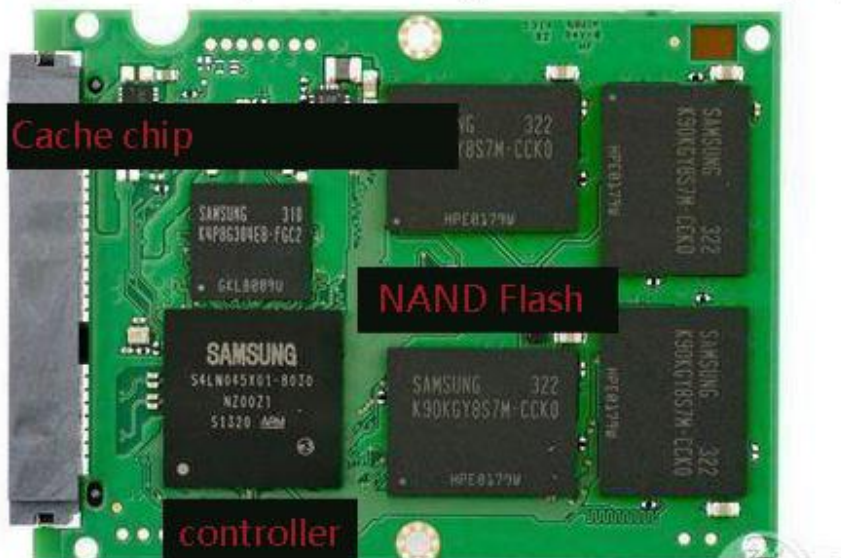


Solid State Drive

An SSD is a Solid-State Drive that is used to store data, files and applications, as well as to run computing devices. Solid state drives offer significant performance advantages over conventional hard disk drives (HDDs). HDD's work by way of a mechanical drive head that must physically move to access locations on a rapidly-spinning magnetic disk. That has rotating magnetic parts, read- write heads etc. When the computer sends a request to retrieve data, the disk and arm must each move to the appropriate location for the data to be collected and sent to the CPU for processing. SSD's, on the other hand, have no moving parts. An SSD is actually just a thin wedge of NAND flash memory - exactly same as the one that is there in your USB drive, but rather than being encased in a thin stick, SSD is put inside a shell or enclosure of 2.5 inches along with a SATA interface for enhanced performance levels.

SSD Internals(Take Samsung 840EVO for Example)



The major components inside an SSD are:

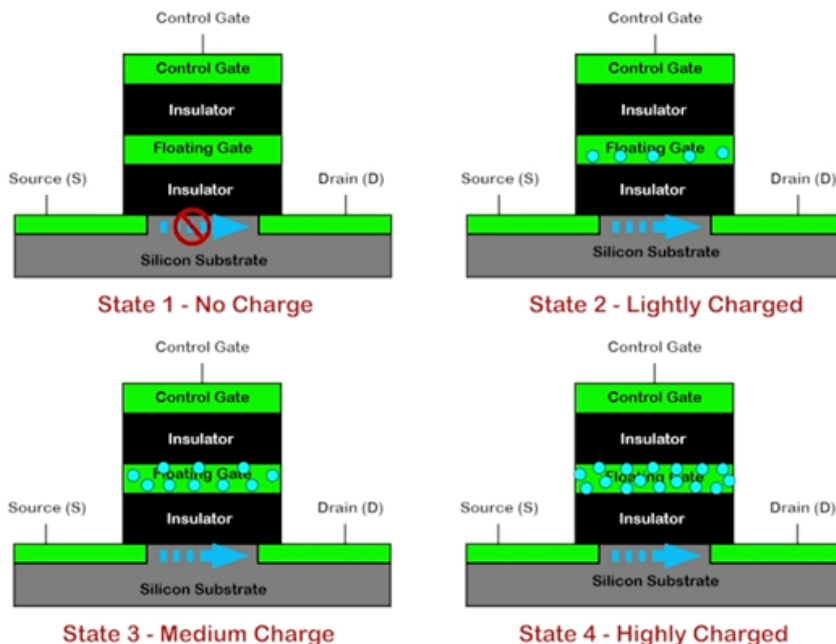
Outer Shell:

- The outer shell could be of metal or plastic and it helps in absorbing most of the heat from inside the flash memory.

- Although, SSDs don't contain any moving parts, they give off very little heat and emit no noise. This particular feature is vital in increased durability of an SSD.
- An SSD can withstand 10X more vibration than an HDD and up to 1500G of shock (compared to less than 70Gs for a typical HDD).
- SSDs exceed expectations in handling shock, vibration, and temperature extremes as well.

NAND FLASH:

Solid-state drives are called that specifically because they don't rely on moving parts or spinning disks. Instead, data is saved to a pool of NAND flash. NAND itself is made up of what are called floating gate transistors. Unlike the transistor designs used in DRAM, which must be refreshed multiple times per second, NAND flash is designed to retain its charge state even when not powered up. This makes NAND a type of non-volatile memory.



The diagram above shows a simple flash cell design. In NAND cell, the data are stored as charged state, Electrons are stored in the floating gate, which then reads as charged “0” or not-charged “1.” Yes, in NAND flash, a 0 means data is stored in a cell — it’s the opposite of how we typically think of a zero or one. NAND flash is organized in a grid. The entire grid layout is referred to as a block, while the individual rows that make up the grid are called a page. Common page sizes are 2K,

4K, 8K, or 16K, with 128 to 256 pages per block. Block size therefore typically varies between 256KB and 4MB.

One advantage of this system should be immediately obvious. Because SSDs have no moving parts, they can operate at speeds far above those of a typical HDD. The following chart shows the access latency for typical storage mediums given in microseconds.

	SLC	MLC	TLC	HDD	RAM
P/E cycles	100k	10k	5k	*	*
Bits per cell	1	2	3	*	*
Seek latency (μs)	*	*	*	9000	*
Read latency (μs)	25	50	100	2000-7000	0.04-0.1
Write latency (μs)	250	900	1500	2000-7000	0.04-0.1
Erase latency (μs)	1500	3000	5000	*	*
Notes	* metric is not applicable for that type of memory				
Sources	P/E cycles [20] SLC/MLC latencies [1] TLC latencies [23] Hard disk drive latencies [18, 19, 25] RAM latencies [30, 52] L1 and L2 cache latencies [52]				

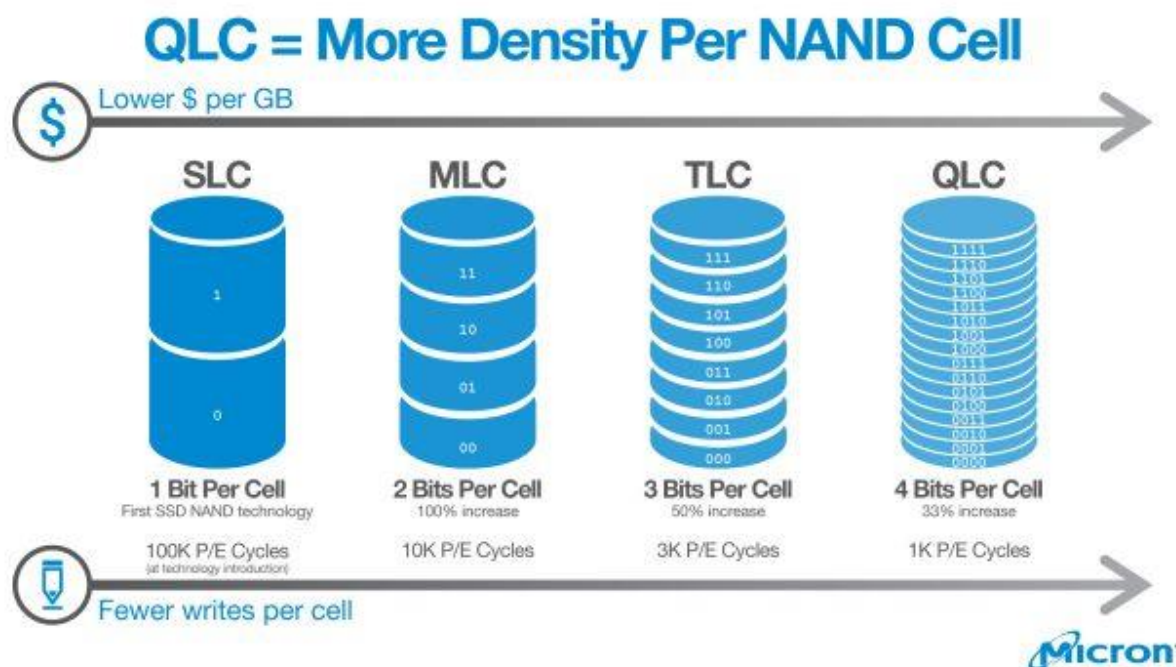
Image by [CodeCapsule](#)

NAND is nowhere near as fast as main memory, but it's multiple orders of magnitude faster than a hard drive. While write latencies are significantly slower for NAND flash than read latencies, they still outstrip traditional spinning media. How to understand SLC, MLC and TLC flash? To be Simple, the data in NAND flash is stored in each storage cell. SLC, MLC, and TLC are different bits of storage.

The difference between single-tier and multi-tier storage is the number of "bits" each NAND storage cell can store at a time. SLC (single-level Cell) stores only 1bit of

data per Cell, while MLC (multi-level Cell) stores 2 bits and TLC (Trinary-Level) stores 3 bits. A storage cell stored the more bits at a time and will has more capacity, which can save the cost of flash memory and increase NAND production. However, the state is difficult to discern when adding more data to each cell. In the same time, reliability, durability and performance will reduce.

There are two things to notice in the above chart. First, note how adding more bits per cell of NAND has a significant impact on the memory's performance. It's worse for writes as opposed to reads — typical triple-level-cell (TLC) latency is 4x worse compared with single-level cell (SLC) NAND for reads, but 6x worse for writes. Erase latencies are also significantly impacted. The impact isn't proportional, either — TLC NAND is nearly twice as slow as MLC NAND, despite holding just 50% more data (three bits per cell, instead of two). This is also true for QLC drives, which store even more bits at varying voltage levels within the same cell.

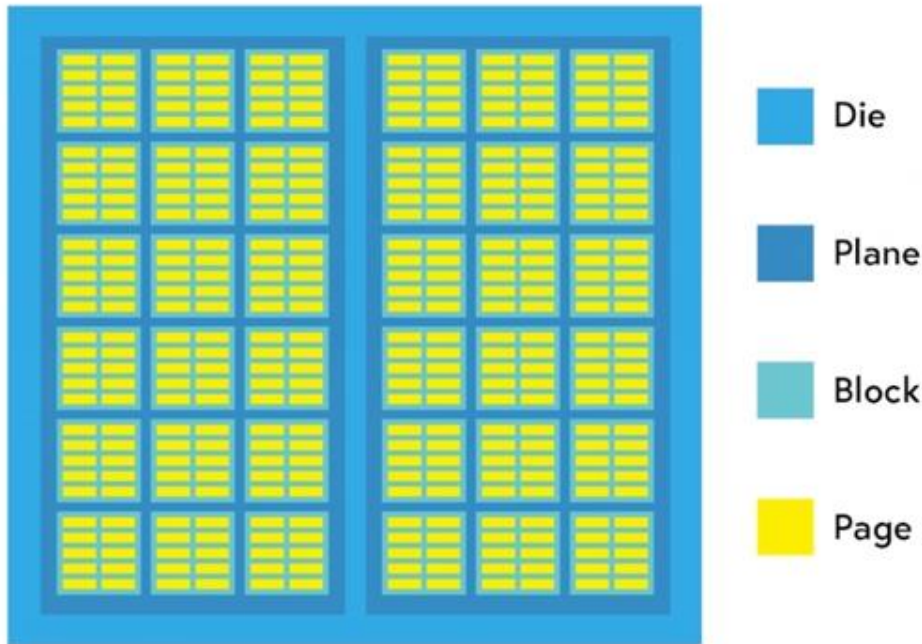


The reason TLC NAND is slower than MLC or SLC has to do with how data moves in and out of the NAND cell. With SLC NAND, the controller only needs to know if the bit is a 0 or a 1. With MLC NAND, the cell may have four values — 00, 01, 10, or 11. With TLC NAND, the cell can have eight values, and QLC has 16. Reading the proper value out of the cell requires the memory controller to use a precise voltage to ascertain whether any particular cell is charged. But the think I like to mention that, in most of the SSD, there are some SLC cells that manufacturer provide with, embedded on the chip to perform faster processing to some level.

Internal Structure of NAND Flash

This is a figure of single die.

NAND Flash Die Layout



The die is marked as blue color here. There are multiple dies in a single NAND flash chip. Die is the smallest stage where a single command or instruction can be executed independently. In a single die, there may be one or two planes. The planes are colored as deep blue In this figure. The planes can run a single process concurrently or you can say run a single process into multiple planes at a time. In planes, there are multiple blocks, here blocks are colored as paste. In this block level, data erase operation is performed. In each single block, there are multiple pages there, in this figure, they are colored as yellow. Data write operation is performed at page level. Remember one thing, data erase operation is done at block level.

Reads, Writes, and Erasure

One of the functional limitations of **SSDs** is while they can read and write data very quickly *to an empty drive*, overwriting data is much slower. This is because while SSDs read data at the page level (meaning from individual rows within the NAND memory grid) and can write at the page level, assuming surrounding cells are empty, they can only erase data at the block level. This is because the act of erasing NAND flash requires a high amount of voltage. While you can theoretically erase NAND at

the page level, the amount of voltage required stresses the individual cells around the cells that are being re-written. Erasing data at the block level helps mitigate this problem.

The only way for an SSD to update an existing page is to copy the contents of the entire block into memory, erase the block, and then write the contents of the old block + the updated page. If the drive is full and there are no empty pages available, the SSD must first scan for blocks that are marked for deletion but that haven't been deleted yet, erase them, and then write the data to the now-erased page. This is why SSDs can become slower as they age — a mostly-empty drive is full of blocks that can be written immediately, a mostly-full drive is more likely to be forced through the entire program/erase sequence.

If you've used SSDs, you've likely heard of something called **“garbage collection.”** Garbage collection is a background process that allows a drive to mitigate the performance impact of the program/erase cycle by performing certain tasks in the background. The following image steps through the garbage collection process.

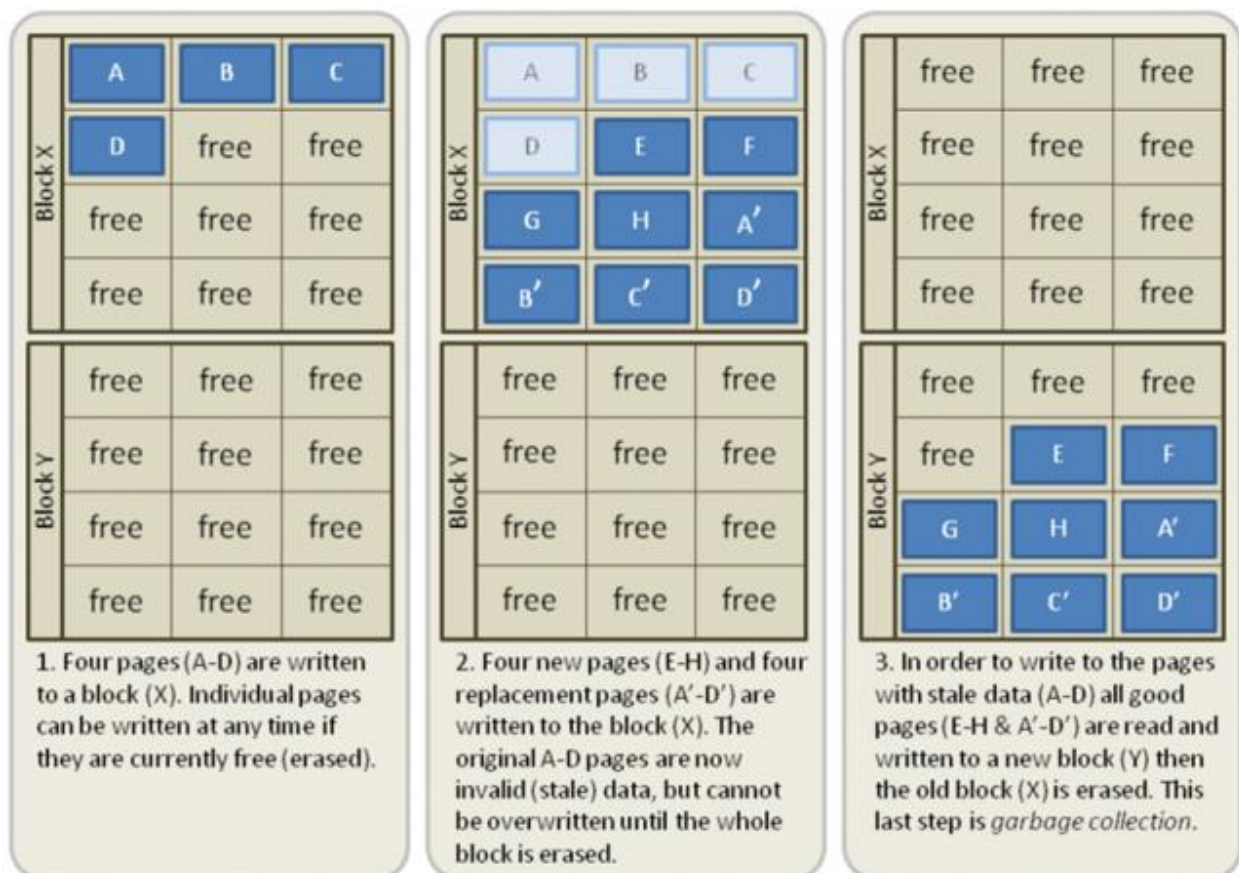


Image courtesy of Wikipedia

The next concept I want to discuss is TRIM. When you delete a file from Windows on a typical hard drive, the file isn't deleted immediately. Instead, the operating system tells the hard drive it can overwrite the physical area of the disk where that data was stored the next time it needs to perform a write. This is why it's possible to undelete files (and why deleting files in Windows doesn't typically clear much physical disk space until you empty the recycling bin). With a traditional HDD, the OS doesn't need to pay attention to where data is being written or what the relative state of the blocks or pages is. With an SSD, this matter. So, what trim command actually does, when we give modify command, trim will mark some pages which are outdated for erased, and will perform the operation later.

The **TRIM command** allows the operating system to tell the SSD it can skip rewriting certain data the next time it performs a block erase. This lowers the total amount of data the drive writes and increases SSD longevity. Both reads and writes damage NAND flash, but writes do far more damage than reads. Fortunately, block-level longevity has not proven to be an issue in modern NAND flash. More data on [SSD longevity](#), courtesy of the Tech Report, can be found here.

The last two concepts we want to talk about are **wear leveling and write amplification**. Because SSDs write data to pages but erase data in blocks, the amount of data being written to the drive is always larger than the actual update. If you make a change to a 4KB file, for example, the entire block that 4K file sits within must be updated and rewritten. Depending on the number of pages per block and the size of the pages, you might end up writing 4MB worth of data to update a 4KB file. Wear leveling is an algorithm that perform a special task that is to find out the pages where there is no such page that are marked for erase or outdated. So, they find out a fresh block of memory, where they don't need to do any erase, rather than they can write the data easily. When, there will be no such fresh storage, then it will again search the outdated pages or pages which were marked for erase, and update them. So, we can easily understand, an SSD performs well when it has been recently bought and almost fresh. But when it is going to be filled up, the blocks are wearing out day by day, then SSD takes more time than previous to perform necessary operation and gets slower to some comparatively. Here, wear leveling helps them to find out fresh blocks and pages.

Wear- leveling

Wear-leveling data (Wear leveling is a process to distribute write/erase cycles evenly among all flash blocks of the SSD to prevent data from becoming unreliable and a particular block from wearing out) is also stored on the cache for some time during

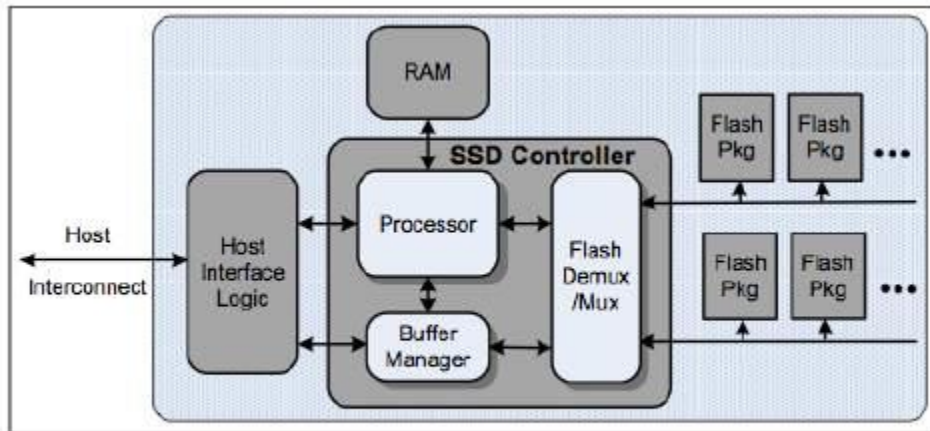
the running of the drive. **Wear leveling** refers to the practice of ensuring certain NAND blocks aren't written and erased more often than others. While wear leveling increases a drive's life expectancy and endurance by writing to the NAND equally, it can actually increase write amplification. In order to distribute writes evenly across the disk, it's sometimes necessary to program and erase blocks even though their contents haven't actually changed. A good wear leveling algorithm seeks to balance these impacts. This is a figure of wear leveling. In SSD, SSD controller operate multiple FLASH particles in parallel through several channels. In first figure it has written data on a single page of 4kb within a block. But when it wants to write data of 16kb, it does not write again pages of the same block, rather it distributes data to multiple pages of multiple blocks. That's why a single block is not wearied out so quickly. This happens because of wear leveling process.

Garbage collection reduces the impact of write amplification, as does the TRIM command. Keeping a significant chunk of the drive free and/or manufacturer over-provisioning can also reduce the impact of write amplification. We know, while writing, to write a 4 KB file, an entire block of 4Mb has to be written or replaced. So, all the blocks of NAND chip get wearied out so quickly. To solve this case, In SSD, there are given some extra storage, which we cannot use physical for storage operation, rather it is used by the OS to reduce the overhead of write amplification. This method is called **over-provisioning**.

The SSD Controller

The controller of SSD is an embedded microchip (such as CPU in computer) with the function of issuing all operation requests ranging from reading and writing data to collecting garbage and depletion equalization algorithm to ensure the speed and cleanliness of SSD. So, the controller is the brain center of SSD. There are major controllers available today such as Marvell, SandForce, Samsung, Indilinx. For example, Marvell is powerful in all directions and the Marvell 88ss9187/89/90

controller is used on the SSD of the Plextone, SanDisk, Crucial and other brands.



SSD Logic Components

A typical SSD controller

NAND flash in an SSD is typically connected to the controller through a series of parallel memory channels, you can think of the drive controller as performing some of the same load-balancing work as a high-end storage array—all operations like wear leveling, garbage collection, and SLC cache management all have parallels in the big iron world.

Some drives also use data compression algorithms to reduce the total number of writes and improve the drive's lifespan. The SSD controller handles error correction, and the algorithms that control for single-bit errors have become increasingly complex as time has passed.

In Some SSD, there is DRAM

A bit of DRAM is included in every SSD for the process of buffering.

Similar to hard drive's cache, data is stored on it for some time temporarily before it is being written to the device.

It increases SSD performance to some level.

Advantages:

- **SSDs are More Durable** Solid-State Drives feature a non-mechanical design of NAND flash mounted on circuit boards, and are shock resistant up to 1500g/0.5ms. Hard Drives consist of various moving parts making them susceptible to shock and damage.
- **SSDs are Faster** SSDs can have 100 times greater performance, almost instantaneous data access, quicker boot ups, faster file transfers, and an overall snappier computing experience than hard drives. HDDs can only access the data faster the closer it is from the read write heads, while all parts of the SSD can be accessed at once. Features
- **SSDs Consume Less Power** SSDs use significantly less power at peak load than hard drives, less than 2W vs. 6W for an HDD. Their energy efficiency can deliver longer battery life in notebooks, less power strain on system, and a cooler computing environment. SSDs are **Lighter** Flash-based SSDs weigh considerably less than hard drives – only 77g vs. 752.5g for HDDs. SSDs won't weigh down your notebook when you're on the go or your desktop when rearranging your office!
- **SSDs are Cooler** As an energy-efficient storage upgrade for your desktop or laptop, SSDs require very little power to operate that translates into significantly less heat output by your system.
- **SSDs are Quieter** With no moving parts, SSDs run at near silent operation and never disturb your computing experience during gaming or movies, unlike loud, whirring hard disc drives.

Disadvantages:

The only single disadvantage of SSDs is that they are quite expensive. But there's a silver lining to it- prices of SSDs are decreasing rapidly. Moreover, with an increase in SSD adoption rates, prices are sure to fall further.