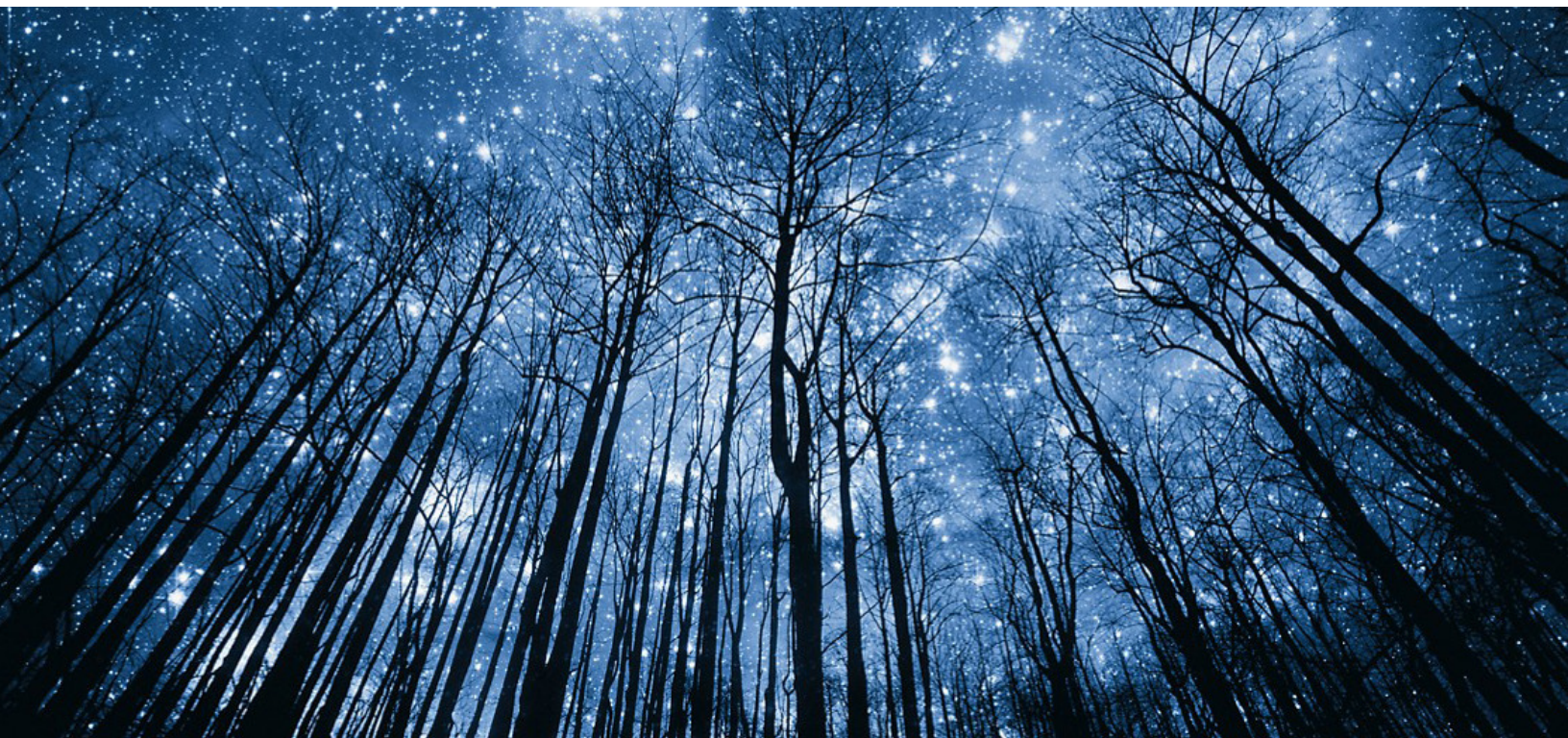


# DATA REDUCTION DEMYSTIFIED



## John Powell

Advisory Systems Engineer, SPS  
Dell Technologies  
[John.s.powell@dell.com](mailto:John.s.powell@dell.com)

## Mark Elliott

Advisory Systems Engineer, SPS  
Dell Technologies  
[Mark.elliott@dell.com](mailto:Mark.elliott@dell.com)

## Ayyaswamy Thangavel

Advisory Systems Engineer, SPS  
Dell Technologies  
[Ayyaswamy.thangavel@dell.com](mailto:Ayyaswamy.thangavel@dell.com)



The Dell Technologies Proven Professional Certification program validates a wide range of skills and competencies across multiple technologies and products.

From Associate, entry-level courses to Expert-level, experience-based exams, all professionals in or looking to begin a career in IT benefit from industry-leading training and certification paths from one of the world's most trusted technology partners.

Proven Professional certifications include:

- Cloud
- Converged/Hyperconverged Infrastructure
- Data Protection
- Data Science
- Networking
- Security
- Servers
- Storage
- Enterprise Architect

Courses are offered to meet different learning styles and schedules, including self-paced On Demand, remote-based Virtual Instructor-Led and in-person Classrooms.

Whether you are an experienced IT professional or just getting started, Dell Technologies Proven Professional certifications are designed to clearly signal proficiency to colleagues and employers.

[Learn more at www.dell.com/certification](http://www.dell.com/certification)

## Table of contents

Introduction.....	4
The History of Data Reduction.....	4
Data compression.....	5
Data deduplication.....	6
The Mechanics Behind Data Reduction.....	6
Compression and Deduplication – Better Together.....	6
Compression and Deduplication – The Math.....	8
Compression and Deduplication – The Write Process.....	9
Rehydration – aka the Read Process.....	10
Sizing for Data Reduction.....	10
The Benefits of Data Reduction.....	11
Reduced data center footprint.....	11
Electricity / Power Costs.....	12
Storage Efficiency.....	12
Data Reduction in Dell Technologies Primary Storage.....	13
PowerStore <sup>5</sup> .....	13
Deduplication.....	13
Data compression.....	14
Thin provisioning.....	14
PowerMax <sup>6</sup> .....	14
Deduplication.....	15
Data compression.....	15
Thin provisioning.....	15
Unity XT <sup>7</sup> .....	16
Deduplication.....	16
Data compression.....	16
Thin provisioning.....	16
PowerFlex <sup>8</sup> .....	17
XtremIO <sup>9</sup> .....	17
SC Series <sup>10</sup> .....	17
Dell Technologies FutureProof Program.....	18
Storage Data Reduction Guarantee <sup>11</sup> .....	18
Remediation.....	19
Conclusion.....	20
Related resources and references.....	23

Disclaimer: The views, processes or methodologies published in this article are those of the authors. They do not necessarily reflect Dell Technologies' views, processes or methodologies.

# Introduction

Increasing amounts of data are created by applications daily. And with no signs of slowing.

According to IDC, data creation leaped forward in 2020 thanks to the COVID-19 pandemic. From 2020 to 2025, IDC forecasts new data creation will grow at a compound annual growth rate (CAGR) of 23%, resulting in approximately 175ZB of data creation by 2025. Ensuring that data is saved efficiently and effectively helps reduce overall solution data cost and resource consumption. Data efficiency not only reduces the amount of data that is stored but also reduces the physical capacity that is required to store the data. Reducing the footprint of the system can also lead to floor space, rack, power, and cooling savings. Virtually all modern storage arrays include some type of data efficiency method to help reduce the total space that is consumed by storage resources created. These methods include compression, deduplication, as well as other inherent efficiencies from snapshots, thin clones, and thin provisioning.

To better educate and enable our pre-sales engineers, partners, and even our customers, this article will explain in an agnostic and product-specific manner why data reduction ratio (DRR) matters, the technology, mechanics, the math, behind it, and most importantly some best practices and examples on how to determine what DRR should be used when sizing a solution. Rather than using a “guesstimate” or “off the shelf” guarantee, knowing what a likely DRR will be should mitigate creating mid-sized configurations which in turn can affect the cost, confidence, customer satisfaction, and the need for post-sales remediation.

## The History of Data Reduction

Data reduction has a long history. Invented in 1838, Morse code is an early example of data compression based on shorter codewords for common letters such as "e" and "t". In 1949 Claude Shannon and Robert Fano devised a systematic way to assign codewords based on probabilities of blocks. An optimal method for doing this was then found by David Huffman in 1951. Early implementations were typically done in hardware, with specific choices of codewords being made as compromises between compression and error correction. The mid-1970s introduced dynamically updating codewords for Huffman encoding, based on the actual data encountered. In the late 1970s, with online storage of text files becoming common, software compression programs began to be developed, almost all based on adaptive Huffman coding. In 1977, Abraham Lempel and Jacob Ziv suggested the basic idea of pointer-based encoding. By the mid-1980s, following work by Terry Welch, the Lempel-Ziv-Welch (LZW) algorithm became the method of choice for most general-purpose compression systems. LZW became popularly used in software such as PKZIP, as well as in hardware devices.

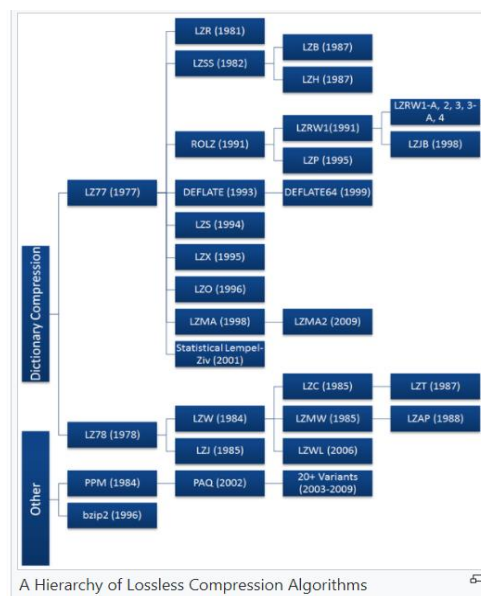


Figure 1: Modern lossless compression algorithms (source: ETHW<sup>1</sup>)

Because disk storage capacity was quite small and very expensive back in the 1970s and 1980s, any algorithm that could save precious bits and bytes was more than welcome. For instance, IBM introduced the first-gigabyte disk drive in 1980. This drive had a 2.5GB capacity, weighed over 100lbs/226Kgs, and cost \$40,000<sup>2</sup>.



Figure 2: 1980s IBM 3380 disk drive (source: Wikipedia<sup>2</sup>)

In the late 1980s, digital images became more common, and because these files were larger than previous text-based files, standards for compressing them like JPEG and GIF began to emerge. Images were quickly followed by digital audio and video and for the same reasons as digital images, compression standards and formats like MPEG were developed.

Despite compression being many decades old and even with the advent of some of these “newer” algorithms, there is still only a small and finite number of algorithms available and in use today<sup>3</sup>.

---

**Takeaway:** Data reduction is by no means a new technology or concept.

---

There are only two ways to reduce the amount of capacity required to store data

1. Reduce the size of the data via compression
2. Remove redundant data blocks via deduplication

Furthermore, and more importantly, there are only so many ways to do compression and there is no “silver bullet” algorithm that offers significant savings over the others. If there was – everyone would be using it.

## Data compression

Software data compression is the process of modifying, encoding, or converting data bits structure to consume less disk space. This approach enables storage size reduction and enables more data storage within the same physical disk space. Data compression is also known as source coding or bit-rate reduction<sup>1</sup>. Whereas deduplication is either yes or no, compression has a spectrum of possibilities. Two factors dictate compression savings:

1. The algorithm is used as discussed earlier in the History section of this article.
2. The type of data being compressed.

While most algorithms used today are general purpose in nature and work well with a multitude of data types, others are “tailored” to a very specific type of data and only produce measurable results for that type of data. For general-purpose algorithms, different types of data will compress differently. Some data types will compress well, others not at all, and others will be in-between. For instance:

- Types of data that compress well include most databases, i.e. Oracle and SQL.

- Types of data that don't compress much if at all include data that is already compressed, i.e., images, audio and video, and anything pre-encrypted.

## Data Deduplication

Data deduplication – often called *intelligent compression* or *single-instance storage* – is a process that eliminates redundant copies of data resulting in reduced storage overhead. Data deduplication techniques ensure that only one unique instance of data is retained<sup>3</sup>.

This technology was first used with backup solutions during the late 1990s and early 2000s. Although similar in concept, in contrast to compression, not only is data deduplication a newer, more recent development, there is truly only one method – storing unique data once by comparing it at the file or block level to a table of what has already been stored. This table is a series of generated HASH IDs. Like compression, the HASH algorithms used will vary.

As previously mentioned, neither compression nor data reduction is new technology. Both methods save space by removing data redundancy. However, it is the combination of the two that yields the best results.

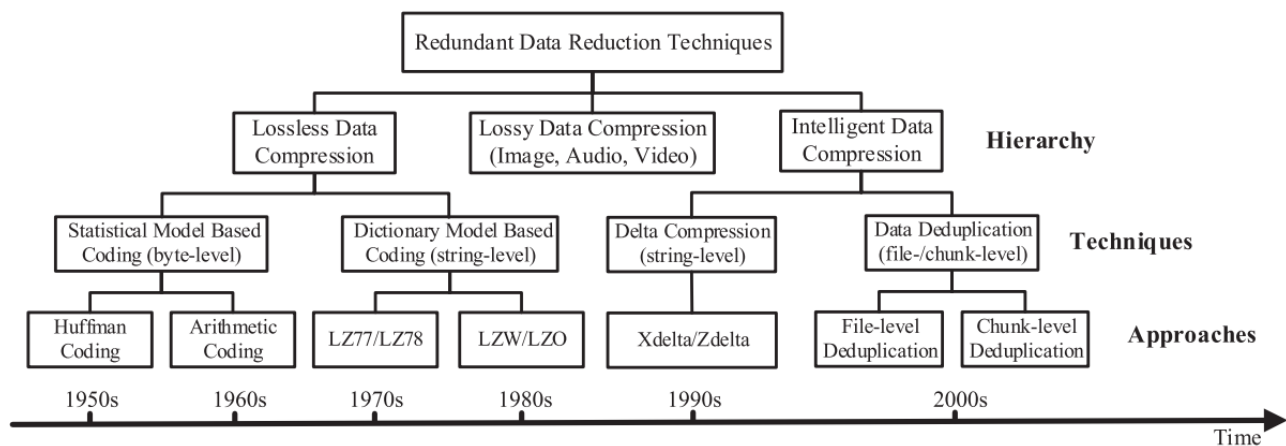


Figure 3: Hierarchy of data reduction techniques (source: IEEE<sup>3</sup>)

## The Mechanics Behind Data Reduction

### Compression and Deduplication – Better Together

As seen in Figure 3, Intelligent Data Reduction combines compression and deduplication and it's this combination that's new to the primary storage market over the past several years. This approach is driven by two factors:

- The introduction of solid-state drive (SSD) technology for persistent data storage. SSDs are more expensive than their spinning counterparts and like in the past, space savings are key to effective storage efficiency.
- Enabling inline, always on data reduction via more powerful processors and offload engines in the storage controller versus non-optimal post-process techniques that were done in the past.

Simply put, modern storage architectures use SSDs as a true persistent storage tier for data and also metadata (data about data), regardless if they are used for caching or not. Moreover, inline data reduction makes SSDs affordable. In many cases, with data reduction services, SSD arrays are more cost-effective than traditional hard disk drive storage. It's this affordability that drives the cost-effective adoption of any technology.

**Question:** How much storage capacity will be saved?

**Answer:** It depends. As mentioned earlier, reducing the size of the data in question and how much of the overall data is the same ultimately dictates how much space will be saved. In simple terms, to compress data, it must be compressible, and to deduplicate data there must be data that can be deduplicated. Consider the following example of creating a new file.

First, if the file is compressible, applying any compression algorithm will reduce the size of the file. With that said, while certain algorithms **MAY** reduce the file more than others, this is not guaranteed. And even if one were better than another – how much better is it? Most modern compression algorithms will produce similar results unless significantly modified to be more “aggressive” (i.e. reducing the size of the file more than an “unmodified”/standard version) – and even then, the extra savings **MAY** be, and typically is, nominal. Plus, the extra is not for free. Most modified algorithms use more CPU cycles. Consequently, is the extra capacity saved worth the extra CPU cycles? More times than not, the answer is **NO** – especially when CPU cycles are finite and are allocated to critical tasks such as processing IO and the advanced data services most modern storage arrays need to provide.

Second, how many copies of the data being stored are duplicated? The more "hits" you have, the more capacity will be saved via deduplication. Beyond user file systems and shares where some duplicate files could be identified, virtualized servers and virtual desktop infrastructure (VDI), or multiple copies/instances of database files are great candidates for deduplication because many of the underlying data blocks are most likely the same.

One thing that will affect reduction savings, specifically at the sub-file/chunk level, is the size of the chunk. The smaller the chunk, the better the savings **CAN** be. However, like the modified compression algorithms discussed above, the more chunks you have, the more CPU cycles it will require to analyze them. As with pretty much everything, there is no "free lunch". To provide effective and efficient data reduction, a balance must always be sought between the pros and cons.

---

**Takeaway:** The following bullets may sound like common sense but as you can see, everything is interrelated.

---

- The more the data is the same/repeated/redundant, the more capacity you **WILL** save
- The longer you keep the data, the more capacity you **WILL** save
- The more static the data, the more capacity you **WILL** save
- The more granular the deduplication, the more capacity you **WILL** save

## Compression and Deduplication – The Math

Although there is no formula to define the savings for compression and deduplication respectively, there is one that shows how both combined to form what is known as the data reduction ratio (DRR) using the following formula:

$$\text{DRR} = \text{Compression Savings} * \text{Deduplication Savings}$$

DRR is typically expressed in a ratio format (e.g. 4:1) and has changed the traditional definitions of capacity – specifically adding the concept of Effective Capacity (TBe) to the typical Raw (TBr) to Usable (TBu) capacity categorizations.

**Raw Capacity (TBr)** = the number of disks ^ their respective size

**Usable Capacity (TBu)** = Raw Capacity \* Efficiency

where **Efficiency** is the percentage of Raw Capacity that will be “converted” to Usable Capacity. This percentage can range anywhere from roughly 30% to 90% depending on the geometry/RAID stripe width

**Effective Capacity (TBe)** = Usable Capacity \* DRR

---

**Takeaway:** DRR certainly dictates the effective capacity of an array. The higher the DRR, the better the capacity savings. And the greater the savings, the less raw capacity that is required. However, do not discount **Efficiency**. Efficiency can “level the playing field” when comparing solutions with different DRRs. (i.e., a higher efficiency can make up for a lower DRR – and in some instances, even more so).

---

**Question:** How do you then quantify DRR?

**Answer:** When it comes to quantifying DRR, “mileage will vary”. It starts with understanding the data to be stored. First, there are “suitable” and “unsuitable” data types – in other words, is the data to be stored amenable to reduction?

Examples of suitable types include databases (i.e. SQL, Oracle), VMs, and virtual desktops.

Examples of unsuitable types include PDFs, audio, and image (picture and video) files. The reason? These types of data are already reduced (compressed) as part of their file format. Any data type that is pre-compressed or pre-encrypted at the host, or “in-flight” will not further reduce.

Beyond identifying the types of data, two other things need to be considered:

- How “mixed” are the data types? For instance, if most of your data is suitable, you will realize a higher DRR than if your data is mixed or even worse, if it’s mostly unsuitable
- The uniqueness of the data being stored. The more unique the data, the less deduplication that will occur. On the other hand, things like VDI or even VSI tend to dedupe extremely well depending on how they are implemented (i.e. not using linked clones or other hosts/virtualization-based data reduction techniques)

Lastly, compression is typically more significant to DRR than deduplication, especially in database-heavy environments. Keep in mind, as stated earlier, your mileage will vary.

---

**Takeaway:** Data reduction savings with primary storage is near impossible to quantify and no formula can be used to calculate the savings from compression and deduplication, respectively. To “effectively estimate” savings from data reduction, the best one can do is attempt to identify the types of data to be stored, how much of each type, and how redundant it is. Also, the more data that is stored and the longer it is stored, the bigger the savings that **CAN** be realized. It typically takes hundreds of TBs of capacity and a mix of workloads and file/data types to generate any measurable and consistent amount of data reduction savings. Therefore, **DO NOT** expect to see immediate savings from data reduction – it may take a while.

---



## Compression and Deduplication – The Write Process

In general, the write data path for compression and duplication is as follows:

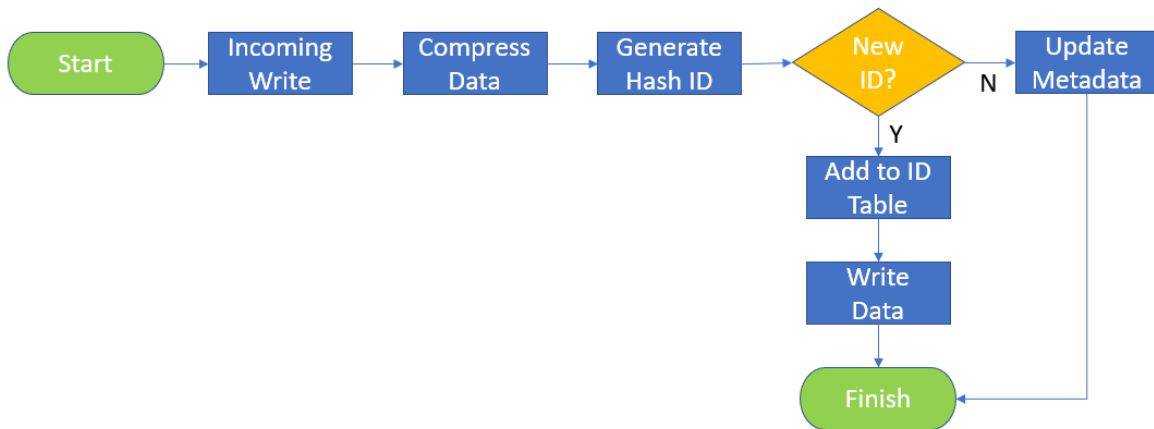


Figure 4: A Typical Data Reduction Process

---

**Takeaway:** Keeping Figure 4 in mind, the following **WILL** ultimately dictate how data reduction is performed, the DRR that will be realized, and what implications may arise when comparing one architecture with another.

---

- Incoming writes **MAY** be broken up into smaller “chunks” and if so, chunk sizes will vary.
  1. The smaller the chunk, the more CPU that is required with potentially higher DRR
  2. The larger the chunk, the less CPU that is required with potentially lower DRR
- Compression and/or dedupe **MAY** be inline or post-process. They both will use the same amount of CPU cycles. The difference being when those cycles are used.
  1. Inline is real-time as the data is coming in which means it **MAY** be harder to “keep up” with IO and data services during the production day.
- Post-process is done later so it’s not as hard to keep up during the production day. However, it requires more capacity as the data must have a “landing space”.
- Compression and/or Deduplication **MAY** be offloaded from the controller CPU. Offloading **WILL** solve the “keeping up” problem as well as enable selection of a smaller chunk size.
- Compression and Hash algorithms **WILL** vary and ultimately dictate what DRR is realized – but the difference in many instances is nominal.
- Compressed and deduped data may be written directly to disk or a page/stripes which will occur once it is full. How writes eventually make it to the disk is key to SSD longevity (drive wear) as well as how garbage collection is performed and what the impact will be, if any.

## Rehydration – aka the Read Process

Rehydration is “reversing” the effect of compression and duplication. For instance, when data that has been written to a disk is requested (e.g. from a host read, or for array-based replication), it must be returned in its original form before compression or deduplication. That requires metadata operations to find where the unique data is physically stored and subsequently uncompress it.

---

**Takeaway:** Although not talked about much if at all, as with compression and deduplication, rehydration will require CPU cycles and accessing metadata. Furthermore, the data rehydration process requires time to perform the required steps. If these operations cannot be performed in a timely fashion, latency **WILL** increase and ultimately **MAY** impact end-user/application performance.

---

## Sizing for Data Reduction

DRR is nearly impossible to quantify in most situations as circumstances will vary from situation to situation and data type to data type. It is **HIGHLY** recommended to be diligent to understand what is being stored vs. simply going with an "off the shelf" DRR (aka the "easy" route) or 'guesstimating'.

---

**Takeaway:** Understanding what a likely DRR will be is key to correctly sizing a solution. If you choose the “easy” route, more than likely one of two things will occur:

---

- If your DRR estimate is too high, you will undersize the solution which means you won't have enough capacity to store the required amount of data. This can lead to significant customer dissatisfaction issues and remediation procedures which will cost someone extra money that was not originally planned for.
- If your DRR estimate is too low, you will oversize the solution which means you specified too much capacity and could potentially price yourself out of an opportunity.

Another thing to remember; because the CPU's impact on data reduction contributes to both reads and writes as discussed earlier, it is essential to quantify the likely R/W ratio along with the sustained IOP profile. If you do not, mis-sizing the capacity as discussed above, you may select a model of the controller that either is too "big" or not "big" enough.

## Benefits of Data Reduction

The primary benefit of data reduction is financial. Regardless of what the DRR is, any reduction means procuring less physical capacity. Fewer disks can lead to:

- Lower procurement costs
- Less floor and/or rack space (based on needing fewer expansion enclosures)
- Reduced power and cooling overhead

Moreover, the resulting savings of DRR will directly affect IT operational efficiency and effectiveness.

- **Efficiency:** Refers to using the most appropriate amount of resources. And at minimal one-time (CapEx) and recurring (OpEx) costs.
- **Effectiveness:** Refers to whether IT can continually deliver services that meet objectives and agreed Service Level Agreements (SLAs).

## Reduced data center footprint

For instance, a single commercial data center rack including dual power costs between \$1,300 and \$1,600 per month. Or \$78,000 or \$96,000 over 5yrs. If data center requirements can be reduced, data center cost savings between 65% and 75% can be achieved.

In the following example, DRR driven efficiency and effectiveness can drive:

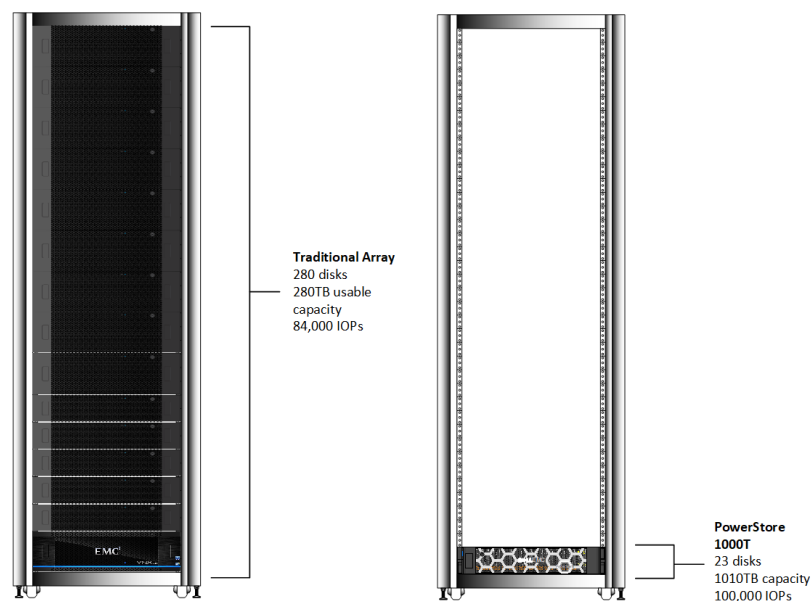


Figure 5: The DRR effect on data center consolidation

- 95% less required rack space
- 92% fewer physical disks
- 260% capacity increase based on a 4:1 DRR
- Valuable rack space freed up and now available for other solutions, or scale storage capacity
- Ability to consolidate servers and backup solutions in a single rack

## Electricity / Power Costs

A typical hard disk based, full-rack storage system with over 200 disks consumes up to 4KW of power per hour 24 hours a day. At \$0.16 per KW/h that's:

$$(\$0.16 \times 4\text{KW}) \times 24 = \$15.36 \text{ per day, or } \$461 \text{ per month}$$

However, a modern SSD-based array benefiting from fewer drives and data reduction will typically use up to 1.5KW per hour.

$$(\$0.16 \times 1.5\text{KW}) \times 24 = \$5.76 \text{ per day, or } \$172 \text{ per month}$$

This equates to saving over \$17,340 over a 5 -year period.

---

**Takeaway:** When combining the savings from the physical footprint, power, and cooling, an organization **CAN** realize savings exceeding \$100,000 over 5 years in addition to a reduction to their carbon footprint.

---

## Storage Efficiency

In addition to the savings from data reduction (compression and deduplication), additional efficiency space savings via snapshots, thin clones, and thin provisioning can be realized in most modern storage arrays. The following represents how PowerStore calculates and reports Data Reduction and system efficiency items such as snapshot/thin clone and thin provisioning savings.

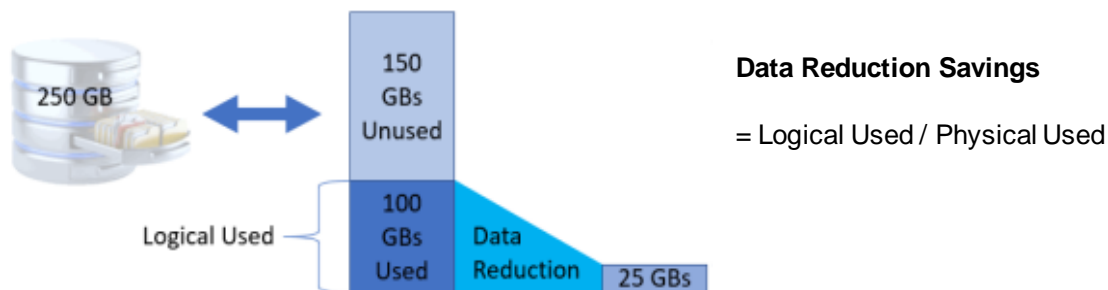


Figure 6: Data Reduction Savings

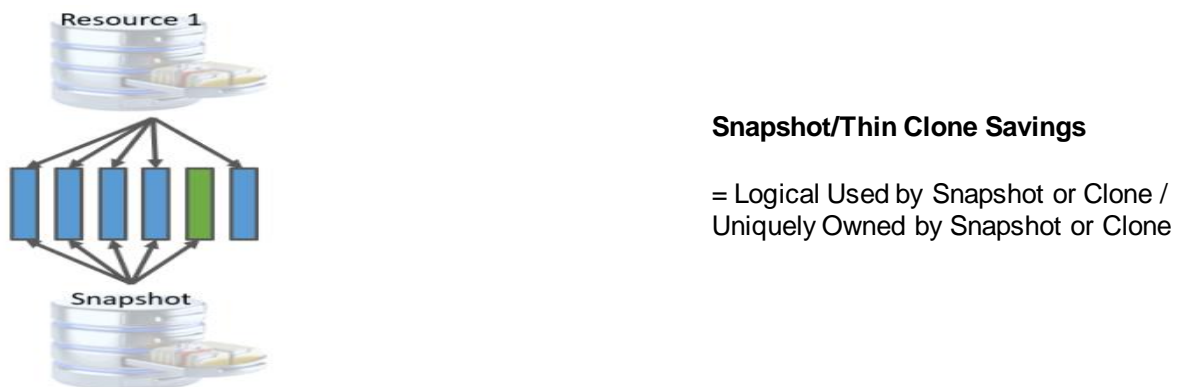
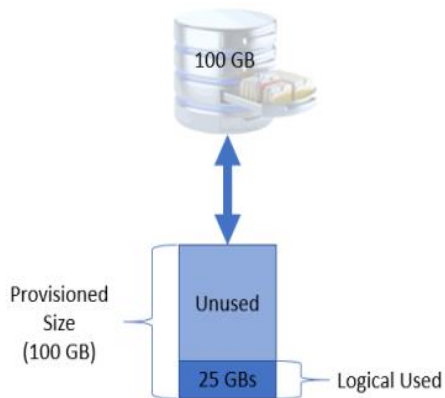


Figure 7: Snapshot Savings



### Thin Provisioning Savings

$$= \text{Provisioned Size} / \text{Logical Used}$$

Figure 8: Thin Provisioning Savings

---

**Takeaway:** There is a difference between Data Reduction savings from compression and deduplication and the Efficiency savings from snapshots, thin clones, and thin provisioning.

---

- Data Reduction should **ONLY** include savings from compression and deduplication.
- Efficiency should include the inherent savings from snapshots, thin clones, and thin provisioning.
- If overall System Efficiency is reported, it probably includes both types of savings.
- Each solution **MAY** calculate and report on each of the above categorizations differently. Consequently, you cannot compare them at face value as they may not be similar. Be sure you know what is reported and how it is calculated.

## Data Reduction in Dell Technologies Primary Storage

### PowerStore<sup>5</sup>

PowerStore data reduction services optimize storage capacity by reducing storage blocks to their simplest form. This approach provides increased effective storage available to hosts, applications, and users. Data reduction leverages many techniques delivered using intelligent software. The basic approach is simple – increase the volumetric efficiency of physical storage as much as possible resulting in maximum usable capacity.

### Deduplication

As data enters PowerStore, each node processes data reduction services inline. Data entering a node is compared and deduplicated against the data received on the peer node. Each node contains a data fingerprint cache. When data enters a node, and a new/unique fingerprint is generated owned by that node. Deduplication is achieved as fingerprints are compared across each node using the internal mid-plane links. The result of this approach is increased data reduction when common data blocks are identified. Not mirroring the fingerprint cache across nodes, but instead having the ability to compare fingerprint cache between nodes enables a larger cache size. Therefore, an increased amount of deduplication is achieved by globally tracking a larger number of fingerprints supported by a larger cache.

## Data compression

Deduplication, then compression, occurs within an appliance when data is written from the NVRAM drives to the data drives. During this process, data is stored within the back-end storage in full stripe writes, which are created once data has passed through deduplication and compression processes.

PowerStore uses the DEFLATE algorithm (LZ77 + Huffman). As compression is CPU-intensive, this process is offloaded to a dedicated Lewisburg chip with Intel® QuickAssist Technology (QAT) that improves performance across lossless data compression.

Offloading compression to a dedicated chip ensures CPUs dedicated to data services are not compromised.

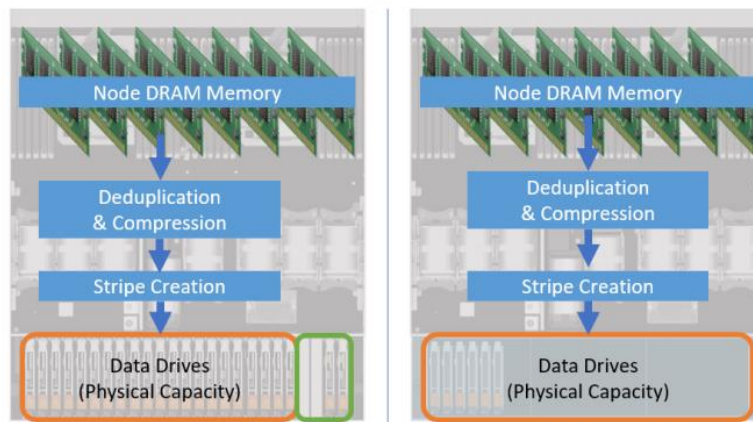


Figure 9: PowerStore I/O flow (1000 – 9000 and 500T)

## Thin provisioning

Thin Provisioning (TP) also improves efficiency. It optimizes space efficiency by pre-allocating disk storage in a flexible manner. When allocating a thin-provisioned volume, PowerStore does not reserve disk space in advance. Instead, storage is allocated dynamically on demand. Free space is released back to the storage system when data in the volume is deleted.

For example, 20 host servers need 1000GB of disk space each. We estimate 20TB disk space will eventually be required. In this situation, we could allocate 20TB of thin-provisioned usable storage. Logically, 20TB has been allocated to host servers. But as not all space is used, thin provisioning allows us to satisfy the needs of the larger storage consumers without having to purchase storage that might never be used by all host servers. Since storage space is not allocated until it is consumed, we can “overcommit” storage and deliver to host servers what’s actually required.

For further details on PowerStore’s data reduction services, please see the Dell EMC PowerStore: Data Efficiencies white paper: <https://www.delltechnologies.com/asset/en-us/products/storage/industry-market/h18151-dell-emc-powerstore-data-efficiencies.pdf>

## PowerMax<sup>6</sup>

PowerMax data reduction combines the Adaptive Compression Engine (ACE) and inline deduplication to provide a space-efficient platform with negligible impact on performance. Executing these functions in parallel allows PowerMax to be extremely capacity-efficient which in turn produces a significantly smaller data center footprint and an overall reduction in TCO. Using PowerMax data reduction is as simple as a single click to either enable or disable.

## Deduplication

In PowerMax systems, dedupe is accomplished through a series of functions and components including hardware acceleration, dedupe algorithm, hash table, and dedupe management object (DMO). Dedupe is an inline process that uses the same data reduction hardware as compression. All data reduction-enabled incoming data is passed through data reduction hardware. In a single pass, the data reduction hardware handles compression, pattern detection and generates a hash ID for deduplication. This produces compressed data with a unique hash ID. Leveraging data reduction hardware for this process allows system resources to be focused on host I/O and other system operations.

## Data compression

PowerMax's Adaptive Compression Engine (ACE) is the combination of multiple core components that work together to achieve maximum system efficiency while delivering optimized performance. These core components include:

- Hardware acceleration
- Optimized data placement
- Activity-Based Compression
- Fine-grain data packing
- Extended Data Compression

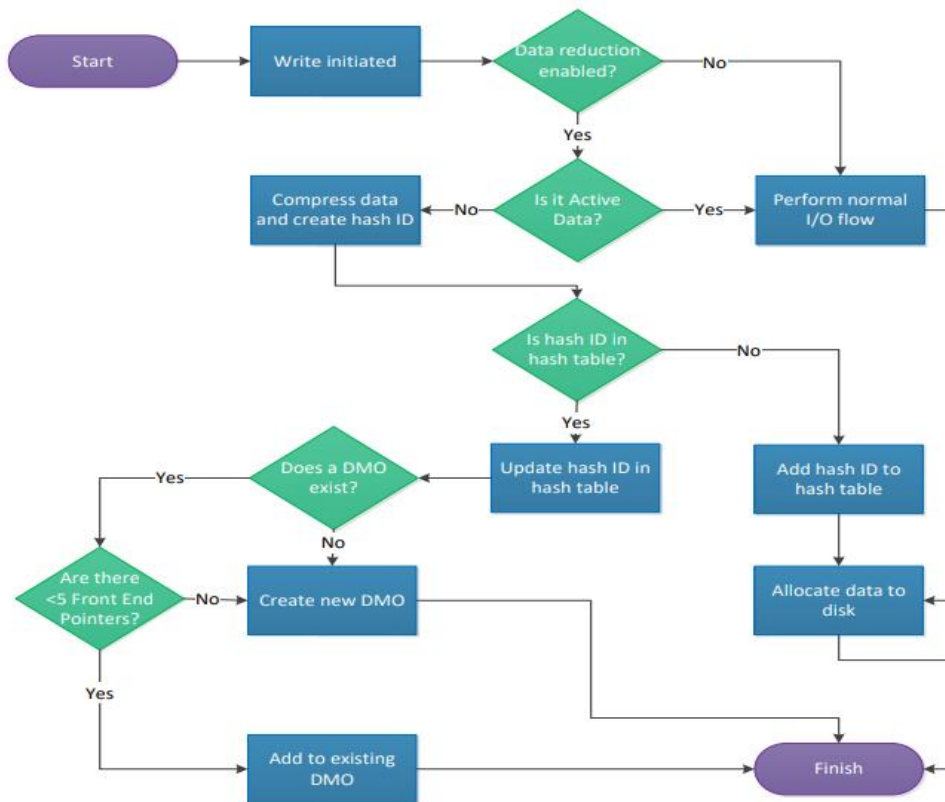


Figure 10: PowerMax Data Reduction IO Flow

## Thin provisioning

Inherent Thin Provisioning also improves overall system efficiency by pre-allocating disk storage in a flexible manner. When allocating a thin-provisioned volume, PowerMax does not reserve disk space in advance. Instead, storage is allocated dynamically on demand. Free space is released back to the storage system when data in the volume is deleted.

For further details on PowerMax's data reduction services, please see the Dell EMC PowerMax: Data Reduction white paper: <https://www.delltechnologies.com/asset/el-gr/products/storage/industry-market/h17072-data-reduction-with-dell-emc-powermax.pdf>

## Unity XT<sup>7</sup>

For systems running OE 4.3 or later, Unity XT inline data reduction provides space savings using data deduplication and compression. Data reduction is easy to manage, and once enabled, is intelligently controlled by the storage system. Configuring data reduction and reporting savings is simple, and can be done through Unisphere, Unisphere CLI, or REST API.

### Deduplication

When new data first enters the data reduction logic, it is first passed through the deduplication algorithm. The deduplication algorithm is a lightweight software algorithm that analyzes the blocks of data for known patterns. The patterns may be a block of zeros written by the host, or common patterns found in Dell EMC Unity's many use cases, such as virtual environments. In the event no pattern is detected, the data is passed to Advanced Deduplication if it is enabled. The Advanced Deduplication algorithm utilizes fingerprints created for each block of data to quickly identify duplicate data within the dataset.

### Data compression

As blocks enter the compression algorithm if savings can be achieved, space is allocated within the Pool which matches the compressed size of the data, the data is compressed, and the data is written to the Pool. When Advanced Deduplication is enabled, the fingerprint for the block of data is also stored with the compressed data on disk. Compression will not compress data if no savings can be achieved.

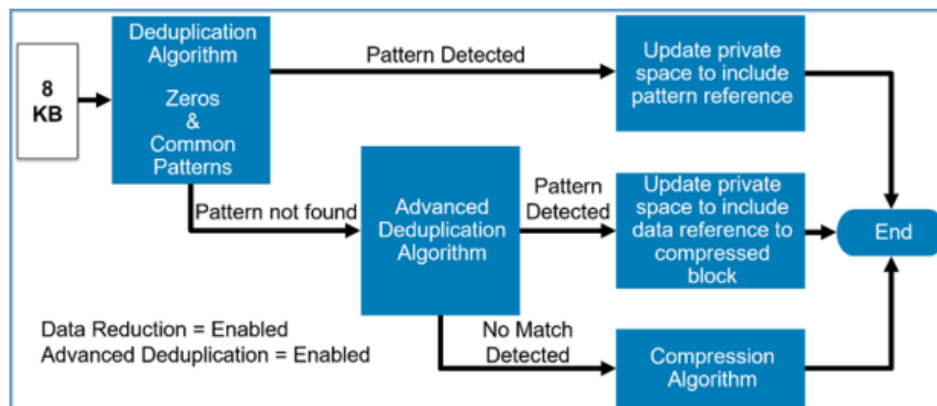


Figure 11: Dell EMC Unity Data Reduction – Advanced Deduplication Enabled

### Thin provisioning

Inherent Thin Provisioning also improves overall system efficiency by pre-allocating disk storage in a flexible manner. When allocating a thin-provisioned volume, Unity does not reserve disk space in advance. Instead, storage is allocated dynamically on demand. Free space is released back to the storage system when data in the volume is deleted.

For further details on Unity XT's data reduction services, please see the Dell EMC Unity: Data Reduction white paper: <https://www.delltechnologies.com/asset/en-us/products/storage/industry-market/h16870-dell-emc-unity-data-reduction.pdf>



## **PowerFlex<sup>8</sup>**

Although PowerFlex is part of the Dell Technologies HCI portfolio, it does have the capability to “mimic” a more traditional tiered architecture as well. It is also a part of the FutureProof Storage Data Reduction Guarantee.

PowerFlex creates a server and IP-based SAN from direct-attached disk storage to deliver flexible and scalable architecture. As an alternative to a traditional SAN infrastructure, it combines diverse storage media to create virtual pools of block storage with varying performance and data services options. PowerFlex exclusively uses inline data compression and thin provisioning to effectively deliver storage savings and efficient capacity management.

For further details on PowerFlex’s inline compression, please see the PowerFlex administrator guide or best practices white paper at: <https://www.delltechnologies.com/asset/en-us/products/storage/industry-market/h18390-dell-emc-powerflex-networking-best-practices-wp.pdf>

## **XtremIO<sup>9</sup>**

The XtremIO X2 all-flash array has a wide range of features including inline data reduction. Data reduction services work by first deduplicating data: If two blocks are the same, deduplication uses metadata to track the duplicate block. Metadata is maintained in memory for fast access, and unique blocks are written to storage, thus saving space on the flash drives. After deduplication, the XtremIO X2 array compresses data, reducing the amount of array space that is used by unique blocks of data. The compression space savings means that data blocks are stored most efficiently. Inline data reduction allows for consolidation of more workloads without cluster expansion.

For further details on XtremIO’s data reduction services, please see the Introduction to XtremIO X2 Storage Array white paper: <https://www.delltechnologies.com/asset/en-us/products/storage/industry-market/h16444-introduction-xtremio-x2-storage-array-wp.pdf>

## **SC Series<sup>10</sup>**

SCOS 7 or higher offers block-level deduplication and compression on the SC Series. Dell Storage Center's implementation of deduplication uses flash in the array for metadata. Integrating deduplication and compression into Data Progression (SC's auto-tiering engine) is a natural fit and provides the effective benefits of inline deduplication while optimizing performance and data protection.

For further details on SC’s data reduction services, please see the SC Series administrator guide or the Compression and Deduplication Feature Brief at: [https://i.dell.com/sites/csdocuments/Shared-Content\\_data-Sheets\\_Documents/en/FB\\_SC\\_Series\\_Dedup\\_Compression.pdf](https://i.dell.com/sites/csdocuments/Shared-Content_data-Sheets_Documents/en/FB_SC_Series_Dedup_Compression.pdf)

# Dell Technologies Future-Proof Program

## Storage Data Reduction Guarantee<sup>11</sup>

The Future-Proof Program is designed to help customers optimize the IT lifecycle through a series of guarantees, offers, and assurances. Future-Proof provides support from beginning to end by guaranteeing outcomes, maximizing investments, and helping customers navigate the future of IT. This program enables customers to focus on critical business needs while Dell Technologies handles the rest.

The program consists of several offers; however, we are going to specifically focus on the Storage Data Reduction Guarantee. For information on the others, please see <https://www.dell.com/en-us/dt/products/future-proof-program.htm>

The Storage Data Reduction Guarantee provides guaranteed storage data reduction for your workloads without an assessment.

---

**Takeaway:** To be eligible, the customer must sign the Terms and Conditions document before an order is placed. <https://www.dell.com/en-us/dt/products/future-proof-program.htm#pdf-overlay=//www.delltechnologies.com/asset/en-us/products/storage/legal-pricing/storage-data-reduction-guarantee-terms-and-conditions.pdf>

---

Table 1 outlines the DRR guarantee as outlined in the Terms and Conditions document.

Array	DRR Guarantee
PowerStore	4:1
PowerMax	4:1
Unity XT x80F models	3:1
XtremIO	3:1
SC Series <sup>*</sup>	2:1
PowerFlex	2:1

**Table 1: FutureProof Storage Data Reduction Guarantee Ratios**

\* SC9000, SC7020F, SC7020, SC5020F, SC5020, SC4020 and SCv3000 Arrays (All Flash configurations) Although not recommended nor necessary in most situations, requests for custom guarantees vs. the standard outlined above can be made by emailing [storage.data.reduction@dell.com](mailto:storage.data.reduction@dell.com). Information surrounding what the specific request is (e.g. a specific DRR) and why it is being requested (i.e. a justification statement) should be included in the request.

## Remediation

In the event your array is not delivering the guaranteed DRR per the Terms and Conditions document, the remediation process below should be followed.

Submit the following information to [storage.data.reduction@dell.com](mailto:storage.data.reduction@dell.com):

1. A copy of the signed guarantee Terms and Conditions document
2. The serial number of the array
3. A screenshot of the respective array's manager GUI showing the array is > **50%** full
4. For virtual environments, a screenshot to show that **UNMAP** and in-guest reclaims are **ON**
5. Information for any volumes that have > **30%** unique data
  - a. Based on the unique/used percentage found in the volume details section of the array's manager GUI or CloudIQ
  - b. Is there any video, audio, PDFs, or anything that is compressed, or encrypted?
  - c. How much of each

Figure 12 represents an example of how you may want to format the data for item 5 outlined above.

Volume	Used	Excluded	%	ZIP	GZ	JPG	PDF	MP4	MOV
AppVol1	14950.4	7680	51.37%	2048		1126.4	4505.6		
AppVol2	11878.4	5647.8	47.55%	1740.8		1017	2560		330
AppVol3	8908.8	4634.2	52.02%	1228.8		620	2150.4	635	
AppVol4	2764.8	523.4	18.93%	177.8		34.5	414.7	74.2	
AppVol5	10956.8	6694.1	61.10%	1228.8		720.6	3379.2	937	428.5
AppVol6	81.2	0	0.00%						
AppVol7	5427.2	3308.8	60.97%	1002		331	1843.2	132.6	
AppVol8	596	258.3	43.34%	58.5		12	187.8	0	
AppVol9	4812.8	2587.9	53.77%	321		500	865	489.3	412.6
AppVol10	2457.6	1488.9	60.58%	446		377.5	519.8	145.6	
AppVol11	1945.6	1211.5	62.27%	298.8		89.7	613	136.5	73.5
	<b>64779.6</b>	<b>34034.9</b>	<b>52.54%</b>						

Figure 12: Example Format for Remediation Submission

The above can be gathered by the customer and/or the local SE. Once received, the information will be reviewed. Responses are typically sent within 24 hours.

- If approved, an email with the process to follow will be sent to the requestor.
- If denied, an email will be sent to the requestor with an explanation.

## Conclusion

Throughout this paper, several Takeaways have been noted in the hope that collectively they will not only “debunk” many of the myths about data reduction (compression and deduplication) – specifically one solution or method being substantially better than another – but they will also serve as best practice guidance. The following is a recap of the Takeaways.

1. Data reduction is by no means a new technology or concept, and when you get down to it, there are only two ways to reduce the amount of capacity required to store data
  - Reduce the size of the data via compression
  - Remove redundant data blocks via deduplication

More importantly, there are only so many ways to do compression and there is no “silver bullet” algorithm that offers significant savings over the others. If there was – everyone would be using it.
2. The following few bullets may sound like common sense but as you can see, everything is interrelated.
  - The more the data is the same/repeated/redundant – the more capacity you **WILL** save
  - The longer you keep the data - the more capacity you **WILL** save
  - The more static the data – the more capacity you **WILL** save
  - The more granular the deduplication – the more capacity you **WILL** save
3. DRR certainly dictates the effective capacity of an array. The higher the DRR, the better the capacity savings and, the greater the savings, the less raw capacity that is required. However, do not discount **Efficiency**. Efficiency can “level the playing field” when comparing solutions with different DRRs. (i.e. a higher efficiency can make up for a lower DRR – and in some instances, even more so).
4. Data reduction savings with primary storage is near impossible to quantify and no formula can be used to calculate the savings from compression and deduplication respectively. To “effectively estimate” the savings from data reduction, the best one can do is attempt to identify the types of data to be stored, how much of each type, and how redundant is it. Also, the more data that is stored and the longer it is stored, the bigger the savings that **CAN** be realized. It typically takes hundreds of TBs of capacity and a mix of workloads and file/data types to generate any measurable and consistent amount of data reduction savings. Therefore, **DO NOT** expect to see immediate savings from data reduction – it may take a while.
5. The following **WILL** ultimately dictate how data reduction is performed, the DRR that will be realized, and what implications may arise when comparing one architecture over another.
  - Incoming writes **MAY** be broken up into smaller “chunks” and if so, chunk sizes will vary.
    1. The smaller the chunk, the more CPU that is required with potentially higher DRR.
    2. The larger the chunk, the less CPU that is required with potentially lower DRR.
  - Compression and/or dedupe **MAY** be inline or post-process. They both will use the same amount of CPU cycles but the difference is when are those cycles required?
    1. Inline is real-time as the data is coming in which means it **MAY** be harder to “keep up” with IO and data services during the production day,

2. Post-process is done later so it's not as hard to keep up during the production day. However, it requires more capacity as the data must have a "landing space".
  - Compression and/or Deduplication **MAY** be offloaded from the controller CPU. Offloading **WILL** solve the "keeping up" problem as well as enable the selection of a smaller chunk size
  - Compression and Hash algorithms **WILL** vary and ultimately dictate what DRR is realized – but the difference in many instances is nominal
  - Compressed and deduped data may be written directly to disk or a page/stripe which will occur once it is full. How writes eventually make it to the disk is key to SSD longevity as well as how garbage collection is performed and what the impact will be, if any.
6. Although not talked about much if at all, as with compression and deduplication, rehydration will require CPU cycles and accessing metadata. Furthermore, the data rehydration process requires time to perform all the required steps. If these operations cannot be performed in a timely fashion, latency **WILL** increase and ultimately **MAY** impact end-user/application performance.
7. Understanding what a likely DRR will be is key to correctly sizing a solution. If you choose the "easy" route, more than likely one of two things will occur:
  - If your DRR estimate is too high, you will undersize the solution which means you won't have enough capacity to store the required amount of data. This can lead to significant customer satisfaction issues and remediation procedures which will cost someone extra money that was not originally planned for.
  - If your DRR estimate is too low, you will oversize the solution which means you specified too much capacity and could potentially price yourself out of an opportunity.

The other thing to remember, because the CPU impact data reduction contributes to both reads and writes, it is essential to quantify the likely R/W ratio along with the sustained IOP profile. If you do not, like mis-sizing the capacity as discussed above, you may select a model of the controller that either is too "big" or not "big" enough.

8. There is a difference between Data Reduction savings from compression and deduplication and the Efficiency savings from snapshots, thin clones, and thin provisioning.
  - Data Reduction should **ONLY** include savings from compression and deduplication.
  - Efficiency should include the inherent savings from snapshots, thin clones, and thin provisioning.
  - If overall System Efficiency is reported, it probably includes both types of savings.
  - Each solution **MAY** calculate and report on each of the above categorizations differently. Consequently, you cannot compare them at face value as they may not be similar. Be sure you know what is reported and how it is calculated.
9. The primary benefit of data reduction is a financial one. For a typical mid-sized array, savings from the physical footprint, power, and cooling **CAN** exceed \$100,000 over 5 years in addition to a significant reduction in carbon footprint.
10. The customer must sign the Terms and Conditions document before an order is placed.

<https://www.dell.com/en-us/dt/products/future-proof-program.htm#pdf->

[overlay=//www.delltechnologies.com/asset/en-us/products/storage/legal-pricing/storage-data-reduction-guarantee-terms-and-conditions.pdf](https://www.delltechnologies.com/asset/en-us/products/storage/legal-pricing/storage-data-reduction-guarantee-terms-and-conditions.pdf)

---

**Final Takeaway:** As with pretty much everything, there is no “free lunch”. To provide effective and efficient data reduction, a balance must always be struck between the pros and cons as it pertains to the data reduction algorithms chosen, the CPU cycles required to perform those algorithms, and the resulting DRR.

---

## Related resources and references

ETHW (2019) 'History of Lossless Data Compression Algorithms', [Online].  
[https://ethw.org/History\\_of\\_Lossless\\_Data\\_Compression\\_Algorithms#History](https://ethw.org/History_of_Lossless_Data_Compression_Algorithms#History)

Think Computers (2013) 'The History of the Hard Drive', [Online].  
<https://thinkcomputers.org/the-history-of-the-hard-drive/>

Bigelow, S (2019), TechTarget 'Data Deduplication', [Online].  
<https://www.techtarget.com/searchstorage/definition/data-deduplication>

Xia, W et al (2016), IEEE 'A Comprehensive Study of the Past, Present, and Future of Data Deduplication', [Online].  
[https://www.researchgate.net/publication/305801856\\_A\\_Comprehensive\\_Study\\_of\\_the\\_Past\\_Present\\_and\\_Future\\_of\\_Data\\_Deduplication](https://www.researchgate.net/publication/305801856_A_Comprehensive_Study_of_the_Past_Present_and_Future_of_Data_Deduplication)

Poulin, Ryan (2021), 'Dell EMC PowerStore: Data Efficiencies', [Online].  
<https://www.delltechnologies.com/asset/en-us/products/storage/industry-market/h18151-dell-emc-powerstore-data-efficiencies.pdf>

Tasker, Robert (2022), 'Dell EMC PowerMax: Data Reduction white paper', [Online].  
<https://www.delltechnologies.com/asset/el-gr/products/storage/industry-market/h17072-data-reduction-with-dell-emc-powermax.pdf>

Poulin, Ryan, (2021), 'Dell EMC Unity: Data Reduction', [Online].  
[https://www.delltechnologies.com/asset/en-us/products/storage/industry-market/h16870-dell\\_emc\\_unity-data\\_reduction.pdf](https://www.delltechnologies.com/asset/en-us/products/storage/industry-market/h16870-dell_emc_unity-data_reduction.pdf)

Dean, Brian, (2021), 'Dell EMC PowerFlex: Networking Best Practices and Design Considerations', [Online].  
<https://www.delltechnologies.com/asset/en-us/products/storage/industry-market/h18390-dell-emc-powerflex-networking-best-practices-wp.pdf>

Dell EMC Engineering, (2019), 'Introduction to XtremIO X2 Storage Array', [Online].  
<https://www.delltechnologies.com/asset/en-us/products/storage/industry-market/h16444-introduction-xtremio-x2-storage-array-wp.pdf>

Dell EMC Engineering, (2016), 'SC Series Compression and Deduplication Feature Brief', [Online].  
[https://i.dell.com/sites/csdocuments/Shared-Content\\_data-Sheets\\_Documents/en/FB\\_SC\\_Series\\_Dedup\\_Compression.pdf](https://i.dell.com/sites/csdocuments/Shared-Content_data-Sheets_Documents/en/FB_SC_Series_Dedup_Compression.pdf)

Dell Technologies, (2022), Future Proof Program Portal, [Online].  
<https://www.dell.com/en-us/dt/products/future-proof-program.htm>

(All references and resources accessed 24 March 2022).

---

Dell Technologies believes the information in this publication is accurate as of its publication date. The information is subject to change without notice.

THE INFORMATION IN THIS PUBLICATION IS PROVIDED "AS IS." DELL TECHNOLOGIES MAKES NO REPRESENTATIONS OR WARRANTIES OF ANY KIND WITH RESPECT TO THE INFORMATION IN THIS PUBLICATION, AND SPECIFICALLY DISCLAIMS IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE.

Use, copying and distribution of any Dell Technologies software described in this publication requires an applicable software license.

Copyright © 2022 Dell Inc. or its subsidiaries. All Rights Reserved. Dell Technologies, Dell, EMC, Dell EMC and other trademarks are trademarks of Dell Inc. or its subsidiaries. Other trademarks may be trademarks of their respective owners.