

Troubleshooting and Health Monitoring for Intel[®] Optane[™] DC SSDs

White Paper



Revision History

Revision	Description	Date
001	<ul style="list-style-type: none">Initial release	June 2020

Tests document performance of components on a particular test, in specific systems. Differences in hardware, software, or configuration will affect actual performance. Consult other sources of information to evaluate performance as you consider your purchase. For more complete information about performance and benchmark results, visit www.intel.com/benchmarks.

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Performance varies depending on system configuration. No product or component can be absolutely secure. Check with your system manufacturer or retailer or learn more at intel.com.

Test and System Configuration: Intel® Server Board S2600WF Family, OS: CentOS 7.5, kernel 5.2.11-1.el7.elrepo.x86_64, CPU 2 x Intel® Xeon® Gold 6244 CPU @ 3.60GHz (8 cores), RAM 192GB DDR@2666MHz, NVMe Driver: Inbox, C-states: Disabled, Hyper Threading: Disabled, CPU Governor (through OS): Performance Mode

The products described in this document may contain design defects or errors known as errata which may cause the product to deviate from published specifications. Current characterized errata are available on request.

© Intel Corporation. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others.



Contents

1	Introduction	4
1.1	Overview	4
1.2	Applicable Intel Solid State Drives	4
1.3	Compatible Operating Systems	4
2	Install and Use Open Sourced Tool	5
2.1	Install nvme-cli Tool to the System	5
2.2	Verify Installation of Intel® Optane™ DC SSDs	5
3	SMART Reading and Health Check	6
3.1	Health Indicators	6
3.2	SMART Log	6
3.3	Health Monitoring	7
3.4	Temperature Statistics	8
4	Troubleshooting Tips	9
4.1	Additional SMART Attribute Monitoring	9
4.2	Thermal Throttling Status	9
4.3	PCIe Link Speed and Link Width	9
4.4	Get Nlog	10
4.5	Retrieve Telemetry Log	10
4.6	Latency Tracking	10
5	Manageability	11
5.1	Firmware Update and Reset	11
5.2	Change Power Governor Mode	11
5.3	Format the SSD with Different Sector Size	12
5.4	Endurance Estimation	12



1 Introduction

1.1 Overview

The purpose of this document is to guide customers in the use of the open source nvme-cli tool with which they can perform health monitoring and troubleshooting on Intel® Optane™ SSD DC P4800X/P4801X. With the nvme-cli tool, customers can develop their own health management software to centrally manage their datacenter environment. This document is intended for IT administrators, engineers, solutions architects, and field sales professionals.

1.2 Applicable Intel Solid State Drives

- Intel® Optane™ SSD DC P4800X
- Intel® Optane™ SSD DC P4801X
- Intel® Optane™ SSD DC D4800X

1.3 Compatible Operating Systems

The example tests conducted in this document were performed on CentOS 7.5 with the latest version of the Linux Kernel available at the time of this document's initial publication. However, the following operating systems with default NVMe drivers should yield similar results:

- Red Hat Enterprise Linux 7.5 or later
- CentOS 7.5 or later
- SUSE Linux Enterprise Server 12 or later
- Windows 2019
- Ubuntu 18.04 or later

§



2 Install and Use Open Sourced Tool

2.1 Install nvme-cli Tool to the System

Download the open source tool, **nvme-cli**, from the following location: <https://github.com/linux-nvme/nvme-cli>

After clicking on the link, you will be given instructions on installing **nvme-cli** on your system.

The commands provided in this document can also be achieved using the Intel® SSD Data Center Tool (DCT). After downloading the Intel® SSD DCT [here](#), issue the following command to install:

```
# rpm -ivh isdct-3.0.24-1.x86_64.rpm
```

2.2 Verify Installation of Intel® Optane™ DC SSDs

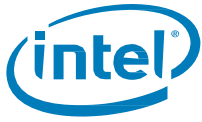
Use the following **nvme-cli** command to display the installed Intel® Optane™ DC SSDs.

Example:

```
# nvme list | grep SSDPE
/dev/nvme1n1      PHKE8156002R100EGN  INTEL SSDPE21K100GA      1
100.03 GB / 100.03 GB   512 B + 0 B   E2010475
```

This will list your SSD's serial number (SN), model number, capacity and current firmware version.

§



3 SMART Reading and Health Check

3.1 Health Indicators

Overall drive health status checking:

```
# nvme intel id-ctrl /dev/nvme0 |grep health
health : healthy
```

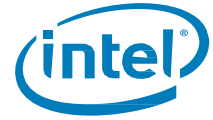
Or, you can use the Intel SSD DCT (isdct) to check:

```
# isdct sgow -intelssd 1 |grep DeviceStatus
DeviceStatus : Healthy
```

3.2 SMART Log

Issue the following command to display the list of SMART Attributes (Log Identifier 02h).

```
# nvme list
# nvme smart-log /dev/nvme1
Smart Log for NVME device:nvme1 namespace-id:ffffff
critical_warning           : 0
temperature                : 29 C
available_spare            : 100%
available_spare_threshold  : 0%
percentage_used            : 0%
data_units_read            : 761,322
data_units_written         : 1,304,568
host_read_commands        : 95,062,499
host_write_commands       : 176,593,139
controller_busy_time      : 19
power_cycles               : 42
power_on_hours             : 7,091
unsafe_shutdowns          : 20
media_errors               : 0
num_err_log_entries       : 0
Warning Temperature Time   : 0
Critical Composite Temperature Time : 0
Thermal Management T1 Trans Count : 0
Thermal Management T2 Trans Count : 0
Thermal Management T1 Total Time : 0
Thermal Management T2 Total Time : 0
```



3.3 Health Monitoring

- “Critical_Warning” Count

“Critical Warning” SMART attribute set by various warning sources:

- Available Spare is below Threshold
- Temperature has exceeded Threshold
- Reliability is degraded due to excessive media or internal errors
- Media is placed in Read-Only Mode
- Volatile Memory Backup System has failed (e.g., enhanced power loss capacity test failure)

Drive Health Indicator defined under bytes 3095-3076 of Identify Controller may still indicate “healthy” status even when the critical warning flag is set.

Note: From time to time it is important to use this “counter” to monitor for pre-directive failure. If Critical Warning “count” number increases above the pre-defined threshold, system administrator may need to take action (possible hardware replacement):

```
# nvme smart-log /dev/nvme1 | grep critical_warning
critical_warning           : 0
```

- “Available Spare” Indicator

This indicator is for P4800X and D4800X, not for P4801X. Available Space indicator will display a value of 100%, 75%, 50% or 25%.

Intel recommends the customer take no action if this value is greater than 0%.

```
# nvme smart-log /dev/nvme1n1 | grep available_spare
available_spare           : 100%
```

If this value indicates 0%, consider calling an Intel representative for hardware replacement.

- “Percentage Used” Indicator

This is Intel® Optane™ SSD endurance indicator. It will be zero when the SSD is first installed, and may remain at zero for a while. For more details on why this occurs, see this Technical Advisory:

<https://www.intel.com/content/www/us/en/support/articles/000033326/memory-and-storage/data-center-ssds.html>

As the SSD endurance is consumed, the indicated value will increase.

A value of 100 indicates that the estimated endurance of the device has been consumed, but may not indicate a device failure, as the value is allowed to exceed 100. Once the value reaches or exceeds 105, the drive will enter write protect mode, in which write bandwidth maxes out at <30MB/sec.

Any value from 95-100 should be interpreted as a warning that endurance consumption is approaching its maximum; once it reaches 105 performance may begin to suffer.

```
# nvme smart-log /dev/nvme1 | grep percentage_used
percentage_used           : 0%
```

- “Unsafe Shutdowns” Indicator

The importance of the Unsafe Shutdowns indicator will depend on your server system, HSBP stability, and hot plug times. Check this indicator number as needed.

```
# nvme smart-log /dev/nvme1n1 | grep unsafe_shutdowns
unsafe_shutdowns         : 31
```



- **“Media Errors”** Indicator

Indicates the number of unrecovered data integrity errors detected by the controller. Errors such as uncorrectable ECC, CRC checksum failure, or LBA tag mismatch are included in this field.

Combine “Healthy” Status and “Critical Warning” indicators to determine drive’s health. If the counter number has increased significantly, central management software should take action (based on pre-defined threshold).

```
# nvme smart-log /dev/nvme1n1 | grep media_errors
media_errors                : 0
```

- **“End-to-End Error Detection Count”** Indicator

Reports number End-to-End errors detected and corrected by the hardware:

```
# nvme intel smart-log-add /dev/nvme1n1 | grep detection
end_to_end_error_detection_count: 100%      0
```

Or you can use ISDCT tool:

```
# isdct show -sensor -intelssd 1 | grep Detection
EndToEndErrorDetectionCount : 0
```

- **“Temperature”** Indicator

Reports the SSD drive composite temperature which will help you understand the effectiveness of your server system’s thermal solution.

```
# isdct show -sensor -intelssd 1 | grep "Temperature - Celsius"
Temperature - Celsius : 30
```

Or you can use nvme-cli as well:

```
# nvme smart-log /dev/nvme1n1 | grep temperature
temperature                : 30 C
```

3.4 Temperature Statistics

The Temperature Statistics log is a vendor specific log page unique to Intel SSDs. To access the parameters listed in this log page please see the following example:

```
# nvme intel temp-stats /dev/nvme0

Intel Temperature Statistics
-----
Current temperature       : 41
Last critical overtemp flag : 0
Life critical overtemp flag : 0
Highest temperature      : 44
Lowest temperature       : 0
Max operating temperature : 70
Min operating temperature : 0
Estimated offset         : 0

#
```

Max operating temperature is defined as “Throttle Start” in P4800X/P4801X/D4800X product specification.

§



4 Troubleshooting Tips

4.1 Additional SMART Attribute Monitoring

- “C7 (CRC Error Count)” Count

Total number of PCIe Interface CRC errors encountered, as specified in PCIe Link Performance Counter Parameter for “Bad TLP.”

```
# nvme intel smart-log-add /dev/nvme1 | grep "key\|crc"
key                               normalized raw
crc_error_count                   : 100%      0
```

Alternatively, you can obtain this information using the Intel SSD Data Center Tool:

```
# isdct show -sensor -intelssd 1 |grep Crc
CrcErrorCount : 0
```

Note: As stated, CRC error could be contributed by PCIe link rather the drive itself. If your SSD already contained CRC an error count record, we suggest you check if the count number increased between two test cases.

4.2 Thermal Throttling Status

- “Thermal Throttle Status”

If the drive has encountered a performance issue, it could be due to many different reasons, once of which can be SSD thermal throttling event.

Here is the way to check if your system has increased thermal throttling event:

```
# nvme intel smart-log-add /dev/nvme1 |grep "key\|thermal"
key                               normalized raw
thermal_throttle_status           : 100%      0%, cnt: 1
```

Alternatively, you can obtain this information using the Intel SSD Data Center tool:

```
# isdct show -intelssd -smart |grep ThrottlingEventCount
ThrottlingEventCount : 1
```

If this count number increased in between test cases, you will have to consider either drive itself issue or server system fan speed control issue.

4.3 PCIe Link Speed and Link Width

- “PCIe Link Speed/ PCIe Link Width”

Another performance related case could be your PCIe link stability.

To check if your drive is running at the expected PCIe link speed or link width, you can use Intel SSD Data Center Tool to check:

```
# isdct show -a -intelssd 1 |grep Link
PCIlinkGenSpeed : 3
PCIlinkWidth : 4
```

With this information, you may able find clues on possible performance issues. It will help you troubleshoot server platform compatiability, especially signal integration related issues.



Additionally, this Linux command will help verify both PCIe speed and PCIe link width:

```
# lspci -s 10002:01:00.0 -vvv |grep Lnk
LnkCap: Port #0, Speed 8GT/s, Width x4, ASPM L0s, Exit Latency L0s <4us, L1 unlimited
LnkCtl: ASPM Disabled; RCB 64 bytes Disabled- CommClk-
LnkSta: Speed 8GT/s, Width x4, TrErr- Train- SlotClk+ DLActive- BWMgmt- ABWMgmt-
LnkCtl2: Target Link Speed: 8GT/s, EnterCompliance- SpeedDis-
LnkSta2: Current De-emphasis Level: -3.5dB, EqualizationComplete+, EqualizationPhase1+
```

“Speed 8GT/s” means your drive is running at PCIe Gen3.0 (PCIe Gen4.0 will be 16GT/s). “Width x4” means your SSD and backplane work in PCIe X4 width. If it drops to X2, it means you should expect to get half the performance.

4.4 Get Nlog

If you cannot root cause an issue based on SMART attribute readings, Intel engineer can help if you provide the Nlog, which is an Intel internal log. This command is only applicable if Telemetry is not supported on your current firmware revision. If Telemetry is supported, please see the “Retrieve Telemetry Log” in the next section.

Example of pulling Nlog using nvme -cli::

```
# nvme intel internal-log /dev/nvme0
Successfully wrote log to Nlog_FUKS7502001W1P5CGN.bin
# ls -al Nlog_FUKS7502001W1P5CGN.bin
-rw-r--r--. 1 root root 89260944 Jan 29 21:24 Nlog_FUKS7502001W1P5CGN.bin
```

You can zip this file with password protection and send to an Intel representative along with the SMART attribute reading mentioned above.

4.5 Retrieve Telemetry Log

P4800X/P4801X (after firmware version E2010475), D4800X supports retrieving a Telemetry log from the drive for debug purposes. nvme-cli, via the following command, can be used to retrieve a telemetry log. Once gathered please send back to Intel for issue escalation. Example:

```
# nvme telemetry-log /dev/nvme0 --output-file=telemetry_log.bin
```

4.6 Latency Tracking

D4800X support latency tracking while P4800X/P4801X do not support this feature.

Latency tracking is useful for users to understand if latency issues are coming from the drive itself or the server system level.

Here is the command to perform the latency tracking:

Enable latency tracking:

```
# isdctl set -intelssd X LatencyTrackingEnabled=True
```

Read latency statics after the test:

```
# isdctl show -intelssd X -latencystatistics
```



5 Manageability

5.1 Firmware Update and Reset

In order to troubleshoot and resolve issues, you'll need to have firmware update activity.

Download firmware update using NVMe-cli tool. Examples:

```
# nvme fw-download /dev/nvme0n1 -f E2010475_EB3B0438_WFEM01M0_signed.bin
# nvme fw-commit /dev/nvme0n1 -s 1 -a 1
```

Perform NVMe reset on the P4800X/P4801X for new firmware to take effect.

```
# nvme reset
```

When updating firmware on DC D4800X, an NVMe subsystem reset is required.

```
# nvme subsystem-reset /dev/nvme0
```

Depending on firmware upgrade requirements, you may need to reboot or power cycle the system for new firmware to take effect.

```
# reboot
```

You can update your Intel® Optane™ DC SSD firmware using Intel SSD Data Center Tool (DCT). After downloading the latest Intel SSD DCTL, update to the latest firmware using this bundle solution.

```
# isdct load -intelssd 1
# nvme reset /dev/nvme1
```

5.2 Change Power Governor Mode

For P4800X/P4801X, change power mode is not generally recommended since it also decreases performance.

Intel offers a power governor as a way for customers to change the power mode of the P4800X and P4801X for troubleshooting *and* testing. P4800X and P4801X support two power modes: 00h & 01h, check for "set feature" function details on exact power consumption of different power governor settings.

Following are examples to set and verify your power governor settings for P4800X and P4801X:

Confirm your current power governor mode:

```
# nvme get-feature /dev/nvme1n1 -f 0xc6
get-feature:0xc6 (Unknown), Current value:00000000
```

Or use ISDCT tool:

```
# isdct show -d PowerGovernorMode -intelssd 1
- Intel Optane(TM) SSD DC P4800X Series PHKE8375000B375AGN -
PowerGovernorMode : 0
```



Set power governor mode from 00h to 01h and reset the controller to take effect:

```
# nvme set-feature /dev/nvme0n1 -f 0xC6 -v 0x1
set-feature:c6 (Unknown), value:0x000001
# nvme reset /dev/nvme1
```

You can achieve same using ISDCT:

```
# isdct set -intelssd 1 PowerGovernorMode=1
# nvme reset /dev/nvme1
```

5.3 Format the SSD with Different Sector Size

By default, P4800X/P4801X/D4800X comes with 512B sector size, but you can format the drive with a different variable sector size.

Display by default variable by nvme-cli, here it is 512B:

```
# nvme list |grep 375
Node          SN          Model          Namespace Usage
Format       FW Rev
-----
/dev/nvme1n1  PHKE8375000B375AGN  INTEL SSDPE21K375GA  1
375.08 GB / 375.08 GB  512 B + 0 B  E2010435
```

Depending on which variable sector size you want to format as, here are examples:

```
# nvme format /dev/nvme0n1 -n 1 -l 0 -t 900000 # formatting to 512B
# nvme format /dev/nvme0n1 -n 1 -l 3 -t 900000 # formatting to 4KB
# nvme format /dev/nvme0n1 -n 1 -l 6 -m 1 -p 0 -i 1 -t 900000 # formatting to 4KB+128B (LBAF=6)
```

Note: For variable sector size other than 512B and 4K, if ms != 1, it will fail for P4800X/P4801X.

Note: Please contact [Intel Customer Support](#) for details on supported variable sector sizes in the Intel® Optane™ DC SSD P4800X/P4801X/D4800X.

5.4 Endurance Estimation

For P4800X/P4801X/D4800X, you can determine how much drive endurance has been consumed using “percentage_used” SMART attribute. For example:

```
# nvme smart-log /dev/nvme0 | grep percent
percentage_used          : 2%
```

By design, this SMART attribute will be zero when the SSD is first installed, and may remain at zero for a while. For more details on why this occurs, see this Technical Advisory:

<https://www.intel.com/content/www/us/en/support/articles/000033326/memory-and-storage/data-center-ssds.html>

As detailed in Technical Advisory referenced above, it may take multiple weeks of continuous writes to reach 1% (“host_bytes_written” estimated >1.8PBW). Once it reaches 1% you can begin using your workload to estimate the endurance capability.



By design SMART attribute ADh in log page 202 does not work for P4800X/P4801X/D4800X. Because the endurance analyzer does not work for P4800X/P4801X/D4800X product series, you'll have to monitor delta of "percentage_used" and delta of "host_bytes_written" between real workload tests. (Intel suggests a minimum of 60 minutes, or even more, for large write workloads).

When real workload is running, use the following suggested examples to perform online monitoring of the "delta" mentioned above:

```
# more smart.sh
for ((i=1;i<200000;+i))
do
echo $i >> smart.log
date >> smart.log
nvme smart-log /dev/nvme2 |grep percentage_used >> smart.log
nvme intel smart-log-add /dev/nvme2 |grep host_bytes_written >> smart.log
sleep 10
done
```

You will get results similar to the following:

```
# more smart.sh
```

Time Consumed (s)	Delta (hours)	Percentage Used	Host Bytes Written	Delta (GB)
837980	233	1	57,894,690	1,809,209
949040	31	2	65,479,150	237,014
1058060	30	3	72,924,904	232,680
1167840	30	4	80,426,164	234,414
1278100	31	5	87,816,863	230,959
1386710	30	6	95,232,844	231,749
1494130	30	7	102,572,275	229,357
1602570	30	8	109,982,031	231,555
1714140	31	9	117,605,519	238,234

Based on the data of "Delta (hours)", "Percentage used", "Delta (GB)" and your SSD's PBW endurance data, you will be able to know the selected Intel® Optane™ SSDs sustained endurance in your unique workload conditions.