--------------------------|----------------|-----------------|----------------|------------------|
| Local APIC Mode          | <i>default</i> | <i>default</i>  | <i>default</i> | <i>default</i>   |
| PCIe Speed PMM Control   | <i>default</i> | <i>default</i>  | <i>default</i> | <i>default</i>   |
| PCIe ARI Support [SRIOV] | <i>default</i> | <i>default</i>  | <i>default</i> | <i>default</i>   |
| PCIe Ten Bit Tag Support | <i>default</i> | <i>default</i>  | <i>default</i> | <i>default</i>   |
| IOMMU                    | <i>default</i> | <i>default</i>  | <i>default</i> | <i>default</i>   |
| TSME                     | <i>default</i> | <i>default</i>  | <i>default</i> | <i>default</i>   |

Table 4-5: I/O settings

## 4.2 Memory and I/O Intensive Workloads

### 4.2.1 Processor Core Settings

| Setting                 | Memory Throughput | Storage I/O Throughput | NIC Throughput |
|-------------------------|-------------------|------------------------|----------------|
| SMT Control             | <i>default</i>    | <i>default</i>         | Disabled       |
| L1 Stream HW Prefetcher | <i>default</i>    | <i>default</i>         | <i>default</i> |
| L1 Stride HW Prefetcher | <i>default</i>    | <i>default</i>         | <i>default</i> |
| L1 Region Prefetcher    | <i>default</i>    | <i>default</i>         | <i>default</i> |
| L2 Stream HW Prefetcher | <i>default</i>    | <i>default</i>         | <i>default</i> |
| L2 Stream HW Prefetcher | <i>default</i>    | <i>default</i>         | <i>default</i> |
| L2 Up/Down Prefetcher   | <i>default</i>    | <i>default</i>         | <i>default</i> |
| BoostFmaxEn             | <i>default</i>    | <i>default</i>         | <i>default</i> |
| BoostFmax               | <i>default</i>    | <i>default</i>         | <i>default</i> |

Table 4-6: Processor core settings

### 4.2.2 Power Management Settings

| Setting                 | Memory Throughput | Storage I/O Throughput | NIC Throughput |
|-------------------------|-------------------|------------------------|----------------|
| Power Profile Selection | <i>default</i>    | <i>default</i>         | <i>default</i> |
| Determinism Control     | Manual            | Manual                 | Manual         |
| Determinism Enable      | Power             | Power                  | Performance    |
| TDP Control             | <i>default</i>    | <i>default</i>         | <i>default</i> |
| TDP                     | OPN Max           | <i>default</i>         | <i>default</i> |
| PPT Control             | <i>default</i>    | <i>default</i>         | <i>default</i> |
| PPT                     | OPN Max           | <i>default</i>         | <i>default</i> |
| CPPC                    | <i>default</i>    | <i>default</i>         | <i>default</i> |

Table 4-7: Power efficiency settings

### 4.2.3 NUMA and Memory Settings

| Setting                           | Memory Throughput | Storage I/O Throughput | NIC Throughput |
|-----------------------------------|-------------------|------------------------|----------------|
| ACPI SRAT L3 Cache as NUMA Domain | Enabled           | <i>default</i>         | <i>default</i> |
| NUMA Nodes per Socket (NPS)       | 2 (if available)  | <i>default</i>         | NPS1           |
| Memory Target Speed               | <i>default</i>    | <i>default</i>         | <i>default</i> |
| Memory Interleaving               | <i>default</i>    | <i>default</i>         | <i>default</i> |

Table 4-8: NUMA and memory settings

### 4.2.4 Infinity Fabric Settings

| Setting                        | Memory Throughput | Storage I/O Throughput | NIC Throughput |
|--------------------------------|-------------------|------------------------|----------------|
| xGM3I Link Max Speed*          | <i>default</i>    | <i>default</i>         | <i>default</i> |
| xGMI Link Width Control*       | <i>default</i>    | <i>default</i>         | <i>default</i> |
| xGMI Max Link Width*           | <i>default</i>    | <i>default</i>         | <i>default</i> |
| xGMI Max Link Width Control*   | <i>default</i>    | <i>default</i>         | <i>default</i> |
| xGMI Force Link Width Control* | <i>default</i>    | <i>default</i>         | <i>default</i> |
| xGMI Force Link Width*         | <i>default</i>    | <i>default</i>         | <i>default</i> |
| APBDIS                         | <i>default</i>    | 1                      | 1              |
| DF C-States                    | <i>default</i>    | <i>default</i>         | <i>default</i> |

Table 4-9: Infinity Fabric settings

\*This option only applies to dual-socket systems.

### 4.2.5 I/O Settings

| Setting                  | Memory Throughput | Storage I/O Throughput | NIC Throughput |
|--------------------------|-------------------|------------------------|----------------|
| Local APIC Mode          | <i>default</i>    | <i>default</i>         | 0x2 x2APIC     |
| PCIe Speed PMM Control   | <i>default</i>    | Static TLS Gen5        | Auto           |
| PCIe ARI Support [SRIOV] | <i>default</i>    | <i>default</i>         | <i>default</i> |
| PCIe Ten Bit Tag Support | <i>default</i>    | <i>default</i>         | <i>default</i> |
| IOMMU                    | <i>default</i>    | <i>default</i>         | <i>default</i> |
| TSME                     | <i>default</i>    | <i>default</i>         | <i>default</i> |

Table 4-10: I/O settings

## 4.3 Virtualization and Containers

### 4.3.1 Processor Core Settings

| Setting                 | VMware vSphere Optimized | OpenStack | Containers |
|-------------------------|--------------------------|-----------|------------|
| SMT Control             | Enabled                  | Enabled   | Enabled    |
| L1 Stream HW Prefetcher | default                  | default   | default    |
| L1 Stride HW Prefetcher | default                  | default   | default    |
| L1 Region Prefetcher    | default                  | default   | default    |
| L1 Burst Prefetch Mode  | default                  | default   | default    |
| L2 Stream HW Prefetcher | default                  | default   | default    |
| L2 Up/Down Prefetcher   | default                  | default   | default    |
| BoostFmaxEn             | default                  | default   | default    |
| BoostFmax               | default                  | default   | default    |

Table 4-11: Processor core settings

### 4.3.2 Power Management Settings

| Setting                 | VMware vSphere Optimized | OpenStack   | Containers |
|-------------------------|--------------------------|-------------|------------|
| Power Profile Selection | default                  | default     | default    |
| Determinism Control     | default                  | Manual      | default    |
| Determinism Enable      | default                  | Performance | default    |
| TDP Control             | default                  | default     | default    |
| TDP                     | default                  | default     | default    |
| PPT Control             | default                  | default     | default    |
| PPT                     | default                  | default     | default    |
| CPPC                    | default                  | default     | default    |

Table 4-12: Power efficiency settings



### 4.3.3 NUMA and Memory Settings

| Setting                           | VMware vSphere Optimized | OpenStack      | Containers     |
|-----------------------------------|--------------------------|----------------|----------------|
| ACPI SRAT L3 Cache as NUMA Domain | <i>default</i>           | Enabled        | <i>default</i> |
| NUMA Nodes per Socket (NPS)       | <i>default</i>           | <i>default</i> | <i>default</i> |
| Memory Target Speed               | <i>default</i>           | <i>default</i> | <i>default</i> |
| Memory Interleaving               | <i>default</i>           | <i>default</i> | <i>default</i> |

Table 4-13: NUMA and memory settings

### 4.3.4 Infinity Fabric Settings

| Setting                        | VMware vSphere Optimized | OpenStack      | Containers     |
|--------------------------------|--------------------------|----------------|----------------|
| xGMI Link Max Speed*           | <i>default</i>           | <i>default</i> | <i>default</i> |
| xGMI Link Width Control*       | <i>default</i>           | <i>default</i> | <i>default</i> |
| xGMI Max Link Width*           | <i>default</i>           | <i>default</i> | <i>default</i> |
| xGMI Max Link Width Control*   | <i>default</i>           | <i>default</i> | <i>default</i> |
| xGMI Force Link Width Control* | <i>default</i>           | <i>default</i> | <i>default</i> |
| xGMI Force Link Width*         | <i>default</i>           | <i>default</i> | <i>default</i> |
| APBDIS                         | <i>default</i>           | <i>default</i> | <i>default</i> |
| DF C-states                    | <i>default</i>           | <i>default</i> | <i>default</i> |

Table 4-14: Infinity Fabric settings

\*This option only applies to dual-socket systems.

### 4.3.5 I/O Settings

| Setting                  | VMware vSphere Optimized | OpenStack      | Containers     |
|--------------------------|--------------------------|----------------|----------------|
| Local APIC Mode          | <i>default</i>           | <i>default</i> | <i>default</i> |
| PCIe Speed PMM Control   | <i>default</i>           | <i>default</i> | <i>default</i> |
| PCIe ARI Support [SRIOV] | <i>default</i>           | <i>default</i> | <i>default</i> |
| PCIe Ten Bit Tag Support | <i>default</i>           | <i>default</i> | <i>default</i> |
| IOMMU                    | <i>default</i>           | <i>default</i> | <i>default</i> |
| TSME                     | <i>default</i>           | <i>default</i> | <i>default</i> |

Table 4-15: I/O settings

## 4.4 OS, Analytics, and Telco Settings

### 4.4.1 Processor Core Settings

| Setting                 | Linux KVM Optimized | IoT Gateway    | Telco (Core)   | Telco (Edge/vRAN) |
|-------------------------|---------------------|----------------|----------------|-------------------|
| SMT Control             | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| L1 Stream HW Prefetcher | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| L1 Stride HW Prefetcher | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| L1 Region Prefetcher    | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| L1 Burst Prefetch Mode  | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| L2 Stream HW Prefetcher | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| L2 Up/Down Prefetcher   | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| BoostFmaxEn             | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| BoostFmax               | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |

Table 4-16: Processor core settings

### 4.4.2 Power Management Settings

| Setting                 | Linux KVM Optimized | IoT Gateway    | Telco (Core)   | Telco (Edge/vRAN) |
|-------------------------|---------------------|----------------|----------------|-------------------|
| Power Profile Selection | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| Determinism Control     | <i>default</i>      | <i>default</i> | Manual         | Manual            |
| Determinism Enable      | <i>default</i>      | <i>default</i> | Power          | Power             |
| TDP Control             | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| TDP                     | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| PPT Control             | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| PPT                     | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| CPPC                    | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |

Table 4-17: Power efficiency settings

### 4.4.3 NUMA and Memory Settings

| Setting                           | Linux KVM Optimized | IoT Gateway    | Telco (Core)   | Telco (Edge/vRAN) |
|-----------------------------------|---------------------|----------------|----------------|-------------------|
| ACPI SRAT L3 Cache as NUMA Domain | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| NUMA Nodes per Socket (NPS)       | <i>default</i>      | <i>default</i> | NPS1           | NPS1              |
| Memory Target Speed               | Auto                | Auto           | Auto           | Auto              |
| Memory Interleaving               | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |

Table 4-18: NUMA and memory settings

### 4.4.4 Infinity Fabric Settings

| Setting                        | Linux KVM Optimized | IoT Gateway    | Telco (Core)   | Telco (Edge/vRAN) |
|--------------------------------|---------------------|----------------|----------------|-------------------|
| xGMI Link Max Speed*           | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| xGMI Link Width Control*       | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| xGMI Max Link Width*           | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| xGMI Max Link Width Control*   | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| xGMI Force Link Width Control* | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| xGMI Force Link Width*         | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| APBDIS                         | <i>default</i>      | <i>default</i> | 1              | 1                 |
| DF C-States                    | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |

Table 4-19: Infinity Fabric settings

\*This option only applies to dual-socket systems.

### 4.4.5 I/O Settings

| Setting                  | Linux KVM Optimized | IoT Gateway    | Telco (Core)   | Telco (edge/vRAN) |
|--------------------------|---------------------|----------------|----------------|-------------------|
| Local APIC Mode          | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| PCIe Speed PMM Control   | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| PCIe ARI Support [SRIOV] | <i>default</i>      | <i>default</i> | Enabled        | Enabled           |
| PCIe Ten Bit Tag Support | <i>default</i>      | <i>default</i> | Enabled        | Enabled           |
| IOMMU                    | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |
| TSME                     | <i>default</i>      | <i>default</i> | <i>default</i> | <i>default</i>    |

Table 4-20: I/O settings



*This page intentionally left blank.*

## Chapter

## 5

## Processor Identification

Figure 5-1 shows the processor naming convention for AMD EPYC 8004 Series Processors and how to use this convention to identify particular processors models:

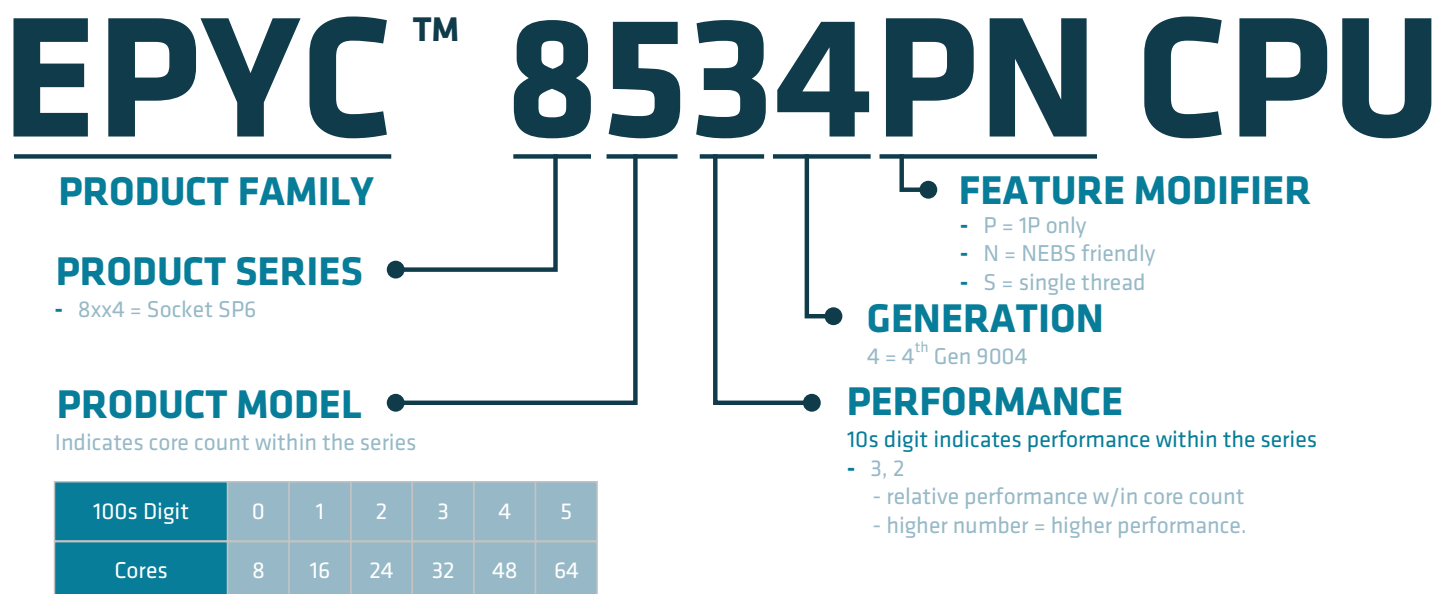


Figure 5-1: AMD EPYC SoC naming convention

## 5.1 CPUID Instruction

Software uses the `CPUID` instruction (`Fn0000_0001_EAX`) to identify the processor and will return the following values:

- **Family:** 19h identifies the “Zen 4” architecture
- **Model:** Varies with product. For example, EPYC Model 10h corresponds to an “A” part “Zen 4” CPU.
  - **8xx4:** Models A0h–AFh
- **Stepping:** May be used to further identify minor design changes

For example, `CPUID` values for Family, Model, and Stepping (decimal) of 25, 17, 1 correspond to a “B1” part “Zen 4” CPU.

## 5.2 New Software-Visible Features

AMD EPYC 8004 Series Processors introduce several new features that enhance performance, ISA updates, provide additional security features, and improve system reliability and availability. Some of the new features include:

- 5-level Paging
- AVX-512 instructions on a 256-bit datapath, including BFLOAT16 and VNNI support.
- Fast Short Rep STOSB and Rep CMPSB

Not all operating systems or hypervisors support all features. Please refer to your OS or hypervisor documentation for specific releases to identify support for these features.

Please also see the latest version of the [AMD64 Architecture Programmer's Manuals](#) or [Processor Programming Reference \(PPR\) for AMD Family 19h](#).

### 5.2.1 AVX-512

AVX-512 is a set of individual instructions supporting 512-bit register-width data (i.e., single instruction, multiple data [SIMD]) operations. AMD EPYC 8004 Series Processors implement AVX 512 by “double-pumping” 256-bit-wide registers. AMD's AVX-512 design uses the same 256-bit data path that exists throughout the Zen4 core and enables the two parts to execute on sequential clock cycles. This means that running AVX-512 instructions on AMD EPYC 8004 Series will cause neither drops on effective frequencies nor increased power consumption. On the contrary, many workloads run more energy-efficiently on AVX-512 than on AVX-256P.

Other AVX-512 support includes:

- Vectorized Neural Network Instruction (VNNI) instructions that are used in deep learning models and accelerate neural network inferences by providing hardware support for convolution operations.
- Brain Floating Point 16-bit (BFLOAT16) numeric format. This format is used in Machine Learning applications that require high performance but must also conserve memory and bandwidth. BFLOAT16 support doubles the number of SIMD operands over 32-bit single precision FP, allowing twice the amount of data to be processed using the same memory bandwidth. BFLOAT16 values mantissa dynamic range at the expense of one radix point.

## Chapter

## 6

# Debugging BIOS Setting Changes

Some BIOS default settings can be changed from the command line by a user with adequate privileges. Some of these settings take effect immediately while others may require rebooting. Some settings may not be available in the BIOS but can be set at the OS level. Some OS-level settings can either override or modify the expected BIOS default behavior. This section lists some of the settings that could have a significant impact on system performance and response times.

*Note: For each of the following examples, if you are running Windows, then please see the Microsoft® Windows® Server Tuning Guide for AMD EPYC™ 8004 Series Processors (available from [AMD EPYC Tuning Guides](#)).*

## NUMA Node Configurations:

The number of NUMA nodes, association of nodes to the memory channels, and binding of processes to specific NUMA nodes and memory nodes plays a vital role in many deployments. The **NPS** and **L3 as NUMA** settings control the number of NUMA nodes and memory channels associated with the NUMA nodes. To do this, a user leverages the NUMA command options as defined for a given OS and platforms and expects the processes to bind to certain CPUs and memory as per default BIOS settings. However, many OS-level commands and daemons may, if enabled, alter this expected behavior. The following methodology is recommended to understand and debug NUMA related issues:

### Understand NUMA topology:

Verify that the number of NUMA nodes and topology are correct. Linux has many commands such as `lstopo`, `hwloc`, and `numactl`.

### Check daemons and services that could alter default behavior:

Each OS and platform provides various tools and services that attempt to optimize a given system. There are daemons and services intended to change or allocate CPU assignments and memory to leverage locality. This can sometimes result in unintended behavior where NUMA-assigned CPUs and memory may migrate after the initial application launch. For example, `numad` is a Linux system daemon that monitors NUMA topology and resource usage. It will attempt to locate processes for efficient NUMA locality and affinity by dynamically adjusting to changing system conditions. Enabling this feature may interfere with and override the initial launch commands of a deployment like `physcpubind`, `membind`, etc. If user is observing unexpected process threads and memory migrations, check if this daemon is enabled.

### Check NUMA optimization policies:

Most platform attempt to allocate CPU and memory resources optimal NUMA leverage. Some OS and virtualization platform may have many available policy options. For example Linux automatic NUMA balancing moves tasks (threads or processes) closer to the memory they are accessing. Most applications and deployments benefit from being close to memory, meaning that this feature is often be enabled by default. NUMA balancing can have undesired effects. Further, a user's ability to bind the process CPUs and memory to different NUMA nodes may cause this setting to interfere with expected behavior. Execute the following command to disable NUMA balancing:

```
echo 0 > /proc/sys/kernel/numa_balancing
```

**OS-level settings that are not available in BIOS:**

An OEM platform may not provide BIOS options to change certain settings, but users may have been able to change them at the OS level. **AMD EPYC Core C6 (CC6) States** (alternately named C2 at the OS level) is one such example.

4th Gen AMD EPYC processors have C-States associated to cores and the Infinity Data fabric (DF). Disabling processor core C-States is highly discouraged. The system BIOS includes options to disable DF-C States for low latency and jitter-sensitive use cases. You can execute the following command to disable the Core C6 (CC6) state for the for all of the CPUs in a given system:

```
cpupower idle-set -d 2
```

You can also selectively disable Core C6. For a dual-socket system with 96-core processors, use 0-191 in the command to disable C2 for all 192 cores by executing the following command:

```
cpupower -c 0-191 idle-set -d 2.
```