



Over-Provisioning NAND-Based Intel® SSDs for Better Endurance

How over-provisioning improves NAND-based Intel® SSD's endurance and performance in real-world workloads.

Authors Executive Summary

Zhdan Bybin
Senior Application Engineer

Mohammed Khandaker
Technical Marketing Engineer

Monika Sane
Design Engineer

Graham Hill
Application Engineer

Over-provisioning (OP) increases SSD endurance by allowing extra space for the flash controller to manage incoming data. Over-provisioning improves wear-leveling and random write performance, and reduces the write amplification factor (WAF), thereby improving the endurance of NAND-based SSDs.

Background

Writing data to an SSD is significantly different than writing data to a hard disk drive (HDD). Writing to a spinning HDD is performed by magnetizing and demagnetizing sectors on a thin layer of metal mounted on a circular platter, and information is written in binary format (value 1 or 0). Overwriting an HDD is simply a matter of changing the magnetization value, allowing the data to be directly overwritten.

However, overwriting NAND-based SSDs is a much more complicated process. The underlying technology media, NAND flash, is primarily an array of floating gate transistors where electrons are trapped and stored in a floating gate. Applying an electromagnetic field causes the trapped electrons to tunnel through the insulator into the substrate.

The presence or absence of electrons on the floating gate determines value 1 or 0 in the case of the single level cell (SLC) NAND, and the number of electrons determines the voltage level in multiple level cell (MLC, TLC, etc.) NAND. Due to the physical and electrical functioning properties of the floating gate, existing data must be erased before new data can be written to the cell. The process to write data to SSD media is called program (P), and the process to erase data from it is called erase (E), together called the P/E cycle of the NAND flash. Every P/E cycle leads to very slight wear of the media. Moreover, there is a finite, predefined number of P/E cycles, after which the media becomes unusable.

Another important difference is how the data is organized and partitioned logically on the media. Flash memory is divided into blocks, each block contains pages, and each page is a collection of memory cells. Data is written to an SSD in pages, however the erasing of data happens in blocks. Thus, in order to overwrite a block containing valid and invalid (discarded by the host) pages, the valid pages must be temporarily moved elsewhere on the media, so that the entire block can be erased, then overwritten. This temporary data movement creates the undesirable phenomenon called write amplification (WA). It is undesirable because the actual amount of data written to the SSD media is larger than the amount of data that the host intended to write, thereby leading to further media wear.

To minimize wear and increase the lifetime of the SSD media, there are various algorithms and techniques, such as wear-leveling, garbage collection, etc. Over-provisioning is one such technique available to users of NAND-based SSDs, and is the focus of this document. Over-provisioning an SSD means to make more of the SSD's capacity available to the controller than what was initially designed in. This technique increases SSD endurance by allowing extra buffer space for the flash controller to manage incoming data—it improves wear-leveling and random write performance, and ultimately decreases the WAF.

Table of Contents:

Executive Summary..... 1
Background..... 1
Introduction to Over-Provisioning..... 2
Methods of Over-Provisioning..... 2
Flexible Capacity and Endurance
Examples..... 4
Conclusion..... 7

In contrast, Intel® Optane™ SSDs don't use garbage collection or wear-leveling techniques because of the entirely different endurance mechanism and superior underlying nature of the Intel® Optane™ technology. In environments where very high endurance is required, Intel recommends using Intel® Optane™ technology-based SSDs such as Intel® Optane™ SSD DC P4800X Series, or Intel® Optane™ SSD 900P Series, with their exceptional endurance rating of up to 60 drive writes per day (DWPD).

Introduction to Over-Provisioning

As previously discussed, SSD endurance is the ability to withstand the repeated writing of data, and the ability to retain that data for a period of time (data retention). NAND endurance is measured in terms of terabytes written (TBW) or DWPD and is dependent on a number of factors: the maximum media P/E cycles specification, capacity, workload, and firmware (FW) techniques that are employed to lower the WAF. Under demanding enterprise workloads, NAND-based SSDs can wear out quicker due to higher WAF which is defined as the ratio of NAND writes/host writes. A minimal WAF—as close to 1.0 as possible—is desirable to minimize wear and increase SSD lifetime. In some cases, WAF can be < 1 if the SSD controller has compression mechanisms built in.

$$WAF = \frac{NAND\ WRITES}{HOST\ WRITES}$$

WAF is different for each workload. One of the more demanding workloads is the one described in the Joint Electron Device Engineering Council (JEDEC) endurance specification. The endurance workload can be obtained at www.jedec.org

Intel measures and publishes its data center SSD endurance in accordance to the JEDEC specification (JESD 218A, JESD219) as well as for sequential write workload. However, because this is not an industry requirement, other SSD vendors may report endurance numbers measured under the 100% 4KB random write workload. Generally this workload has a WAF that is ~10% lower than a JESD219 workload.

$$\text{Endurance (Peta Bytes Written, PBW)} = \frac{\text{Raw Drive Capacity} \cdot \text{NAND PE Cycles}}{WAF}$$

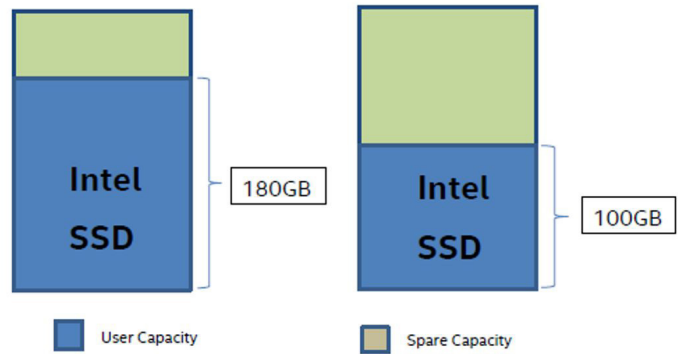
$$\text{Endurance (DWPD for 5 years)} = \frac{\text{Endurance (PBW)}}{\text{Rated drive density} \cdot 5 \cdot 365}$$

To provide an SSD controller the spare capacity needed to move the data around when programming or erasing partially invalid pages or blocks, each SSD has “factory over-provisioned” area. This area is not addressable by the host/user and may vary in size depending on SSD model and capacity. Effective size of this area can influence the WAF for a workload that uses 100% logical block address (LBA) span of the drive.

Sacrificing user-addressable capacity by manually increasing the effective SSD's spare area allocation will result in endurance gains due to WAF decrease, as well as improvements in random write performance and quality of service (QoS) due to fewer SSD controller housekeeping activities. This method of OP is similar to the HDD concept called “short stroking” the drive.

Methods of Over-Provisioning

Over-provisioning should be performed on an SSD in a completely clean state. This can be an SSD that is fresh out-of-the-box or an SSD on which a secure erase has been performed.



The three most common methods of over-provisioning an SSD are as follows:

1. Limiting the logical volume capacity during partitioning in OS (user will see full capacity in Disk Manager or fdisk).
2. Limiting the Maximum LBA on the drive level, so that in OS, it will appear as a lower-capacity drive.
3. Limiting an application to use only a certain LBA range.

Please note, method 3 above will not work for the scenario in which the filesystem is deployed on full LBA range.

For method 2 above, while working with Intel® SSDs for the Data Center, customers can use either Intel or 3rd-party tools. Following are the tools that work with native over-provisioning:

- a. The easiest tool to use for over-provisioning is Intel® SSD Data Center Tool (Intel® SSD DCT), which supports both SATA and NVMe* Intel® SSDs, in Linux* and Windows* environments - <https://www.intel.com/content/www/us/en/support/articles/000006289/memory-and-storage.html>

Intel SSD DCT User Guide can be found here: <https://www.intel.com/content/www/us/en/support/memory-and-storage/000020016.html?wapkw=data+center+tool+user+guide>

The command sequence for over-provisioning:

```
# sudo isdct show -intelssd (To get index of the drive)
# sudo isdct delete -intelssd 0 (Note: ATA security needs to be NOT frozen and NOT in the locked state, please refer to https://www.intel.com/content/www/us/en/support/articles/000006094/memory-and-storage.html ).
# sudo isdct set -intelssd X
MaximumLBA=(xGB\x%\LBA\'native')
(example: isdct set -intelssd 0
MaximumLBA=80%)
```

Power cycle the drive or reboot the system.

```
[root@localhost ~]# isdct show -intelssd
- Intel SSD DC S4500 Series PHYS73120349240AGN -
Bootloader : Property not found
DevicePath : /dev/sg0
DeviceStatus : Healthy
Firmware : SCV10111
FirmwareUpdateAvailable : The selected Intel SSD contains current firmware as of
this tool release.
Index : 0
ModelNumber : INTEL SSDSC2KB240G7
ProductFamily : Intel SSD DC S4500 Series
SerialNumber : PHYS73120349240AGN

[root@localhost ~]# isdct set -intelssd 0 maximumlba=200GB
Set MaximumLBA successful. Please power cycle the device.
[root@localhost ~]#
```

The way we recognize an over-provisioned drive is illustrated in the picture below (notice how MaximumLBA < NativeMaxLBA) :

```
# sudo isdct show -all -intelssd 0
```

```
MaximumLBA : 390625000
MediumPriorityWeightArbitration : Device does not support this command set.
ModelNumber : INTEL SSDSC2KB240G7
NVMePowerState : Device does not support this command set.
NativeMaxLBA : 468862127
OEM : Generic
PLITestTimeInterval : 10000 minutes
PhySpeed : 6.0 Gbps
PhysicalSectorSize : 4096 bytes
PhysicalSize : 200000000512
PowerGovernorAveragePower : 4000 milliwatts
```

An over-provisioned SSD can be reverted to its native capacity as shown in image below. Also, it is best practice to put the SSD in standby mode before power cycling it. This can be done as follows:

```
# sudo isdct set -intelssd 0 maximumlba=native
```

```
[root@localhost ~]# isdct set -intelssd 0 maximumlba=native
Set MaximumLBA successful. Please power cycle the device.
[root@localhost ~]# isdct start -intelssd 0 -standby

- StandbyImmediate PHYS73120349240AGN -
Status : Completed successfully.
```

- b. HDParm* (latest version) 9.49 and above - 3rd party tool, Linux* only, only SATA drives

Just as when using Intel SSD DCT, secure erase must be performed on the SSD before over-provisioning it.

```
# sudo hdparm --user-master u --security-set-pass 123 /dev/sdX
# sudo hdparm --security-erase 123 /dev/sdX
```

The command for over-provisioning would be:

```
# sudo hdparm -N /dev/sdX (To find the maximum sector count)
# sudo hdparm -NpXXXXXXXXXX -yes-i-know-what-i-am-doing /dev/sdX (This enables host protected area and sets the number of visible sectors to the count appearing immediately after "-Np")
```

- c. nvme-cli – open source tool developed by an Intel engineer, Linux only, only NVMe drives

Endurance Calculations for NAND-Based Intel® SSDs

Intel SSDs for the data center have timed workload SMART indicators (E2h, E3h, E4h) that allow users to easily calculate the lifetime of the SSD under any given real-world workload using actual NAND wear statistics. The E2h attribute measures the wear endured by the SSD during the timed workload; E3h tracks the workload's read/write ratio; and E4h reports the number of minutes spent during the workload.

These attributes must be reset prior to applying the characteristic workload, by issuing the SMART Execute Immediate ATA command (for SATA SSDs) or NVMe Vendor Unique Set Features D5h command (for NVMe SSDs). Intel recommends applying a full characteristic cycle of the expected workload, i.e., 8-hours, 24-hours, 1 week etc., while the minimum requirement is a 60 minute workload with enough write pressure to make E2h increase its counter by 1, which is equivalent to ~0.001% overall media wear.

If Intel SSD DCT is used, the tool will do all the calculations on the user's behalf.

The commands are as follows:

```
# sudo isdct set -intelssd 0 EnduranceAnalyzer=reset (to reset the attributes prior to the test)
```

Once reset, raw values for E2h/E3h/E4h will show '65535' and will stay the same until the representative workload is applied for at least 60 minutes.

```
[root@localhost fio]# isdct set -intelssd 0 enduranceanalyzer=reset
Set EnduranceAnalyzer successful. Completed successfully.
[root@localhost fio]# isdct show -smart E2 -intelssd 0

- SMART Attributes PHYS73120349240AGN -

- E2 -

Action : Pass
Description : Timed Workload - Media Wear
ID : E2
Normalized : 100
Raw : 65535
Status : 50
Threshold : 0
Worst : 100
```

After the test is done, Intel SSD DCT can provide the estimated life calculation using the following command:

```
# sudo isdct show -all -intelssd 0
```

Without Intel tools, the same calculations can be made using E2h/E3h/E4h 'raw' values. Here is a real-world application example after testing a 240GB Intel® SSD DC S4500:

Attribute	Attribute Name	Value @ test start	Value @ test end
E1	Host Writes	14308	26173
E2	Timed Workload Media Wear	65535	41
E3	Timed Workload Host Read/Write Ratio	65535	65
E4	Timed Workload Timer	65535	8624

Based on this data, we can calculate following:

$$\text{Host writes} = E1_{\text{end}} - E1_{\text{start}} = 11865 \text{ units} * 32\text{MB per unit} = 370 \text{ GB}$$

$$\text{Time spent during test} = 8624 \text{ min} / 60 / 60 = 6 \text{ days}$$

$$\text{Workload read/write ratio} = 65\%/35\% \text{ R/W}$$

$$\text{Media wear} = 0.040 \%$$

$$\text{Estimated drive's life remaining} = 6 \text{ days} * 100 \% / 0.04 \% = 15000 \text{ days} = 41.1 \text{ years}$$

Flexible Capacity and Endurance Examples

In the next section we're providing examples of the capacity/endurance tradeoffs for latest Intel® SSDs. The following tables provide over-provisioning (OP) values for popular endurance levels, based on simulation calculations. The tables provide answers to questions like how to get from 1 DWPD to 3 DWPD, or from 3 DWPD to 5 DWPD, or for popular over-provisioning and capacity levels, i.e. 10%, 20% or 400GB, etc. Different endurance levels are shown if sequential workload is dominant compared to JEDEC workloads or 4K random write (RW) workloads.

Flexible Capacity and Endurance Calculations for Intel® SSD D3-S4510 Series:

SKU	Size (GB)	No OP JEDEC Endurance		No OP Seq Write Endurance		Size (GB)	10% OP 4K RW Endurance	
D3-S4510		PBW	DWPD	PBW	DWPD	10%OP	PBW	DWPD
2.5"	240	0.9	1.9	2.4	5.5	216	1.2	3.04
	480	1.2	1.3	4.7	5.4	350	3.18	4.98
	960	3.4	1.9	10.3	5.9	864	4.94	3.13
	1920	7.1	2	21.1	6	1728	9.49	3.01
	3840	13.1	1.9	42	6	3456	15.8	2.5

Size (GB)	2DWPD 4K RW Endurance		Size (GB) 20% OP	20% OP 4K RW Endurance		Size (GB)	3DWPD 4K RW Endurance	
	PBW	DWPD		PBW	DWPD		PBW	DWPD
237	0.87	2	200	1.44	3.93	183	1.67	4.99
448	1.64	2	400	2.44	3.34	349	3.2	5.01
950	3.48	2.01	800	5.9	4.04	740	6.75	4.99
1875	6.86	2	1600	11.34	3.88	1450	13.35	5.04
3550	13.1	2.02	3000	21.57	3.94	2700	24.97	5.06

Flexible Capacity and Endurance Calculations for Intel® SSD D3-S4510 (M.2) Series:

SKU	Size (GB)	No OP JEDEC Endurance		No OP Seq Write Endurance		Size (GB)	10% OP 4K RW Endurance	
D3-S4510		PBW	DWPD	PBW	DWPD	10%OP	PBW	DWPD
M.2	240	0.9	1.9	2.4	5.5	216	1.2	3.04
	480	1.2	1.3	4.7	5.4	432	1.91	2.41
	960	2.3	1.3	9	5.1	864	3.69	2.34

Size (GB)	3DWPD 4K RW Endurance		Size GB 20% OP	20% OP 4K RW Endurance		Size (GB)	5DWPD 4K RW Endurance	
	PBW	DWPD		PBW	DWPD		PBW	DWPD
216	1.2	3.04	200	1.44	3.93	182	1.68	5.06
410	2.28	3.04	400	2.44	3.34	349	3.19	5.01
815	4.54	3.05	800	4.77	3.27	690	6.36	5.05

Flexible Capacity and Endurance Calculations for Intel® SSD D3-S4610 Series:

SKU	Size (GB)	No OP JEDEC Endurance		No OP Seq Write Endurance		Size (GB)	10% OP 4K RW Endurance	
D3-S4610		PBW	DWPD	PBW	DWPD	10%OP	PBW	DWPD
2.5"	240	1.4	3.3	2.7	6.2	216	1.76	4.46
	480	3	3.5	6	6.8	350	7.47	4.77
	960	6	3.4	11	6.3	864	7.16	4.54
	1920	10	3.1	22	6.3	1728	12.67	4.01
	3840	22	3.1	44	6.3	3456	21.63	3.43

Size (GB)	3DWPD 4K RW Endurance		Size GB 20% OP	20% OP 4K RW Endurance		Size (GB)	5DWPD 4K RW Endurance	
	PBW	DWPD		PBW	DWPD		PBW	DWPD
207	1.88	4.98	200	1.98	5.41	143	2.61	9.99
420	3.85	5.01	400	4.12	5.64	294	5.39	10.04
830	7.63	5.04	800	8.04	5.5	583	10.66	10.01
1590	14.51	4.99	1600	14.38	4.92	1120	20.49	10.02
2950	26.91	4.99	3000	26.42	4.82	2085	38.13	10.01

Flexible Capacity and Endurance Calculations for Intel® SSD P4510 Series:

SKU	Size (GB)	No OP JEDEC Endurance		Sequential Write Endurance		Size (GB)	10% OP 4K RW Endurance	
P4510		PBW	DWPD	PBW	DWPD	10% OP	PBW	DWPD
U.2 15mm	1000	1.92	1.05	5.62	3.00	900	2.78	1.69
	2000	2.61	0.7	10.90	2.90	1800	4.43	1.35
	4000	6.3	0.85	21.84	2.80	3600	10.02	1.52
	8000	13.88	0.9	44.25	3.00	7200	21.3	1.62
M.2 110mm	1000	0.98	0.54	5.13	2.81	900	1.91	1.16
	2000	1.95	0.53	10.26	2.81	1800	3.79	1.15

SKU	Size (GB)	20% OP 4K RW Endurance		Size (GB)	3DWPD 4K RW Endurance	
P4510	20% OP	PBW	DWPD	Size (GB)	PBW	DWPD
U.2 15mm	800	3.56	2.44	735	4.02	3
	1600	6.13	2.1	1400	7.63	2.98
	3200	13.29	2.27	2870	15.72	3
	6400	27.74	2.37	5840	31.86	2.99
M.2 110mm	800	2.78	1.9	680	3.69	2.97
	1600	5.52	1.89	1350	7.4	3

Flexible Capacity and Endurance Calculations for Intel® SSD P4610 Series:

SKU	Size (GB)	No OP JEDEC Endurance		No OP Seq Write Endurance		Size (GB)	10% OP 4K RW Endurance	
P4610		PBW	DWPD	PBW	DWPD	10%OP	PBW	DWPD
2.5" 15mm	1600	12.25	4.19	21.18	7.25	1440	14.69	5.59
	3200	21.85	3.74	41.03	7.02	2880	26.85	5.1
	6400	36.54	3.13	78.09	6.68	5760	47.01	4.47
	7680	44.25	3.15	93.97	6.7	6900	56.94	4.52

Size (GB)	5DWPD 4K RW Endurance		Size (GB)	20% OP 4K RW Endurance		Size (GB)	10DWPD 4K RW Endurance	
	PBW	DWPD	20% OP	PBW	DWPD	Size (GB)	PBW	DWPD
1500	13.8	5.04	1200	17.92	8.18	1060	19.52	10.08
2900	26.55	5.01	2500	32.26	7.07	2050	37.64	10.05
5500	50.99	5.08	5000	58.17	6.37	3925	71.68	10
6650	60.76	5	6000	70.07	6.39	4720	86.21	10

Flexible Capacity and Endurance Calculations for Intel® SSD D5-P4320/D5-P4420/D5-P4326 Series:

SKU	Size (GB)	No OP JEDEC / 16K ¹ Endurance		No OP Seq Write Endurance		Size (GB)	10% OP 4K RW Endurance	
		PBW	DWPD	PBW	DWPD	10%OP	PBW	DWPD
P4320	7680	2.8	0.2	12.3	0.88	6900	4.94	0.39
P4420	7680	6.7	0.48	28.1	2.01	6900	9.89	0.78
P4326	15360	6.2	0.18	25.4	0.9	13800	9.34	0.37

Size (GB)	5DWPD 4K RW Endurance		Size (GB)	20% OP 4K RW Endurance		Size (GB)	10DWPD 4K RW Endurance	
	PBW	DWPD	20% OP	PBW	DWPD	Size (GB)	PBW	DWPD
5000	9.31	1.02	6100	6.94	0.62	3300	11.98	1.99
6500	11.9	1	6100	13.89	1.25	5050	18.41	2
10000	18.26	1	12200	13.39	0.6	6500	23.9	2.01

The Intel SSD DC D5-P4320/D5-P4420/D5-P4326 Series is designed with QLC-based NAND, compared to previous generations of drives with TLC-based NAND. The endurance values for these SSDs differ from those of the previous generations, because of these variations we've provided more data in the charts and graphs below. This additional information is included to provide clearer understanding of the true endurance of these drives under different real-world workloads.

Write Workload	D5-P4320 7.68TB		D5-P4420 7.68TB		D5-P4326 15.36TB ¹	
	DWPD	TBW	DWPD	TBW	DWPD	TBW
100% 128KB Sequential	1.01	14,100	2.01	28,100	1.02	28,500
90% 128KB Sequential; 10% 4KB Random	1.01	14,100	2.01	28,100	1.01	28,300
80% 128KB Sequential; 20% 4KB Random	1.01	14,100	2.01	28,100	1.00	28,000
70% 128KB Sequential; 30% 4KB Random	1.01	14,100	2.01	28,100	0.98	27,400
50% 128KB Sequential; 50% 4KB Random	1.01	14,100	2.01	28,100	0.94	26,300
100% 16KB Random	0.24	3,300	0.48	6,700	0.21	5,800
100% 8KB Random	0.25	3,500	0.50	7,000	0.10	2,800
100% 4KB Random	0.28	3,900	0.58	8,100	0.05	1,400

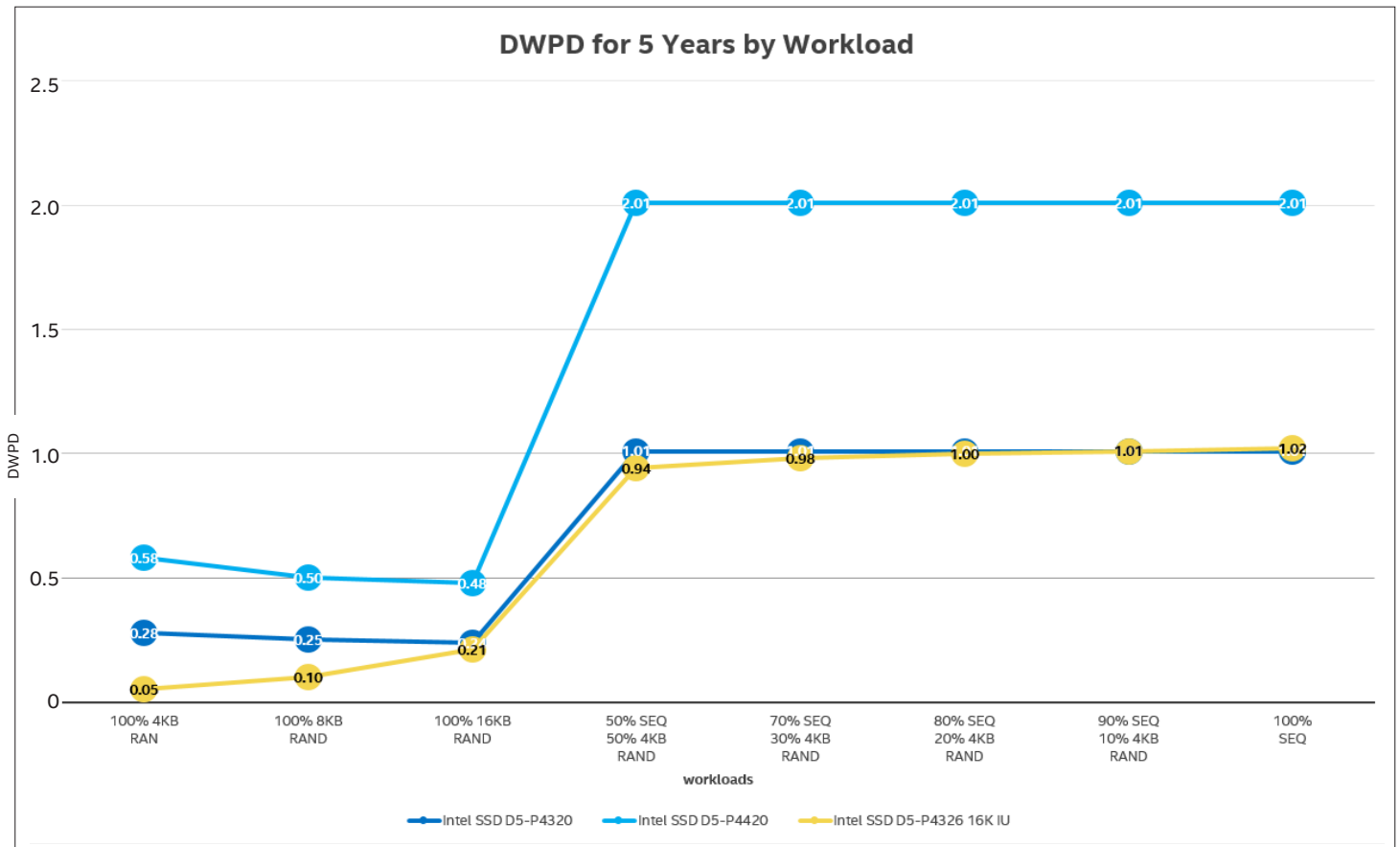


Chart 1. QLC Based SSDs Endurance per Workload

As briefly mentioned previously, over-provisioning will also positively affect random write performance of a NAND-based SSD if 100% LBA span is used, or if filesystem is applied to the whole addressable capacity of the drive. The smaller the default factory over-provisioned area is, the higher the impact will be. The chart below shows the impact of OP percentage on random write performance using a 240GB Intel® SSD DC S4600 as an example.

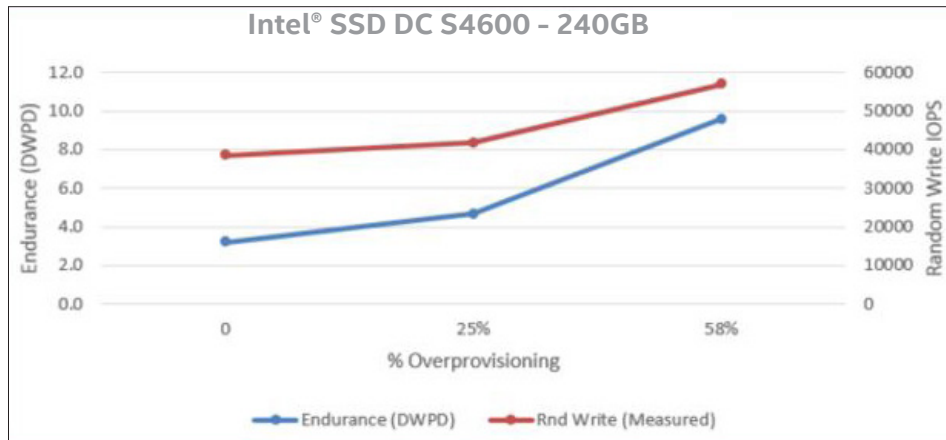


Chart 2. Intel® SSD DC S4600 Series Random Write Performance

Conclusion

As shown in the preceding formulas and tables, while sacrificing user-addressable capacity, over-provisioning provides a positive effect on NAND-based SSD endurance, WAF, and random write performance. In general, over-provisioning allows flexibility in an SSD's endurance and capacity where the user can go from a 1 DWP/D-rated SSD to 3 DWP/D, or from 3 DWP/D to 5, or even up to 10 DWP/D.

If higher endurance levels are required – for example, for hot data tier, or caching, or DRAM displacement—Intel® Optane™ SSDs can be used. Due to the different media nature, Intel Optane™ SSDs do not gain benefits with over-provisioning; they already have very high endurance, with up to 60 DWP/D rating.

ENDURANCE ↑

RANDOM WRITE PERF ↑

USER DRIVE CAPACITY ↓



For more information, visit intel.com/ssd

1. Intel® SSD D5-P4326 has a 16KB indirection unit (IU)

System Configuration for all performance testing: Intel® Xeon® CPU E5-2699 v4 @ 2.20GHz on Intel® S2600WT motherboard, Intel® C612 Chipset, BIOS Version SE5C610.6B.01.01.0019.101220160604 32GB DDR4, FIO version 2.18, CentOS 7.0, Kernel 4.8.6 (DAS patch). Testing by Intel.

Performance results are based on testing as of February 25, 2019 and may not reflect all publicly available security updates. See configuration disclosure for details. No product or component can be absolutely secure.

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors.

Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more complete information 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 computer system can be absolutely secure. Check with your system manufacturer or retailer or learn more at intel.com.

Results have been estimated or simulated using internal Intel analysis or architecture simulation or modeling, and provided to you for informational purposes. Any differences in your system hardware, software, workload, or configuration may affect your actual endurance numbers.

Intel, the Intel logo, Intel Optane and Intel Xeon are trademarks of Intel Corporation or its subsidiaries in the U.S. and/or other countries.

*Other names and brands may be claimed as the property of others.