

The Lenovo logo is displayed in white text on a black rectangular background.

Tuning UEFI Settings for Performance and Energy Efficiency on Intel Xeon Scalable Processor-Based ThinkSystem Servers

Covers the ThinkSystem servers with 1st, 2nd and 3rd generation Intel Xeon Scalable processors

Defines UEFI preset operating modes for performance and efficiency

Provides detailed information on key UEFI tuning parameter settings

Includes hidden UEFI parameters for fine tuning performance of workloads and hardware configurations

John Encizo

Kin Huang

Joe Jakubowski

Robert R. Wolford



Abstract

Properly configuring UEFI parameters in a server is important for achieving a desired outcome such as maximum performance or maximum energy efficiency. This paper defines preset UEFI operating modes for Lenovo® ThinkSystem™ servers running first, second and third generation Intel Xeon Scalable processors that are optimized for each of these outcomes. In addition, a thorough explanation is provided for each of the UEFI parameters so they can be fine-tuned for any particular customer deployment.

This paper is for customers and for business partners and sellers wishing to understand how to optimize UEFI parameters to achieve maximum performance or energy efficiency of Lenovo ThinkSystem servers with Intel Xeon Scalable processors.

At Lenovo Press, we bring together experts to produce technical publications around topics of importance to you, providing information and best practices for using Lenovo products and solutions to solve IT challenges.

See a list of our most recent publications at the Lenovo Press web site:

<http://lenovopress.com>

Do you have the latest version? We update our papers from time to time, so check whether you have the latest version of this document by clicking the **Check for Updates** button on the front page of the PDF. Pressing this button will take you to a web page that will tell you if you are reading the latest version of the document and give you a link to the latest if needed. While you're there, you can also sign up to get notified via email whenever we make an update.

Contents

Introduction	3
Summary of operating modes	3
How to use OneCLI and Redfish to access these settings	7
UEFI menu items	8
Hidden UEFI Items	31
Red Hat Enterprise Linux (RHEL) and derivatives	37
Additional References	38
Authors	38
Notices	40
Trademarks	41

Introduction

The Lenovo ThinkSystem UEFI provides an interface to the server firmware that controls boot and runtime services. The server firmware contains numerous tuning parameters that can be set through the UEFI interface. These tuning parameters can affect all aspects of how the server functions and also how well the server performs.

Lenovo ThinkSystem UEFI offers four convenient operating modes that set pre-defined tuning parameters:

- ▶ Minimal Power
- ▶ Efficiency - Favor Power
- ▶ Efficiency - Favor Performance
- ▶ Maximum Performance

This paper describes the tuning parameter settings for each operating mode and other tuning parameters to consider for performance and efficiency. The servers that this paper covers are as follows:

- ▶ ThinkSystem servers with 1st Gen Intel Xeon Scalable (Skylake) processors
- ▶ ThinkSystem servers with 2nd Gen Intel Xeon Scalable (Cascade Lake) processors
- ▶ ThinkSystem servers with 3rd Gen Intel Xeon Scalable (Cooper Lake and Ice Lake) processors

Note: Some settings are removed or are only valid for servers with 3rd Gen Intel Xeon Scalable processors. Except where noted, we are referring to both Cooper Lake (4-socket) and Ice Lake (2-socket) processors.

Summary of operating modes

The ThinkSystem servers with Intel processors offer four preset operating modes, Minimal Power, Efficiency – Favor Power, Efficiency – Favor Performance and Maximum Performance. These modes are a collection of predefined low-level UEFI settings that simplify the task of tuning the server for either minimal power, maximum performance or efficiency.

The four pre-defined modes are as follows:

- ▶ Minimal Power: Minimize the absolute power consumption of the system.
- ▶ Efficiency – Favor Power: Maximize the performance/watt efficiency with a bias towards power savings. It is expected that will be the favored mode for SPECpower benchmark testing.
- ▶ Efficiency – Favor Performance (default): Maximize the performance/watt efficiency with a bias towards performance. It is the favored mode for Energy Star certification.
- ▶ Maximum Performance: Achieves maximum performance but with higher power consumption and lower energy efficiency.

Table 1 on page 4 summarizes the settings that are made for each mode selected for the ThinkSystem server Intel platforms. Table 2 on page 5 summarizes additional settings not in the preset operating modes that can be tuned for performance and efficiency. The values in the Category column (column 3) in each table are as follows:

- ▶ Recommended: Settings follow Lenovo's best practices and should not be changed without sufficient justification.

- ▶ Suggested: Settings follow Lenovo’s general recommendation for a majority of workloads but these settings can be changed if justified by workload specific testing.
- ▶ Test: The non-default values for the Test settings can optionally be evaluated because they are workload dependent.

Table 1 UEFI Settings for operating modes on ThinkSystem servers with Intel processors.

Menu Item	Page	Category	Minimal Power	Efficiency – Favor Power	Efficiency – Favor Performance	Maximum Performance
Operating Mode	9	Suggested	Minimal power	Efficiency – Favor Power	Efficiency – Favor Performance	Maximum Performance
CPU P-State Control	10	Recommended	Autonomous	Autonomous	Autonomous	None
C-States	11	Recommended	Autonomous	Autonomous	Autonomous	Disabled
MONITOR/MWAIT	12	Suggested	Enabled	Enabled	Enabled	Disabled
C1E Enhanced Mode	12	Recommended	Enabled	Enabled	Enabled	Disabled
UPI Link Frequency	13	Recommended	Minimal power	Minimal Power	Max Performance	Max performance
UPI Link Disable	13	Recommended	Minimum Number of Links Enabled	Minimum Number of Links Enabled	Enabled All Links	Enabled all links
UPI Power Management	14	Recommended	L1: Enabled L0p: Enabled	L1: Enabled L0p: Enabled	L1: Enabled L0p: Enabled	L1: Disabled L0p: Disabled
Turbo Mode	14	Suggested	Disabled	Disabled	Enabled	Enabled
Energy Efficient Turbo	15	Recommended	Enabled	Enabled	Enabled	Disabled
Power/ Performance Bias	15	Suggested	Platform Controlled	Platform Controlled	Platform Controlled	Platform Controlled
Platform Controlled Type	16	Suggested	Minimal Power	Efficiency – Favor Power	Efficiency – Favor Performance	Maximum Performance
Memory Speed	17	Recommended	Minimal power	Balanced	Max Performance	Max performance
Memory Power Management	17	Recommended	Automatic	Automatic	Disabled	Disabled
Page Policy	18	Recommended	Adaptive	Adaptive	Adaptive	Adaptive
ADDDC Sparing	18	Recommended	Disabled	Disabled	Disabled	Disabled

Table 2 lists additional UEFI settings that you should consider for tuning for performance or energy efficiency. These settings are not part of the preset operating modes.

Table 2 Other UEFI settings to consider for performance and efficiency

Menu Item	Page	Category	Comments
Hyper-Threading	19	Suggested	This setting allows two separate instruction streams to run simultaneously on each processor core. Most applications benefit from hyperthreading so this parameter should be left Enabled unless it is known that performance degrades with a specific application. Applications that can fully utilize a core with a single instruction stream should set hyperthreading to Disabled.
Cores in CPU Package	20	Recommended	This setting logically powers off a set number of cores for each processor in a system.
CPU Frequency Limits	20	Suggested / Test	Enable one of the CPU Frequency Limit options to restrict the turbo frequency increase. This will minimize performance variability among a population of servers caused by core frequency jitter at the expense of achieving peak performance and maximum turbo frequency upside. This may be a better option than disabling turbo mode.
Processors X to Y Cores Active	21	Test	Related to CPU Frequency Limits and used for restricting the maximum core frequency.
Hardware Prefetcher	22	Recommended	Also known as the MLC Streamer Prefetcher. When enabled, this parameter fetches the next cache line from memory into the processor's L2 cache if two consecutive cache lines were read. This parameter should be Enabled unless experiments have been run selectively disabling the prefetchers one at a time and performance is improved when Disabled.
Adjacent Cache Prefetch	22	Recommended	Also known as the MLC Spatial Prefetcher. When enabled, this parameter fetches both cache lines that make up a 128 byte cache line pair even if the requested data is only in the first cache line. This parameter should be Enabled unless experiments have been run selectively disabling the prefetchers one at a time and performance is improved when Disabled.
DCU Streamer Prefetcher	22	Recommended	When enabled, this parameter fetches the next cache line into the L1 data cache when multiple loads from the same cache line are executed in a certain time limit. This parameter should be Enabled unless experiments have been run selectively disabling the prefetchers one at a time and performance is improved when Disabled.
DCU IP Prefetcher	22	Recommended	When enabled, this parameter fetches the next cache line into the L1 data cache if there is a sequential load history of cache line accesses. This parameter should be Enabled unless experiments have been run selectively disabling the prefetchers one at a time and performance is improved when Disabled.
XPT Prefetcher	23	Recommended	Recommend this parameter remain Enabled unless experiments have been run showing performance is improved when Disabled.
UPI Prefetcher	23	Suggested	Recommend this parameter remain Enabled unless experiments have been run showing performance is improved when Disabled.
Direct Cache Access (DCA)	24	Suggested	Experiments can be run with this parameter Enabled or Disabled depending on whether the IO device supports Direct Cache Access.

Menu Item	Page	Category	Comments
Uncore Frequency Scaling (UFS)	24	Suggested / Test	For Intel Xeon Scalable Processors: When Disabled, the processor uncore frequency is set to the maximum speed. When Enabled, the uncore frequency is allowed to float depending on the load to the uncore devices, e.g. memory and IO operations. Some workloads perform better with UFS disabled. User can change Uncore Frequency Limit via OneCLI variable Processors.UncoreFrequencyLimit. Refer to "Uncore Frequency Limit" on page 37 for details.
Sub-NUMA Clustering (SNC)	25	Suggested / Test	SNC can be Enabled for applications that permit and can benefit from a high degree of NUMA optimization. Applications that are not NUMA-friendly should leave SNC Disabled.
Snoop Preference	26	Suggested / Test	The default Home Snoop Plus setting provides the highest performance for most application workloads. Moderately NUMA optimized applications that require maximum local and remote bandwidth can benefit from the Home Snoop setting.
Socket Interleave	27	Recommended	Recommend this parameter be set to NUMA unless experiments have been run showing performance is improved when set to Non-NUMA.
Patrol Scrub	27	Recommended	Recommend this parameter be left Enabled for improved memory reliability. Lower memory latency can be achieved if this parameter is Disabled but with the risk of lower memory reliability.
Memory Data Scrambling	28	Recommended	Recommend this parameter remain Enabled to avoid memory data-bit errors.
Workload Configuration	16	Suggested	Tunes the system's I/O bandwidth profile between a balanced profile and an IO sensitive profile by adjusting the processor's core and uncore frequencies.
Intel Virtualization Technology	28	Suggested	Set to Enabled when the server will be configured as a virtualized host. For non-virtualized applications where high frequency and low latency code execution is required, this parameter can be Disabled.
Intel VT for Directed I/O (VT-d)	29	Suggested	Set to Enabled when the server will be configured as a virtualized host. For non-virtualized applications where low latency code execution and low jitter are required, this parameter can be Disabled.
PCIe Power Brake	29	Suggested	PCIe power brake can be set to different modes to tailor the system's response when PCIe devices draw large amounts of power. This setting is only available on servers with 3rd Gen Intel Xeon Scalable processors.
ASPM	30	Suggested	ASPM is a PCIe power saving feature. It puts the PCIe link into a low power mode when the link is idle. This setting is only available on servers with 3rd Gen Intel Xeon Scalable processors.

How to use OneCLI and Redfish to access these settings

In addition to using UEFI Setup, Lenovo also provides OneCLI/ASU variables and Redfish UEFI Setting Attribute names for managing system settings.

► OneCLI/ASU variable usage

Show current setting:

```
Onecli config show "<OneCLI/ASU Var>" --override --log 5 --imm  
<userid>:<password>@<IP Address>
```

Example:

```
onecli config show "OperatingModes.ChooseOperatingMode" --override --log 5  
--imm USERID:PASSWORD@10.240.218.89
```

Set a setting:

```
Onecli config set "<OneCLI/ASU Var>" "<choice>" --override --log 5 --imm  
<userid>:<password>@<IP Address>
```

Example:

```
onecli config set "OperatingModes.ChooseOperatingMode" "Maximum Efficiency"  
--override --log 5 --imm USERID:PASSWORD@10.240.218.89
```

► Redfish Attributes configure URL

Setting get URL: <https://<BMC IP>/redfish/v1/Systems/1/Bios>

Redfish Value Names of Attributes

If no special description, choice name is same as possible values. If there are space character (' '), dash line ('-') or slash line ('/') in the possible values, remove them. Because Redfish standard doesn't support those special characters.

If you use OneCLI to configure the setting, OneCLI will automatically remove them. But if you use other Redfish tools, then you may need to remove those characters by yourselves.

Example:

"Operating Mode" has three choices: Maximum Efficiency, Maximum Performance and Custom Mode, their Redfish value names are MaximumEfficiency, MaximumPerformance and CustomMode.

For more detailed information on the BIOS schema, please refer to the DMTF website:

https://redfish.dmtf.org/redfish/schema_index

Usually, postman can be used for get/set BIOS schema:

<https://www.getpostman.com/>

Enable in Skylake vs Enabled in Cascade Lake: In servers with 1st Gen Xeon Scalable (Skylake) processors, UEFI settings have values of "Enable" or "Disable", however in 2nd Gen Xeon Scalable (Cascade Lake) processors, these have been changed to "Enabled" and "Disabled". In this paper, we will simply use "Enabled" and "Disabled" for these items.

The remaining sections in this paper provide details about each of these settings. We describe how to access the settings via System Setup (Press F1 during system boot).

UEFI menu items

The following items are provided to server administrators in UEFI menus that are accessible by pressing F1 when a server is booted, through the XClarity Controller (XCC) service processor, or through command line utilities such as Lenovo's Advanced Settings Utility™ (ASU) or OneCLI.

These parameters are made available because they are regularly changed from their default values to fine tune server performance for a wide variety of customer use cases.

Menu items described in this paper for ThinkSystem servers with Intel Xeon Scalable processors:

- ▶ "Operating Mode" on page 9
- ▶ "Processor P-States" on page 10
- ▶ "C-States" on page 11
- ▶ "MONITOR/MWAIT" on page 12
- ▶ "C1 Enhanced Mode" on page 12
- ▶ "UPI Link Frequency" on page 13
- ▶ "UPI Link Disable" on page 13
- ▶ "UPI Power Management" on page 14
- ▶ "Turbo Mode" on page 14
- ▶ "Energy Efficient Turbo" on page 15
- ▶ "Power/Performance Bias" on page 15
- ▶ "Platform Controlled Type" on page 16
- ▶ "Workload Configuration Bias" on page 16
- ▶ "Memory Speed" on page 17
- ▶ "Memory Power Management" on page 17
- ▶ "Page Policy" on page 18
- ▶ "ADDDC Sparing" on page 18
- ▶ "Hyper-Threading" on page 19
- ▶ "Cores in CPU Package" on page 20
- ▶ "CPU Frequency Limits" on page 20
- ▶ "Processors X to Y Cores Active" on page 21
- ▶ "Processor Prefetchers" on page 22
- ▶ "XPT Prefetcher" on page 23
- ▶ "UPI Prefetcher" on page 23
- ▶ "DCA (Direct Cache Access)" on page 24
- ▶ "Uncore Frequency Scaling" on page 24
- ▶ "SNC (Processor Sub-NUMA Clustering)" on page 25
- ▶ "Snoop Preference (Processor Snoop Mode)" on page 26
- ▶ "Socket Interleave (Memory Socket Interleave)" on page 27
- ▶ "Patrol Scrub" on page 27
- ▶ "Memory Data Scrambling" on page 28
- ▶ "Intel Virtualization Technology" on page 28
- ▶ "Intel Virtualization Technology for I/O" on page 29
- ▶ "PCIe Power Brake" on page 29
- ▶ "ASPM (PCIe Active State Power Management)" on page 30

Hidden menu items:

- ▶ "CR QoS Configuration and CR FastGo Configuration" on page 31
- ▶ "L2 RFO (Request of Ownership) Prefetch Disable" on page 32
- ▶ "Local/Remote Threshold" on page 32
- ▶ "Stale AtoS" on page 34
- ▶ "Snoop Response Hold Off" on page 35

- ▶ “LLC dead line allocation” on page 36
- ▶ “LLC (Last Level Cache) Prefetch” on page 36
- ▶ “Uncore Frequency Limit” on page 37

Operating Mode

This setting is used to set multiple processor and memory variables at a macro level.

Choosing one of the predefined Operating Modes is a way to quickly set a multitude of processor, memory, and miscellaneous variables. It is less fine grained than individually tuning parameters but does allow for a simple “one-step” tuning method for two primary scenarios.

Tip: Prior to optimizing a workload for maximum performance, it is recommended to set the Operating Mode to “Maximum Performance” and then reboot rather than simply starting from the “Efficiency – Favor Performance” default mode and then modifying individual UEFI parameters. If you don’t do this, some settings may be unavailable for configuration.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Operating Modes** → **Choose Operating Mode**
- ▶ OneCLI/ASU variable: `OperatingModes.ChooseOperatingMode`
- ▶ Redfish attribute: `OperatingModes_ChooseOperatingMode`

Possible values:

- ▶ **Minimal Power**

Minimal Power mode strives to minimize the absolute power consumption of the system while it is operating. In addition, if a customer wants additional power savings, they can set a power cap or fan speed restriction. The tradeoff with the minimal power mode and/or power cap or fan speed restriction is that performance may be reduced depending on the application that is running.

- ▶ **Efficiency – Favor Power**

Efficiency - Favor Power mode maximizes the performance/watt efficiency with a bias towards power savings. It provides the best features for reducing power and increasing performance in applications where maximum bus speeds are not critical. It is expected that this will be the favored mode for SPECpower testing. This mode maintains backwards compatibility with systems that included the preset operating modes before Energy Star for servers was released.

- ▶ **Efficiency – Favor Performance (Default)**

Efficiency - Favor Performance mode optimizes the performance/watt efficiency with a bias towards performance. It is the favored mode for Energy Star. Note that this mode is slightly different than the Favor Power mode. In Favor Performance mode, no bus speeds are derated as they are in Favor Power mode.

- ▶ **Maximum Performance**

Maximum Performance mode will maximize the absolute performance of the system without regard for power. In this mode, power consumption is a “don't care”. Things like fan speed and heat output of the system may increase in addition to power consumption. Efficiency of the system may go down in this mode but the absolute performance can go up depending on the benchmark that is run.

► **Custom Mode**

Custom Mode allows the user to customize the performance & power settings. Custom Mode will inherit the UEFI settings from the previous preset operating mode. For example, if the previous operating mode was the Maximum Performance operating mode and then Custom Mode was selected, all the settings from the Maximum Performance operating mode will be inherited. Note that there are certain settings that may be mutually exclusive or interdependent. For these settings an error will be surfaced if one of the pre-requisite or interrelated settings is set in such a way as to make configuration of the setting in question non-valid.

Processor P-States

P-states (Intel Enhanced SpeedStep) dynamically adjusts core frequency and voltage dependent on processor utilization. This setting allows for a processor to have a variable core voltage which allows for either a low power or base rated clock speed frequency. P-states work by adjusting processor core voltage and frequency.

All modern operating systems have drivers that support P-states. The amount of adjustment to the core voltage and frequency (either up or down) is dependent on the type and stepping of a CPU, the P-states mode and requests from operating system drivers.

Autonomous mode is used for normal power savings and serves well for most typical business applications. Those applications which are clock frequency sensitive it is recommended to test with Cooperative or Legacy mode, while for those workloads which are latency sensitive Lenovo recommends setting the value to None.

This setting is accessed as follows:

- System setup: System **Settings** → **Processors** → **CPU P-state Control**
- OneCLI/ASU variable: Processors.CPUPstateControl
- Redfish attribute: Processors_CPUPstateControl

Possible values:

► **Autonomous** (Default)

All CPU P-state management is handled automatically in the background without any OS intervention.

► **Cooperative without Legacy**

On ThinkSystem, the setting choice is named “Cooperative”.

UEFI doesn't provide legacy P-States. OS provides hints to the processor's PCU on the desired P-state min / max levels. Requires Windows Server 2016 and Linux kernel v4.2 and higher. PCU runs in Autonomous mode until the OS sets the desired frequency.

► **Cooperative with Legacy**

This setting is available on servers with 3rd Gen Intel Xeon Scalable processors.

UEFI leaves the legacy P-states interface initially enabled until/if later an OS that is aware of Intel Hardware P-states (HWP) native mode sets the bit. Legacy P-states will be used until OS sets the HWP native mode. After that, P-states will switch to same behavior as “Cooperative without Legacy”.

- ▶ **Legacy**

Legacy control mechanisms currently implemented for systems with processors prior to the Intel Xeon Scalable Processor codenamed Skylake. Uses standard ACPI interface. Use for applications which benefit from OS level power controls.

- ▶ **None**

No ACPI table entries for P-states are created. P-states are disabled. Use this setting to minimize latency caused by P-state transitions.

C-States

This setting specifies the highest level of C-state allowed on a system. The higher a C-state is, as denoted by its C-state modifier being numerically larger (i.e., C1 has a higher power and activity state than C3) the less power is drawn from that core / CPU. Deeper (i.e., "higher") C-states achieve their power savings by progressively shutting down parts of a processor package, from idling inactive cores all the way to shutting down idle cores.

C-states are found at the package and core level. Core C-states resolve to the highest power threads running on a core, whereas package C-states resolve to the highest power C-state of all cores in a processor package. Package C-states affect the entire processor and entail a deeper low power state and even higher exit latencies than core C-states.

Package C-states are engaged when all cores on a processor package are at a C-state level \geq package C-state level requested (i.e. - all cores must be at \geq C1e for the package C-state to go to C1e).

It is important to understand that as C-states go deeper, a processor's instruction execution latency increases and the probability of cache coherency decreases.

The C-States are beneficial if either power savings are important or if used in conjunction with Turbo Mode and *bursty* applications (i.e., applications that have cyclical periods of inactivity followed by periods of large processing demands) to provide additional power savings for Turbo Mode to draw from. Cache coherent applications (Database & HPC-type applications) will suffer performance hits if they are highly threaded due to the low power states flushing processor and package caches. If an application typically utilizes less than all the CPU cores and reaching the maximum CPU frequency is important, c-states should be enabled so that the active cores can reach maximum turbo frequency.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Processors** → **C-States**
- ▶ OneCLI/ASU variable: `Processors.CStates`
- ▶ Redfish attribute: `Processors_CStates`

Possible values:

- ▶ **Legacy** (Default)

C-states are presented to OS via ACPI table entries. OSPM controls core C-state transitions. PCU controls package C-state transitions.

- ▶ **Autonomous**

This choice has been removed on servers with 3rd Gen Intel Xeon Scalable processors.

When autonomous C-states are enabled the PCU converts HALT and MWAIT (C1) instructions to MWAIT (C6). No C-state objects are created in the ACPI table.

- ▶ **Disabled**

C-states are disabled.

MONITOR/MWAIT

Some operating systems (for example, some Linux distributions) engage C-states by using MONITOR/MWAIT instructions and not the ACPI table. These operating systems will still enter higher C-states even if the C-States UEFI parameter is Disabled. To prevent this, you need to disable MONITOR/MWAIT.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Processors** → **MONITOR/MWAIT**
- ▶ OneCLI/ASU variable: Processors.MONITORMWAIT
- ▶ Redfish attribute: Processors_MONITORMWAIT

Possible values:

- ▶ **Enabled** (Default)
Enable MONITOR/MWAIT instructions support.
- ▶ **Disabled**
Disable MONITOR/MWAIT instructions support.

C1 Enhanced Mode

C1 Enhanced mode (C1E) is a processor power saving feature that halts cores not in use and maintains cache coherency. C1E maintains all of the C1 halt state functionality, but the core voltage is reduced for enhanced power savings. If all cores in a package are in C1 state, the package itself will enter C1E unless C1E is disabled.

C1E can help to provide power savings in those circumstances where cache coherency is paramount. Those applications which thread well and can maintain utilization of processor cores (virtualization, HPC and database workloads) do not benefit and under certain circumstances may be hindered by C1E. If a user is attempting to achieve maximum opportunity for Turbo Mode to engage, C1E is recommended. C1E is not recommended for latency sensitive workloads.

Note: This item is only valid if C-states are set to legacy. If C-states are set to Autonomous the C1E item is statically set to Enabled.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Processors** → **C1 Enhanced Mode**
- ▶ OneCLI/ASU variable: Processors.C1EnhancedMode
- ▶ Redfish attribute: Processors_C1EnhancedMode

Possible values:

- ▶ **Disabled**
- ▶ **Enabled** (Default)

UPI Link Frequency

UPI Link frequency determines the rate at which the UPI processor interconnect link will operate.

If a workload is highly NUMA aware, sometimes lowering the UPI link frequency can free up more power for the cores and result in better overall performance. But there is also interplay between this menu item and UPI link disable and Uncore Frequency Scaling.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Processors** → **UPI Link Frequency**
- ▶ OneCLI/ASU variable: Processors.UPILinkFrequency
- ▶ Redfish attribute: Processors_UPILinkFrequency

Possible values:

- ▶ **Maximum Performance** (Default)
Set UPI interface to rated speed for CPU SKU
- ▶ **Balanced**
Set UPI interface to N-1 of rated speed for CPU SKU
- ▶ **Minimal Power**
Set UPI interface to minimum rated speed for CPU family

UPI Link Disable

Disable one of the available UPI processor interconnect links on the CPU.

Use UPI Link disable when attempting to conserve power on a platform. Disabling a UPI Link will impact performance for any workload which does cross socket communication. Disabling a UPI Link is not recommended for high frequency, low jitter, low latency, virtualization or database workloads.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Processors** → **UPI Link Disable**
- ▶ OneCLI/ASU variable: Processors.UPILinkDisable
- ▶ Redfish attribute: Processors_UPILinkDisable

Possible values:

- ▶ **Enabled All Links** (Default)
Enable all available UPI links between CPUs.
- ▶ **Minimum Number of Links Enabled**
Enables the minimum number of UPI links required by the processor architecture. The actual number of UPI links may vary depending on the server's design, such as 2-socket compared to 4-socket mesh.

UPI Power Management

UPI power management is used when attempting to conserve power on a platform. Low power mode reduces UPI frequency and bandwidth. It is recommended for power saving, however, UPI power management is not recommended for high frequency, low jitter, low latency, virtualization or database workloads. For these workloads Lenovo recommends setting both L1 and L0p to Disabled.

L1 saves the most power but has the highest impact on latency and bandwidth. L1 takes a UPI link and allows it to transition from Full Width>x8 width>Link down. L1 is the deepest power savings state.

L0p allows for a partial link down. A subset of all of the lanes will remain awake.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Processors** → **UPI Power Management**
- ▶ OneCLI/ASU variable:
 - Processors.L1
 - Processors.L0p
- ▶ Redfish attribute:
 - Processors_L1
 - Processors_L0p

Possible values:

- ▶ **Enabled** (Default)
- ▶ **Disabled**

Turbo Mode

Turbo Mode implements Intel Turbo Boost Technology and allows the processor to opportunistically increase a set of CPU cores higher than the CPU's rated base clock speed based on the number of active cores, power and thermal headroom in a system. It is important to understand that this is not a guarantee of a CPU frequency increase, rather it is enabling the opportunity to run at a higher clock frequency. The performance of Turbo Mode increases when fewer cores are active, dynamic power management is enabled, and the system is running below the thermal design limits for the platform.

Turbo Mode is denoted as a series of numbers in the form a/b/c/d/x, where each number is a MHz increment in CPU frequency. This corresponds to the number of cores active in the package. It can generally be interpreted as:

all cores / a-1 cores / a-2 cores / a-x cores

Use Turbo Mode when you have applications which can benefit from clock frequency enhancements. Avoid using this feature with latency sensitive or clock frequency sensitive (low jitter) applications, or if power consumption is a concern.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Processors** → **Turbo Mode**
- ▶ OneCLI/ASU variable: Processors.TurboMode
- ▶ Redfish attribute: Processors_TurboMode

Possible values:

- ▶ **Enabled** (Default)
- ▶ **Disabled**

Energy Efficient Turbo

When Energy Efficient Turbo is enabled, the CPU's optimal turbo frequency will be tuned dynamically based on CPU utilization. The actual turbo frequency the CPU is set to is proportionally adjusted based on the duration of the turbo request.

Memory usage of the OS is also monitored. If the OS is using memory heavily and the CPU core performance is limited by the available memory resources, the turbo frequency will be reduced until more memory load dissipates and more memory resources become available. The power/performance bias setting also influences energy efficient turbo.

Energy Efficient Turbo is best used when attempting to maximize power consumption over performance.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Processors** → **Energy Efficient Turbo**
- ▶ OneCLI/ASU variable: `Processors.EnergyEfficientTurbo`
- ▶ Redfish attribute: `Processors_EnergyEfficientTurbo`

Possible values:

- ▶ **Enabled** (Default)
Enable Energy Efficient Turbo to optimize power consumption over performance.
- ▶ **Disabled**
Disable Energy Efficient Turbo.

Power/Performance Bias

Power/Performance Bias controls how aggressively the CPU will be power managed and placed into turbo mode. Power/Performance bias has no effect on CPU frequencies when P-states are disabled, but it can still influence CPU power management features.

Lenovo recommends enabling Platform Controlled for workloads requiring low jitter and low latency. For High Frequency workloads, it is suggested to set the parameter to Platform Controlled, however some multi-threaded applications may benefit from setting it to OS Controlled, allowing the OS to control power management due to thread scheduling allowing for non-used cores to be placed in a halt / lower power mode.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Power** → **Power/Performance Bias**
- ▶ OneCLI/ASU variable: `Power.PowerPerformanceBias`
- ▶ Redfish attribute: `Power_PowerPerformanceBias`

Possible values:

- ▶ **Platform Controlled** (Default)

The system UEFI configuration determines if and how aggressively the platform will enable power management and turbo mode.

- ▶ **OS Controlled**

The operating system controls how aggressively power management and Turbo Mode will be used via legacy ACPI calls requesting processor performance governors. Not all operating systems support OS control of power management and performance governors. In the event that an OS does not support power management instrumentation, the processor's power control unit (PCU) defaults to Maximum Performance until an OS makes power management calls.

Platform Controlled Type

Controls how aggressively the processor's Power Control Unit (PCU) will engage power management and how the CPU cores are placed into turbo mode. When set to Maximum Performance, turbo mode can be engaged opportunistically before it is requested. In addition, the lowest engagement of uncore power management features occurs.

Lenovo recommends enabling Maximum Performance for all systems where workload performance is more important than power savings.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Power** → **Platform Controlled Type**
- ▶ OneCLI/ASU variable: `Power.PlatformControlledType`
- ▶ Redfish attribute: `Power_PlatformControlledType`

Possible values:

- ▶ **Maximum Performance**

Turbo mode can be engaged opportunistically before it is requested and uncore power management features (Memory, UPI, C-state demotion, I/O bandwidth limit and UFS) are aggressively disabled.

- ▶ **Efficiency – Favor Performance** (Default)

Turbo is quickly engaged but not opportunistically. Some light uncore and UPI power management features are engaged to provide low cost power savings with an overall weighting towards per core and per package performance.

- ▶ **Efficiency - Favor Power**

Turbo is engaged more slowly but uncore and UPI power management policies are quickly engaged to balance the systems towards power savings.

- ▶ **Minimal Power**

Turbo is not engaged and aggressive power management policies on the uncore and UPI are engaged to drive lowest performance / watt.

Workload Configuration Bias

Workload configuration bias is used to tune the system's I/O bandwidth profile. This setting tunes how aggressively the system will allocate processor core and uncore frequency to handle I/O requests.

I/O sensitive workloads will benefit from this setting as it enables higher I/O bandwidth. Lenovo recommends using I/O sensitive for low latency workloads and to evaluate the impact for high frequency and low jitter workloads.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Power** → **Workload Configuration**
- ▶ OneCLI/ASU variable: `Power.WorkloadConfiguration`
- ▶ Redfish attribute: `Power_WorkloadConfiguration`

Possible values:

- ▶ **Balanced** (Default)
CPU core and uncore frequency is balanced to provide equal performance weighting between I/O tasks and application workload threads.
- ▶ **I/O sensitive**
CPU core and uncore frequency is weighted to allocate enough resources to provide high I/O bandwidth when CPU cores are at low utilization.

Memory Speed

The Memory Speed setting determines the frequency at which the installed memory will run.

Memory speed setting could be changed if the objective is to conserve power, since lowering the clock frequency to the installed memory will reduce overall power consumption of the DIMMs.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Memory** → **Memory Speed**
- ▶ OneCLI/ASU variable: `Memory.MemorySpeed`
- ▶ Redfish attribute: `Memory_MemorySpeed`

Possible values:

- ▶ **Maximum Performance** (Default)
Maximum Performance sets the memory to the maximum allowed frequency as dictated by the type of CPU installed and the memory DIMMs installed.
- ▶ **Balanced**
Balanced attempts to balance power and performance for the memory subsystem. The DIMMs will run at a frequency of N-1 of the rated frequency for the CPU and memory DIMMs installed. Memory voltage will be set to the lowest supported value to drive N-1 frequency.
- ▶ **Minimal Power**
Minimal power optimizes memory power savings over performance. Memory is set to run at the minimum support memory frequency supported by the platform. For the Intel Xeon Scalable processors (all generations), this is 1866 MHz.

Memory Power Management

Memory power management allows the platform to put the memory into a lower refresh rate, thus saving power.

Enable memory power management when attempting to maximize platform power consumption. For high frequency, low latency, low jitter, virtualization or database workloads Lenovo generally recommends disabling memory power management.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Memory** → **Memory Power Management**
- ▶ OneCLI/ASU variable: `Memory.MemoryPowerManagement`
- ▶ Redfish attribute: `Memory_MemoryPowerManagement`

Possible values:

- ▶ **Automatic** (Default)
Automatic attempts to optimize for the best performance / watt.
- ▶ **Disabled**
Maximum performance and maximum power draw. Useful for high performance, low latency workloads.

Page Policy

Page Policy setting determines whether the memory controller keeps the last accessed page open. With the Open Page policy, memory access is faster in case of a Page hit but slower in case of a Page Miss.

Lenovo recommends using the Adaptive mode page policy for low latency and high-performance workloads due to lower loaded latencies and better overall performance. Closed Page mode works well for workloads such as databases that may do many streaming reads / writes to the memory subsystem.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Memory** → **Page Policy**
- ▶ OneCLI/ASU variable: `Memory.PagePolicy`
- ▶ Redfish attribute: `Memory_PagePolicy`

Possible values:

- ▶ **Closed**
Closed page policy maps memory pages out to avoid “hot spots” on DIMMs, it is useful for workloads such as databases that may do many memory transactions.
- ▶ **Adaptive** (Default)
Adaptive page policy balances mapping out memory pages across the DIMM infrastructure to provide both performance and latency.

ADDDC Sparing

ADDDC (Adaptive Double DRAM Device Correction) Sparing is a RAS function that provides more reliability of memory error correction in virtual lockstep mode. When the memory correctable error count reaches a pre-defined threshold, ADDDC will trigger virtual lockstep mode, which can significantly reduce performance.

For high frequency, low latency, or low jitter workloads, Lenovo generally recommends disabling ADDDC sparing for high-performance.

ADDDC supported with x4 DRAM only: This setting is not available with x8 DIMMs.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Memory** → **ADDDC Sparing**
- ▶ OneCLI/ASU variable: `Memory.ADDDCSparing`
- ▶ Redfish attribute: `Memory_ADDDCSparing`

Possible values:

- ▶ **Disabled** (Default)
It is static Disable when “Page Policy” is Adaptive.
- ▶ **Enabled**
Enable ADDDC Sparing to increase memory reliability.

Hyper-Threading

Hyper-threading (HT) is the capability of a single core to execute two threads simultaneously. An OS will register a hyper-thread as a logical CPU and attempt to schedule instruction threads accordingly. All processor cache is shared between the physical core and its corresponding hyper-thread. Many operating systems and hypervisors are able to schedule instructions such that both threads execute on the same core.

Hyper-threading takes advantage of out-of-order execution, deeper execution pipelines and improved memory bandwidth in today’s processors to be an effective way of garnering many of the benefits of additional logical CPUs without having to supply the power necessary to drive a physical core.

Most workloads can improve performance with hyper-threading enabled or at least not degrade performance. Virtualization and database workloads generally benefit from HT as they can spawn threads and intelligently place them on the same physical core.

Disable hyper-threading if one instruction stream (thread) can fully utilize a physical core. HPC and scientific workloads can be adversely impacted due to sharing of processor cache resources with the hyper-thread (reducing the logical CPU to cache ratio by 50%) and the fact that a hyper-thread is still dependent on sharing the same instruction execution engines in a single core. Applications that are thread bound may also see a performance impact as they tend to have larger cache to core ratios. Disable hyper-threading for latency sensitive applications, e.g. some high frequency trading and some high-performance computing environments.

The benefit of hyper-threading is workload dependent so evaluate on a case-by-case basis.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Processors** → **Hyper-Threading**
- ▶ OneCLI/ASU variable: `Processors.HyperThreading`
- ▶ Redfish attribute: `Processors_HyperThreading`

Possible values:

- ▶ **Enabled** (Default)
- ▶ **Disabled**

Cores in CPU Package

This setting logically powers off a set number of cores for each processor in a system. As opposed to restricting the number of logical processors an OS will run on, this setting directly affects the number of cores physically powered on by turning off the core level power gates on each processor.

Manipulating the number of physically powered cores is primarily used in three scenarios:

- ▶ Where users have a licensing model that supports a certain number of active cores in a system.
- ▶ Where users have poorly threaded applications but require the additional LLC (last level cache) available to additional processors, but not the core overhead.
- ▶ Where users are looking to limit the number of active cores in an attempt to reclaim power and thermal overhead to increase the probability of Turbo Mode being engaged.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Processors** → **Cores in CPU Package**
- ▶ OneCLI/ASU variable: `<value> Processors.CoresinCPUPackage`
- ▶ Redfish attribute: `Processors_CoresinCPUPackage`

Possible values:

- ▶ **All** (Default)
All cores are enabled.
- ▶ **1**
Only 1 core is enabled in each CPU package.
- ▶ **2**
Only 2 cores are enabled in each CPU package.
- ▶ **3**
Only 3 cores are enabled in each CPU package.
- ▶ **N-1**
N is the maximum cores available in CPU package.

CPU Frequency Limits

Turbo ratio limits can be set between the rated frequency and the default turbo frequencies programmed for the number of cores active. Frequency is set from (Max - 1 bin) to (Max - 4 bins) for number of active cores, in approximately 100MHz increments. A *bin* is an ~100 MHz increment of core frequency.

It is useful when one desires control over turbo mode frequency uplift to minimize potential frequency fluctuations and reduce jitter.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Processors?** → **CPU Frequency Limits**
- ▶ OneCLI/ASU variable: `Processors.CPUFrequencyLimits`
- ▶ Redfish attribute: `Processors_CPUFrequencyLimits`

Possible values:

- ▶ **Full turbo uplift** (Default)
Turbo could run to maximum frequency supported by CPU when Turbo Mode is enabled.
- ▶ **Restrict maximum frequency**
Turbo is only able to run to limited frequency number when Turbo Mode is enabled.

Processors X to Y Cores Active

This setup menu is available when the parameter CPU Frequency Limits is set to **Restrict maximum frequency** as described in the previous section.

The maximum frequency (turbo, AVX, and non turbo) can be restricted to a frequency that is between the maximum turbo frequency for the CPU installed and 1.2GHz. This can be useful for synchronizing CPU tasks.

Note: The max frequency for N+1 cores cannot be higher than N cores. If an unsupported frequency is entered, it will automatically be limited to a supported value. If the CPU frequency limits are being controlled through application software, leave this menu item at the default, full turbo uplift.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Processors?** → **CPU Frequency Limits**
- ▶ OneCLI/ASU variable: `Processors.ProcessorsXtoYcoresactive`
Example: `Processors.Processors17to18coresactive`.
- ▶ Redfish attribute: `Processors_ProcessorsXtoYcoresactive`

X and Y are used more as engineering numbers.

- ▶ X means the start core number of turbo limitation, has a different value SKU to SKU, and can be changed if some CPU cores are disabled.
- ▶ Y means the end core number of turbo limitation, equal to actual enabled CPU cores.

User doesn't need to care the meaning of X and Y, just need to pay attention on the specific variable name when using OneCLI or Redfish to configure this setting.

Possible values:

- ▶ **Maximum Turbo Frequency - 1 bin** (Default)
In setup menu, it shows the actual available highest frequency value that can be restricted. For example, 3.40GHz
- ▶ **Maximum Turbo Frequency - 2 bins**
- ▶ **Maximum Turbo Frequency - 3 bins**
- ▶ ...
- ▶ **Maximum Turbo Frequency - N bins**

In setup menu, it shows the actual available lowest frequency value that can be restricted. For example, 3.2GHz.

Processor Prefetchers

A processor prefetcher attempts to predictively retrieve information from main memory or a higher level processor cache for storage in a processor's L1 or L2 cache. Prefetchers allow a CPU to do more work per clock cycle by maintaining a CPU core fed with information.

When a prefetcher successfully pre-stages information in a CPU's cache, it can serve to lower latency and improve performance. When a prefetcher pre-stages invalid information, there is an extra latency and performance impact as the cache must be flushed of information and new information must be loaded in.

Prefetchers come in two types and two varieties. Prefetchers either pre-stage information to a core's L1 or L2 cache, and either retrieve instructions or data from main memory.

There are up to four tunable prefetchers in ThinkSystem server Intel platforms:

- ▶ **Data Cache Unit (DCU) IP Prefetcher** – a processor L1 instruction cache prefetcher that uses a predictive instruction prefetching scheme. This setting fetches the next cache line into the L1 data cache if there is a sequential load history of cache line accesses. The next cache line is fetched from either L2 or main memory.
- ▶ **Data Cache Unit (DCU) Streamer Prefetcher** – a processor L1 data cache prefetcher that fetches the next cache line into the L1 data cache when multiple loads from the same cache line are executed in a certain time limit. The next cache line is fetched from the L2 cache or main memory.
- ▶ **Hardware Prefetcher (MLC Streamer Prefetcher)** – a processor L2 cache instruction prefetcher that prefetches extra cache lines for every memory request issued. The hardware prefetcher retrieves additional cache lines into the L2 cache of a core based on current requests.
- ▶ **Adjacent Cache Line Prefetcher (MLC Spatial Prefetcher)** – a processor L2 data cache prefetcher that fetches both cache lines that make up a 128- byte cache line pair even if the requested data is only in the first cache line. When disabled, this setting only fetches the required data from a single 64-byte cache line.

Application information access patterns, which tend to be relatively predictable, benefit greatly from prefetching. Most typical Line of Business (LOB) applications, virtualization and scientific applications benefit from having pre-fetching enabled. However, due to the non-inclusive LLC structure of Intel's Intel Xeon Scalable processors, it is recommended to carefully evaluate application performance to determine the benefit of disabling prefetching. Keep these settings enabled unless low-level cache analysis has been done with Intel tools or experiments have been run selectively disabling the prefetchers one at a time.

The Processor Prefetchers setting is accessed as follows:

- ▶ System setup:
 - **System Settings** → **Processors** → **Hardware Prefetcher**
 - **System Settings** → **Processors** → **Adjacent Cache Prefetch**
 - **System Settings** → **Processors** → **DCUStreamer Prefetcher**
 - **System Settings** → **Processors** → **DCU IP Prefetcher**
- ▶ OneCLI/ASU variable:
 - Processors.HardwarePrefetcher
 - Processors.AdjacentCachePrefetch
 - Processors.DCUStreamerPrefetcher
 - Processors.DCUIPPrefetcher

- ▶ Redfish attribute:
 - Processors_HardwarePrefetcher
 - Processors_AdjacentCachePrefetch
 - Processors_DCUStreamerPrefetcher
 - Processors_DCUIPPrefetcher

Possible values:

- ▶ **Enabled** (Default)
- ▶ **Disabled**

XPT Prefetcher

Extended Prediction Table (XPT) prefetcher (memory prefetch from the core) is a mechanism that enables a read request that is being sent to the last level cache to speculatively issue a copy of that read to the memory controller prefetching. It is designed to reduce local memory access latency.

XPT prefetcher is an “LLC miss predictor” in each core that will issue a speculative DRAM read request in parallel to an LLC lookup, but only when XPT predicts a “miss” from the LLC lookup. Data from an XPT prefetcher will wait for a short time in the memory and will be returned to the core only if there is an actual LLC miss. Local memory access latency is reduced because the LLC lookup has been bypassed by the speculative DRAM read.

Note: If **XPT Prefetcher** is set to **Enabled** and an unbalanced memory population is installed and **SNC** is set to **Enabled** (see “SNC (Processor Sub-NUMA Clustering)” on page 25), the prefetcher will remain disabled. Balanced memory means the same capacity/speed/number of DIMMs are installed in the memory channels among all of the memory controllers in the CPU sockets.

XPT prefetcher is best used for those workloads which have good NUMA locality. This setting is recommended for low jitter / low latency workloads.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Processors** → **XPT Prefetcher**
- ▶ OneCLI/ASU variable: Processors.XPTPrefetcher
- ▶ Redfish attribute: Processors_XPTPrefetcher

Possible values:

- ▶ **Enabled** (Default)
- ▶ **Disabled**

UPI Prefetcher

Ultra Path Interconnect (UPI) prefetch enables an early memory read on the memory bus. The UPI receive path spawns a memory read to the memory controller prefetcher.

Note: If **UPI Prefetcher** is set to Enabled and an unbalanced memory population is installed AND SNC is Disabled (see “SNC (Processor Sub-NUMA Clustering)” on page 25), the prefetcher will remain disabled. Balanced memory means all populated memory channels should have the same total memory capacity and the same total number of ranks. In addition, all memory controllers on a processor socket should have the same configuration of memory DIMMs.

UPI prefetcher is best used for those workloads which have good NUMA locality. This setting is recommended for low jitter / low latency workloads. For high frequency workloads it is recommended to leave enabled unless a workload analysis has been performed.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Processors** → **UPI Prefetcher**
- ▶ OneCLI/ASU variable: Processors.UIPrefetcher
- ▶ Redfish attribute: Processors_UIPrefetcher

Possible values:

- ▶ **Enabled** (Default)
- ▶ **Disabled**

DCA (Direct Cache Access)

DCA capable I/O devices such as network controllers can place data directly into the CPU cache, which improves response times.

If a network controller supports DCA, this might provide lower packet transfer latency. However, if applications are managing the last-level cache (LLC) directly, this functionality may pollute the LLC state. It is recommended to enable this setting for High Frequency workloads and to evaluate its functionality for low latency and low jitter workloads.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Processors** → **DCA**
- ▶ OneCLI/ASU variable: Processors.DCA
- ▶ Redfish attribute: Processors_DCA

Possible values:

- ▶ **Enabled** (Default)
- ▶ **Disabled**

Uncore Frequency Scaling

When Uncore Frequency Scaling is enabled, the CPU uncore will dynamically change speed based on the workload. This involves a separate voltage and frequency scaling for Core /LLC complex. For the Intel Xeon Scalable generation of processors, enabling this feature turns on uncore power savings features.

After enabling this setting, you can use the Uncore Frequency Limit setting to set the maximum frequency when scaling. See “Uncore Frequency Limit” on page 37.

Lenovo recommends you disable this setting for maximum performance for servers with Intel Xeon Scalable processors.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Processors** → **Uncore Frequency Scaling**
- ▶ OneCLI/ASU variable: `Processors.UncoreFrequencyScaling`
- ▶ Redfish attribute: `Processors_UncoreFrequencyScaling`

Possible values:

- ▶ **Enabled** (Default)
- ▶ **Disabled**

SNC (Processor Sub-NUMA Clustering)

SNC (Processor Sub-NUMA Clustering) partitions Intel Xeon Scalable processor cores and last-level cache (LLC) into disjoint clusters with each cluster bound to a set of memory controllers in the system. SNC improves average latency to the LLC and memory. SNC is a replacement for the cluster on die (COD) feature found in previous processor families such as the Intel Xeon E5 v4 and E7 v4 processors. For a multi-socketed system, all SNC clusters are mapped to unique NUMA domains.

When workloads are highly NUMA affinitized, SNC can enable better memory performance and utilization of the LLC. Recommend testing an application workload prior to enabling.

Note: OS support that recognizes each cluster and a separate NUMA node is necessary to take advantage of SNC. Consult your OS documentation to determine if SNC is supported.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Processors** → **SNC**
- ▶ OneCLI/ASU variable: `Processors.SNC`
- ▶ Redfish attribute: `Processors_SNC`

Possible values for ThinkSystem V1 servers:

- ▶ **Enabled**

Create 2-way Sub-NUMA Clusters for each Xeon Scalable Family and next generation of Xeon Family CPU. Utilizes LLC capacity more efficiently and reduces latency due to core/Integrated Memory Controller proximity. This may provide performance improvement on NUMA-aware operating systems.

- ▶ **Disabled** (Default)

The LLC is treated as one cluster when this option is disabled

Possible values for ThinkSystem V2 servers:

- ▶ **Disabled** (Default)

The LLC is treated as one cluster when this option is disabled

► **SNC**

SNC (2-way Sub NUMA) partitions the last level cache (LLC) into two clusters called NUMA nodes that contain an equal number of cores, equal number of LLC slices in close proximity to the cores, an equal amount of socket address space, and with each cluster bound to a subset of the memory controllers in the socket.

This requires full DIMM symmetry (left-to-right and top-to-bottom). If there is not full DIMM symmetry, the SNC will be forced to Disabled.

When SNC is selected, the OS sees two NUMA nodes per physical CPU package, except for 38-core and less than 12-core 3rd Gen Intel Xeon Scalable Processors where the SNC is not supported.

► **Virtual NUMA**

When Virtual NUMA is enabled, UEFI creates 2 NUMA nodes per physical CPU package without changing memory controller/channel interleaving and LLC grouping. There is a potential memory bandwidth advantage in Virtual NUMA mode. The latency between these two virtual NUMA nodes is identical to its local latency

Snoop Preference (Processor Snoop Mode)

Snoop mode determines the behavior in which a processor will validate an instruction or piece of data before performing an operation on it. The processor snoop function keeps caches coherent across the Intel processor interconnect fabric thus guaranteeing that data reads from cache or memory obtain the current copy of the data.

Snoop modes should be configured based on the specific instruction and memory access patterns of a particular workload. In general, low latency workloads tend to benefit from the Home Snoop w. Directory + OSB + HitME cache snoop mode, known as Home Snoop Plus, though the NUMA optimization of one's application should be evaluated.

Note: Setting the snoop mode preference does not always guarantee that it will be selected. The mode will be changed if the current hardware configuration does not support the desired mode

This setting is accessed as follows:

- System setup: **System Settings** → **Processors** → **Snoop Preference**
- OneCLI/ASU variable: `Processors.SnoopPreference`
- Redfish attribute: `Processors_SnoopPreference`

Possible values:

► **Home Snoop Plus** (Default)

Home Snoop Plus = Home Snoop with Directory lookup + Opportunistic Snoop Broadcast (OSB) + HitME cache

Best overall for most workloads. Speculative home snoop broadcasts are done under light UPI load conditions to avoid directory lookups from memory and reduce memory bandwidth. OSB is used only for local InvltoE (invalidate state to exclusive state) requests, which are generated due to full-line writes from the core or IO. Local InvltoE requests don't require a data read for this operation; hence the opportunistic broadcast feature is only used with writes to local memory to avoid memory access due to directory lookup. HitME cache is in the Caching and Home Agent (CHA) and caches directory information to speed up cache-to-cache transfers. HitME cache resources scale with number of CHAs in the system.

- ▶ **Home Snoop**

Best for bandwidth sensitive workloads. Home snoop may increase performance on those workloads requiring maximum local and remote bandwidth for cache / memory lookups.

Socket Interleave (Memory Socket Interleave)

Socket interleave determines how the memory map will be laid out within the system. Memory is either laid out such that each CPU has a map of local attached memory (NUMA) or in a flat memory model with no NUMA nodes (Non-NUMA).

Under most circumstances, customers should leave this setting at the default. Only in certain legacy applications which do not handle NUMA memory should this be changed.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Memory** → **Socket Interleave**
- ▶ OneCLI/ASU variable: `Memory.SocketInterleave`
- ▶ Redfish attribute: `Memory_SocketInterleave`

Possible values:

- ▶ **NUMA** (Default)

Each CPU socket is interleaved into separate NUMA nodes, improves average latency to LLC and memory.

- ▶ **Non-NUMA**

Memory mapping is a flat model. BIOS will attempt to interleave all sockets into one NUMA node. If an OS or system application is not NUMA-aware, user may need to choose this option.

Patrol Scrub

Memory scrubbing technology runs at the hardware level to prevent the accumulation of correctable errors within a block of memory that could result in uncorrectable errors. Patrol scrub works in the background to proactively check DIMMs for memory errors. Enabling patrol scrub will negatively impact performance per watt.

Memory scrubbing is strongly encouraged to be implemented on production systems, as it will aid in maintaining the resiliency of the memory subsystem. Patrol scrub will provide the most consistent verification, though at the cost of a slight loss of memory bandwidth and increase in latency. In those workloads where memory latency is vital to application performance customers may elect to disable all memory scrubbing features to gain additional performance.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Memory** → **Patrol Scrub**
- ▶ OneCLI/ASU variable: `Memory.PatrolScrub`
- ▶ Redfish attribute: `Memory_PatrolScrub`

Possible values:

- ▶ **Enabled** (Default)

Memory patrol scrubbing is enabled with 24 hours interval. If the correctable error count reaches the error threshold, then the issuing memory page would be retired with an error record in XCC.

- ▶ **Disabled**

Memory patrol scrubbing is disabled.

Memory Data Scrambling

Memory traffic on the data bus is not random and can cause current “hot spots” on the DIMM. Memory data scrambling uses a data scrambling feature in the memory controller to create pseudo-random patterns on the DDR4 data bus to reduce possibility of data-bit errors due to the impact of excessive current fluctuations.

Lenovo recommends leaving this parameter enabled to avoid data-bit errors. Recommended for all workloads.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Memory** → **Memory Data Scrambling**

- ▶ OneCLI/ASU variable: `Memory.MemoryDataScrambling`

- ▶ Redfish attribute: `Memory_MemoryDataScrambling`

Possible values:

- ▶ **Enabled** (Default)

- ▶ **Disabled**

Intel Virtualization Technology

Intel Virtualization Technology (VTx) is a set of technologies within a system that accelerate virtualization instructions by carrying them out directly in the silicon rather than as code to be binarily translated from a Guest OS, through a Hypervisor to a physical server platform.

VTx must be enabled when a system will be used for virtualization. For general business applications, VTx has little to no performance impact. In those scenarios where high frequency low latency code execution is paramount, VTx should be disabled. Financial, scientific and HPC applications often benefit from disabling VTx.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Processors** → **Intel Virtualization Technology**

- ▶ OneCLI/ASU variable: `Processors.IntelVirtualizationTechnology`

- ▶ Redfish attribute: `Processors_IntelVirtualizationTechnology`

Possible values:

- ▶ **Enabled** (Default)

- ▶ **Disabled**

Intel Virtualization Technology for I/O

Intel Virtualization Technology for I/O (VT-d) includes four key capabilities:

- ▶ I/O device assignment. This feature allows an administrator to assign I/O devices to VMs in any desired configuration.
- ▶ DMA remapping. Supports address translations for device DMA data transfers.
- ▶ Interrupt remapping. Provides VM routing and isolation of device interrupts.
- ▶ Reliability features. Reports and records system software DMA and interrupt errors that may otherwise corrupt memory and impact VM isolation.

VT-d is used in Virtualization environments. For low latency and low jitter environments, we recommend disabling VT-d.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Devices and I/O Ports** → **Intel VT for Directed I/O (VT-d)**
- ▶ OneCLI/ASU variable: `DevicesandIOPorts.IntelVTforDirectedIOVTd`
- ▶ Redfish attribute: `DevicesandIOPorts_IntelVTforDirectedIOVTd`

Possible values:

- ▶ **Enabled** (Default)
- ▶ **Disabled**

PCIe Power Brake

PCIe Power Brake quickly reduces the power consumption and performance of high powered PCIe devices (e.g.: GPU, FPGA). Performance of PCIe devices that are low power are not impacted by this setting. A high powered PCIe device is one that is rated at 75W TDP or greater.

Lenovo recommends enabling Proactive for all workloads.

3rd Gen Intel Xeon Scalable processors: This setup option is available starting with ThinkSystem V2 servers with 3rd Gen Intel Xeon Scalable processors. This setup option is not available with Intel ThinkSystem server V1 products.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Power** → **PCIe Power Brake**
- ▶ OneCLI/ASU variable: `Power.PCIePowerBrake`
- ▶ Redfish attribute: `Power_PCIePowerBrake`

Possible values:

- ▶ **Reactive**

The system performs PCIe throttling when a PSU indicates that a momentary overconsumption or high temperature warning is occurring. The system also proactively performs PCIe power throttling to thermally protect high powered PCIe devices in unsupported high ambient temperatures and during loss of cooling redundancy.

The high temperature warning setting is different among different platforms (typical setting is 35°C degree), and it would be adjusted to a lower value to avoid the loss of cooling redundancy.

► **Proactive** (Default)

The system performs PCIe throttling based on the maximum power rating of the installed high power PCIe adapters. High powered PCIe device performance is reduced when the total power from high powered PCIe devices is greater than 1/3 of the sum of active power supply (PSU) wattages. Proactive mode also includes the PCIe throttling features of Reactive mode.

To de-assert PCIe throttling for this case, user needs to remove some GPUs to down the total TDP of GPU, or replace a proper higher watt PSU.

► **Disabled**

The system won't perform PCIe throttling. Proactive thermal protection provided under Proactive/Reactive modes are not operational. PCIe throttling is limited to the self-throttling capabilities of the high powered PCIe devices.

ASPM (PCIe Active State Power Management)

ASPM is a PCIe power saving feature. It puts the PCIe link into a low power mode when the link is idle. Two low power “standby” Link states are defined for ASPM – L0 low power Link state and L1 Link state:

- The L0s low power Link state is optimized for short entry and exit latencies, while providing substantial power savings. If there is a retimer on the link path, L0s should be disabled.
- The L1 Link state is optimized for maximum power savings at a cost of longer entry and exit latencies. L1 reduces Link power beyond the L0s state for cases where very low power is required and longer transition times are acceptable.

There are two pitfalls to bear in mind when configuring ASPM:

- Some older PCIe devices didn't support L0s for both transmit and receive cycles. There is a newer PCIe engineering change notice (ECN) that now requires software to only enable L0s, if both ends of the link support L0s. L1 already required both ends to be in L1 and software has to check it. The ECN now also allows hardware to support L1 state with support for L0s.

Members of the PCI SIG can view more details on the change at:

http://www.pcisig.com/specifications/pciexpress/specifications/ECN_ASPM_Optionality_2009-08-20.pdf

- Under Linux, the ASPM policy can be set either by writing to a sysfs file/object (/sys/modules/pcie_aspm/parameters/policy) or from the Linux kernel command line (boot with `pcie_aspm.policy=X`). In either case, the valid values for “policy” are `powersave`, `performance`, and `default`. When `default` is chosen, it means Linux will use the policy specified by the platform firmware.

The `pcie_aspm` option simply provides runtime control for disabling/enabling ASPM. There are 2 valid options for `pcie_aspm`: `off` and `force`. The default setting for ASPM in the Linux kernel is enabled. The 2 values for `pcie_aspm` are:

- `off` – disables the ASPM in OS
- `force` – acts an override to platform firmware settings. In other words, this forces ASPM to be enabled, regardless of how the platform firmware settings are configured.

Lenovo recommends disabling ASPM for all systems where workload performance is more important than power savings. Meanwhile, after enabling ASPM, some PCIe devices could encounter uncorrectable errors.

This setup option is available starting with Intel platform based ThinkSystem V2 products.

This setting is accessed as follows:

- ▶ System setup: **System Settings** → **Power** → **ASPM**
- ▶ OneCLI/ASU variable: Power.ASPM
- ▶ Redfish attribute: Power_ASPM

Possible values:

- ▶ Auto
Enable ASPM on PCIe endpoint adapters that support it.
- ▶ Disabled (Default)
Disable ASPM for all PCIe endpoints. User need to disable ASPM is problems occur.

Hidden UEFI Items

The following UEFI items are more limited in their applicability to customer use cases and are not exposed in UEFI menus but can be accessed using the command line utilities such as Lenovo's Advanced Settings Utility (ASU) or OneCLI.

No Redfish support: These UEFI items also have no Redfish attribute and can't be accessed via the Redfish interface.

The UEFI menu items described here are as follows:

- ▶ "CR QoS Configuration and CR FastGo Configuration"
- ▶ "L2 RFO (Request of Ownership) Prefetch Disable" on page 32
- ▶ "Local/Remote Threshold" on page 32
- ▶ "Stale AtoS" on page 34
- ▶ "Snoop Response Hold Off" on page 35
- ▶ "LLC dead line allocation" on page 36
- ▶ "LLC (Last Level Cache) Prefetch" on page 36
- ▶ "Uncore Frequency Limit" on page 37

CR QoS Configuration and CR FastGo Configuration

One of 4 variables (IRQThreshold, StaleAtoS, CRQoSConfiguration/CRFastGoConfiguration, L2RFOPrefetchDisable) used to optimize performance for SAP HANA on servers with 2-hop memory configurations such as 4-socket ring, 6-socket and 8-socket configurations.

Beginning with 2nd Gen Intel Xeon Scalable Processors (Cascade Lake), the default value of **Auto** corresponds to a setting of **Option 5** when a 2-hop memory configuration is detected. Option 5 makes the L2 prefetcher less aggressive and lowers NT (Non-Temporal) write bandwidth. This has the effect of limiting bursty local and remote traffic.

This setting is accessed as follows:

- ▶ OneCLI/ASU variable: `Memory.CRQoSConfiguration` or `Memory.CRFastGoConfiguration` with current UEFI builds. Some older UEFI builds may still use the variable name `Memory.CRQoSConfiguration`.

Possible values:

- ▶ **Auto** (Default)
The **Auto** value corresponds to **Option 5** for 2-hop memory configurations and **None** for all other configurations.
- ▶ **None**
No optimization
- ▶ **Option 1**
- ▶ **Option 2**
- ▶ **Option 3**
- ▶ **Option 4**
- ▶ **Option 5**
Higher option number makes L2 prefetcher less aggressive and lowers NT write bandwidth.

L2 RFO (Request of Ownership) Prefetch Disable

One of 4 variables (`IRQThreshold`, `StaleAtoS`, `CRQoSConfiguration/CRFastGoConfiguration`, `L2RFOPrefetchDisable`) used to optimize performance for SAP HANA on servers with 2-hop memory configurations such as 4-socket ring, 6-socket and 8-socket configurations.

The Enabled option makes L2 prefetcher less aggressive and lowers NT (Non-Temporal) write bandwidth. Disable limits burstiness and reducing snooping.

This setting is accessed as follows:

- ▶ OneCLI/ASU variable: `Processors.L2RFOPrefetchDisable`

Possible values:

- ▶ **Disabled** (Default)
Enable L2 RFO Prefetch.
- ▶ **Enabled**
Disable L2 RFO Prefetch.

Local/Remote Threshold

In multi-socket environments that generate remote traffic, the Table Of Requests (TOR) services both remote and local requests across the Intel UPI links. Remote reads allocate the TOR through the RRQ (Remote Request Queue) ingress and typically send out a read request to the memory controller. When the data returns from the memory controller, the TOR will send a data response back to the requesting socket via the Intel UPI links.

The rate at which these requests and responses can be sent out is governed by Intel UPI bandwidth via credits. The rate is adjustable to limit the number of incoming transactions to

avoid overloading a remote socket. This bandwidth is substantially less than local memory bandwidth or TOR service bandwidth.

When the Intel UPI link is the bottleneck, these remote reads stay in the TOR for a long time. When one socket's TOR is saturated with remote reads, it cannot service its own local requests well enough. The other socket is relatively free to continue sending remote requests to the saturated socket, which may cause an imbalance. If the imbalance shifts around then the performance averages out over time and can converge to a balance, but if a socket is saturated, and continues to get requests from the better performing socket, it may not build enough momentum to overturn the imbalance in its favor. Application workloads with UMA and remote traffic mixes are likely candidates for this imbalance and could affect overall performance.

The UEFI item called *local/remote threshold* is available to control this rate and help avoid the imbalance. This option controls two features simultaneously:

- ▶ RRQ (Remote Request Queue) throttling limits remote reads in the TOR by specifying a threshold.
- ▶ IRQ (Incoming Request Queue) throttling limits the number of local-to-remote reads by specifying a threshold (the requests are from IRQ, but the home node of the requested address is a remote socket).

These features are used to control bandwidth of a multi-socket system, to maximize the throughput of each single socket while trying to make the performance across sockets as balanced as possible. The default threshold values set by BIOS are defined based on internal experiments and are believed to be optimal for typical workloads, but other options are provided to customers if their workload is showing noticeably different behavior from the internal workloads.

The threshold values set by the UEFI are dependent on the system configuration. The UEFI will first detect the platform configuration in terms of the number of sockets and the number of Intel UPI links and then automatically set the default mode accordingly.

To simplify the UEFI setting, the two features RRQ throttling and IRQ throttling are a unified option and the value of each option, based on the number of processors installed and the value of `IrqThreshold` selected, is listed in Table 3.

Table 3 Setting RRQ throttling and IRQ throttling using the `IrqThreshold` parameter

Value of <code>IrqThreshold</code>	1 socket (1S)	2S/4S-Mesh	4S Ring	8S
Auto (default)	Disabled mode	Medium mode	Medium mode	Medium mode
Disabled	RRQ=Disabled, IRQ=Disabled	RRQ=Disabled, IRQ=Disabled	RRQ=Disabled, IRQ=Disabled	RRQ=Disabled, IRQ=Disabled
Low Threshold	N/A	RRQ=6, IRQ=Disabled	RRQ=7, IRQ=2	RRQ=7, IRQ=2
Medium Threshold	N/A	RRQ=7, IRQ=Disabled	RRQ=7, IRQ=7	RRQ=7, IRQ=7
High Threshold	N/A	RRQ=7, IRQ=Disabled	RRQ=9, IRQ=10	RRQ=8, IRQ=10

This setting is accessed as follows:

- ▶ OneCLI/ASU variable: `Processors.IrqThreshold`

Possible values (see Table 3):

- ▶ **Auto** (Default)
- ▶ **Disabled**
- ▶ **Low**
- ▶ **Medium**
- ▶ **High**

Stale AtoS

Stale AtoS is one of four variables (IRQThreshold, StaleAtoS, CRQoSConfiguration/CRFastGoConfiguration, L2RFOPrefetchDisable) used to optimize performance for SAP HANA on servers with 2-hop memory configurations such as 4-socket ring, 6-socket and 8-socket configurations.

The in-memory directory within the Intel memory controller that tracks the usage state of all cache lines has three states: I, A and S:

- ▶ I (Invalid) state means the data is clean and does not exist in any other socket's cache.
- ▶ A (snoop All) state means the data may exist in another socket in exclusive or modified state.
- ▶ S (Shared) state means the data is clean and may be shared across one or more socket's caches.

When doing a read to memory, if the directory line is in the A state we must snoop all the other sockets because another socket may have the line in the modified state. If this is the case, the snoop will return the modified data. However, it may be the case that a line is read in A state and all the snoop come back a miss. This can happen if another socket read the line earlier and then silently dropped it from its cache without modifying it.

If Stale AtoS feature is enabled, in the situation where a line is in A state returns only snoop misses, the line will transition to S state. That way, subsequent reads to the line will encounter it in S state and not have to snoop, saving latency and snoop bandwidth.

Stale AtoS may be beneficial in a workload where there are many cross-socket reads.

This setting is accessed as follows:

- ▶ OneCLI/ASU variable: Processors.StaleAtoS

Possible values:

- ▶ **Auto** (Default), which corresponds to:
 - **Disable** when 2 CPUs are installed or if 1 CPU is installed and Intel persistent memory (Pmem) is not installed.
 - **Enable** if Intel Pmem is installed or if more than 2 CPUs are installed.
- ▶ **Disabled**
- ▶ **Enabled**

Snoop Response Hold Off

Set the I/O delay in response to snoop request. This parameter is used to fix a low InfiniBand performance problem. This hidden parameter is accessible on ThinkSystem servers that use Intel Xeon Scalable processors.

For some workloads in which throughput and latency are critical, it is better to constantly poll the status of an I/O device rather than use an interrupt. Network adapter device drivers commonly use a thread to continuously poll in a fast loop so that incoming requests can be handled as fast as possible.

Continuous polling can create contention between a processor core running the polling thread and the processor's *Integrated I/O* feature (IIO) for an I/O-owned line in cache. This contention can cause an I/O operation to lose ownership of the cache line it has just acquired. It must then spend more time reacquiring the cache line to write it back.

When there are a large number of network ports each servicing small packets, the system may not be able to achieve the full throughput required due to excessive I/O and core contentions of cache lines. For this situation, the I/O operation should delay its response to core snoops and hold onto its cache lines until it successfully completes its write.

The Snoop Response Hold Off parameter allows the I/O operation to delay its snoop response by a selected amount to achieve this delay.

This setting is accessed as follows:

- ▶ OneCLI/ASU variable: `Processors.SnoopResponseHoldOff`

Possible values:

- ▶ 0, 1, 2 ... F

Set the Hold Off IO delay value. The value entered is a hex digit and the valid range is 0 to F. The default (initial value) is 9 for ThinkSystem servers. The value corresponds to a number cycles as shown in Table 4. The cycles in the table are IIO clock cycles which are 2 ns per cycle. For Intel Omni-Path network adapter use, setting this parameter to 9 is recommended as a starting point. Network performance tests should be performed to determine the most optimal value for each workload.

Table 4 Possible values of the parameter in the `Processors.SnoopResponseHoldOff` command

Value	Number of cycles (1 cycle = 2 nanoseconds)
0	Disabled
1	8
2	16
3	32
4	64
5	128
6	256
7	512
8	1K
9	2K

Value	Number of cycles (1 cycle = 2 nanoseconds)
A	4K
B	8K
C	16K
D	32K
E	64K
F	128K

LLC dead line allocation

When enabled, opportunistically fill dead lines in LLC. In the Intel Xeon Scalable family non-inclusive cache scheme, MLC evictions are filled into the LLC. When lines are evicted from the MLC, the core can flag them as “dead” (in other words, not likely to be read again). The LLC has the option to drop dead lines and not fill them in the LLC. If the Dead Line LLC Allocation feature is disabled, dead lines will always be dropped and will never fill into the LLC.

Not filling the LLC with dead lines can help save space in the LLC and prevent the LLC from evicting useful data. However, if the Dead Line LLC Allocation feature is enabled, the LLC can opportunistically fill dead lines into LLC if there is free space available.

This setting is accessed as follows:

- ▶ OneCLI/ASU variable: `Processors.LLCdeadlinealloc`

Possible values:

- ▶ **Enabled** (Default)
- ▶ **Disabled**
- ▶ **Auto**

Map to Enabled.

LLC (Last Level Cache) Prefetch

The LLC prefetcher is a new prefetcher added to the Intel Xeon Scalable family as a result of the non-inclusive cache architecture. The LLC prefetcher is an additional prefetch mechanism on top of the existing prefetchers that prefetch data into the core DCU and the MLC.

Enabling LLC prefetch gives the core prefetcher the ability to prefetch data directly into the LLC without necessarily filling into the MLC. In some cases, setting this option to disabled can improve performance.

This setting is accessed as follows:

- ▶ OneCLI/ASU variable: `Processors.LLCPrefetch`

Possible values:

- ▶ **Disabled** (Default)
- ▶ **Enabled**

Uncore Frequency Limit

When uncore frequency scaling is enabled (see “Uncore Frequency Scaling” on page 24), the Uncore Frequency Limit setting is used to set the upper frequency limit of the uncore in the CPU.

Note that this setting is not available if Uncore Frequency Scaling is disabled. When Uncore Frequency Scaling is disabled, the uncore runs at the maximum frequency all the time.

This setting is accessed as follows:

- ▶ OneCLI/ASU variable: `Processors.UncoreFrequencyLimit`

Possible values:

- ▶ Hex value 0, 1, 2 ... 3F

Specifies the CPU uncore frequency limit value. The value is entered as a hex value in the range of 0 to 3F. The default value is 3F (maximum). The value, when converted to decimal, corresponds to 1/10th of a GHz. For example, 3F is 63 in decimal and corresponds to a limit of 6.3 GHz, and a value of 14 hex (20 in decimal) corresponds to a limit of 2.0 GHz.

Don't include any prefix or suffix when entering the hex value.

The default value of 3F (63 in decimal, resulting in a limit of 6.3 GHz) is the maximum possible value that can be set. However, if a value is set that is higher than what the installed CPU SKU supports, the CPU will clip the value to the highest supported value. For example, if uncore frequency limit is set to 3F (6.3 GHz) but the CPU SKU has a maximum uncore frequency of 2.0 GHz, then the limit will be clipped to 2.0 GHz.

Red Hat Enterprise Linux (RHEL) and derivatives

The following are general considerations for optimizing systems for low jitter / low latency under Red Hat Enterprise Linux / CentOS Linux

- ▶ Disable unnecessary system services
- ▶ Avoid non-data disk access by:
 - Consider disabling filesystem journaling
 - Consider running w/o swap partition / swap file
- ▶ For a quick set of OS level tunings for low latency / high frequency workloads consider running tuned-adm profile latency-performance
- ▶ Disable power management in both kernel and UEFI
 - Use **powertop** or **turbostat** to see if there are P-state and C-state transitions
 - Set both **intel_idle.max_cstate=0** and **processor.max_cstate=n** (n = 0 to 3)
CAUTION: `intel_idle.max_cstate=0` doesn't mean maximum state=0(zero). This option disables `intel_idle`.
 - If you set **intel_idle.max_cstate=n** (n= 1 to 6), then `intel_idle` is enabled and maximum cstate is set to `n`.
- ▶ Consider using the `deadline` elevator for I/O tuning
ELEVATOR = "deadline"

- ▶ Enforce NUMA locality for applications / threads:
`numactl -C1 -m1 ./command`
`numad` user-level daemon (new since RHEL 6.3 and disabled by default) attempts automatically co-locate a process and its memory
- ▶ Consider using control groups (Cgroups) for CPU/Memory/Network/Disk to manage multi-application and / or large NUMA machines
- ▶ Make sure adapters are using MSI-X for IRQ deployment
 For example, use the following and then look for PCI-MSI-X in the output
`cat /proc/interrupts | grep eth`
- ▶ Consider pinning adapter IRQs on large NUMA systems (requires deep knowledge of systems and applications – not for the faint of heart). This requires disabling `irqbalance` daemon
- ▶ Disable `usb0` to avoid USB device polling and unplug and USB devices (keyboards, mice, USB keys, etc)

Additional References

The following resources provide additional information:

- ▶ Low Latency Performance Tuning for Red Hat Enterprise Linux 7
<https://access.redhat.com/sites/default/files/attachments/201501-perf-brief-low-latency-tuning-rhel7-v1.pdf>
- ▶ 3rd Generation Intel Xeon Scalable Processors
<https://ark.intel.com/content/www/us/en/ark/products/series/204098/3rd-generation-intel-xeon-scalable-processors.html>
- ▶ Lenovo Press paper, *Intel Xeon Scalable Family Balanced Memory Configurations*
<https://lenovopress.com/lp0742>
- ▶ Lenovo Press paper, *Second Generation Intel Xeon Scalable Family Balanced Memory Configurations*
<https://lenovopress.com/lp1089>
- ▶ Lenovo ThinkSystem Memory Configurator
https://dcsc.lenovo.com/#/memory_configuration
- ▶ WikiChip Chip & Semi - User can get CPU detailed information by search on any CPU SKU
<https://en.wikichip.org/wiki/intel>

Authors

This paper was produced by the following team of specialists:

John Encizo is a Senior Technical Sales Engineer on the Brand Technical Sales team for Lenovo across the Americas region. He is part of the field engineering team responsible for pre-GA support activities on the ThinkSystem server portfolio and works closely with customers and Lenovo development on system design, BIOS tuning, and systems

management toolkits. He has also been the lead field engineer for high-end server product launches. When not in the field doing technical design and onboarding sessions with customers or pre-GA support activities, he works in the BTS lab to do remote customer proof-of-technology engagements, consult on systems performance and develop documentation for the field technical sellers.

Kin Huang is a Principle Engineer for firmware in Lenovo Infrastructure Solutions Group in Beijing. He has 2 years of experience in software application programming, and 15 years of experience in x86 server, focusing on BIOS/UEFI, management firmware and software. His current role includes firmware architecture for AMD platforms, boot performance optimization, RAS (Reliability, Availability and Serviceability). Kin holds Bachelor of Applied Chemistry from Jilin University, and he is working on a Masters of Engineering Management of Computer Science and Technology from Wuhan Institute of Technology.

Joe Jakubowski is the Principal Engineer and Performance Architect in the Lenovo Infrastructure Solutions Group Performance Laboratory in Morrisville, NC. Previously, he spent 30 years at IBM. He started his career in the IBM Networking Hardware Division test organization and worked on various token-ring adapters, switches and test tool development projects. For more than 25 years, he has worked in the x86 server performance organization focusing primarily on database, virtualization and new technology performance. His current role includes all aspects of Intel and AMD x86 server architecture and performance. Joe holds Bachelor of Science degrees with Distinction in Electrical Engineering and Engineering Operations from North Carolina State University and a Master of Science degree in Telecommunications from Pace University.

Robert R. Wolford is the Principal Engineer for power management and efficiency at Lenovo. He covers all technical aspects that are associated with power metering, management, and efficiency. Previously, he was a lead systems engineer for workstation products, video subsystems building block owner for System x, System p, and desktops, signal quality and timing analysis engineer, gate array design, and product engineering. In addition, he was a Technical Sales Specialist for System x. Bob has several issued patents centered around computing technology. He holds a Bachelor of Science degree in Electrical Engineering with Distinction from the Pennsylvania State University.

Thanks to the following people for their contributions to this project:

- ▶ David Watts
- ▶ Charles Stephan
- ▶ Larry Cook
- ▶ Sampson Chien
- ▶ Kevin Huang

Notices

Lenovo may not offer the products, services, or features discussed in this document in all countries. Consult your local Lenovo representative for information on the products and services currently available in your area. Any reference to a Lenovo product, program, or service is not intended to state or imply that only that Lenovo product, program, or service may be used. Any functionally equivalent product, program, or service that does not infringe any Lenovo intellectual property right may be used instead. However, it is the user's responsibility to evaluate and verify the operation of any other product, program, or service.

Lenovo may have patents or pending patent applications covering subject matter described in this document. The furnishing of this document does not give you any license to these patents. You can send license inquiries, in writing, to:

Lenovo (United States), Inc.
1009 Think Place - Building One
Morrisville, NC 27560
U.S.A.
Attention: Lenovo Director of Licensing

LENOVO PROVIDES THIS PUBLICATION "AS IS" WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESS OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. Some jurisdictions do not allow disclaimer of express or implied warranties in certain transactions, therefore, this statement may not apply to you.

This information could include technical inaccuracies or typographical errors. Changes are periodically made to the information herein; these changes will be incorporated in new editions of the publication. Lenovo may make improvements and/or changes in the product(s) and/or the program(s) described in this publication at any time without notice.

The products described in this document are not intended for use in implantation or other life support applications where malfunction may result in injury or death to persons. The information contained in this document does not affect or change Lenovo product specifications or warranties. Nothing in this document shall operate as an express or implied license or indemnity under the intellectual property rights of Lenovo or third parties. All information contained in this document was obtained in specific environments and is presented as an illustration. The result obtained in other operating environments may vary.

Lenovo may use or distribute any of the information you supply in any way it believes appropriate without incurring any obligation to you.

Any references in this publication to non-Lenovo Web sites are provided for convenience only and do not in any manner serve as an endorsement of those Web sites. The materials at those Web sites are not part of the materials for this Lenovo product, and use of those Web sites is at your own risk.

Any performance data contained herein was determined in a controlled environment. Therefore, the result obtained in other operating environments may vary significantly. Some measurements may have been made on development-level systems and there is no guarantee that these measurements will be the same on generally available systems. Furthermore, some measurements may have been estimated through extrapolation. Actual results may vary. Users of this document should verify the applicable data for their specific environment.

This document was created or updated on June 15, 2021.

Send us your comments via the **Rate & Provide Feedback** form found at <http://lenovopress.com/lp1477>

Trademarks

Lenovo and the Lenovo logo are trademarks or registered trademarks of Lenovo in the United States, other countries, or both. These and other Lenovo trademarked terms are marked on their first occurrence in this information with the appropriate symbol (® or ™), indicating US registered or common law trademarks owned by Lenovo at the time this information was published. Such trademarks may also be registered or common law trademarks in other countries. A current list of Lenovo trademarks is available from <https://www.lenovo.com/us/en/legal/copytrade/>.

The following terms are trademarks of Lenovo in the United States, other countries, or both:

Advanced Settings Utility™
Lenovo®

Lenovo(logo)®
System x®

ThinkSystem™

The following terms are trademarks of other companies:

Intel, Xeon, and the Intel logo are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States and other countries.

Linux is a trademark of Linus Torvalds in the United States, other countries, or both.

Windows, Windows Server, and the Windows logo are trademarks of Microsoft Corporation in the United States, other countries, or both.

Other company, product, or service names may be trademarks or service marks of others.