

Realtime
publishers

The Essentials Series: New Techniques
for Creating Better Backups

Reducing Backups with Data Deduplication

sponsored by

ARCserve®
More than Backup

by Eric Beehler

Reducing Backups with Data Deduplication	1
Explaining Data Deduplication.....	2
Data Deduplication for Backups: Breakdown of Benefits and Risks	3
How to Integrate Deduplication	4
Deduplication Helps with More than Daily Backups.....	5
Deduplication Is Now an Industry Standard Technology.....	6

Copyright Statement

© 2009 Realtime Publishers. All rights reserved. This site contains materials that have been created, developed, or commissioned by, and published with the permission of, Realtime Publishers (the “Materials”) and this site and any such Materials are protected by international copyright and trademark laws.

THE MATERIALS ARE PROVIDED “AS IS” WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE, TITLE AND NON-INFRINGEMENT. The Materials are subject to change without notice and do not represent a commitment on the part of Realtime Publishers or its web site sponsors. In no event shall Realtime Publishers or its web site sponsors be held liable for technical or editorial errors or omissions contained in the Materials, including without limitation, for any direct, indirect, incidental, special, exemplary or consequential damages whatsoever resulting from the use of any information contained in the Materials.

The Materials (including but not limited to the text, images, audio, and/or video) may not be copied, reproduced, republished, uploaded, posted, transmitted, or distributed in any way, in whole or in part, except that one copy may be downloaded for your personal, non-commercial use on a single computer. In connection with such use, you may not modify or obscure any copyright or other proprietary notice.

The Materials may contain trademarks, services marks and logos that are the property of third parties. You are not permitted to use these trademarks, services marks or logos without prior written consent of such third parties.

Realtime Publishers and the Realtime Publishers logo are registered in the US Patent & Trademark Office. All other product or service names are the property of their respective owners.

If you have any questions about these terms, or if you would like information about licensing materials from Realtime Publishers, please contact us via e-mail at info@realtimepublishers.com.

Reducing Backups with Data Deduplication

It seems, at least on a personal technology level, we don't think about the limitations of data storage much anymore. With cell phones having more storage than servers and personal computers did only 10 years ago, it seems the need to trim and conserve storage bits has become almost irrelevant. However, this outlook becomes a problem in an IT data center. With budgets being cut back and disk storage outrunning backup capacity at every turn, administrators need a way to keep up with the need to put all that ever-expanding data somewhere.

Why should you even worry about the size of data? The first reason is limited resources. Although data storage is quite a bit cheaper than it used to be, SANs and large tape backup libraries are still expensive to purchase. In fact, with many companies cutting IT budgets, it is nearly impossible to get a requisition signed for new hardware in many companies. The second reason is inefficiencies. Any systems administrator will recognize the fact that the same file, with the same information, is usually stored multiple times, being backed up multiple times, and taking up extra room that is unnecessary. When it continues to happen, it can have an exponential effect on storage growth. Tape storage tends to lag hard disk storage in capacity, so the issue is even more critical when looking at the backup scenario.

When presenting the benefits of data deduplication, the easiest target to point at is traditional backup hardware. When looking at it from a backup perspective, take 20TB of data that you back up every night. Using a normal rotation, and assuming there is 2% growth with 3% change in that data every week, you'll be backing up 100TB to 110TB of data over a 5-week period. That is 100TB worth of tapes or disk-based backup for 20TB of disk data. Now consider the fact that those tapes, whether used as primary backup or as an off-site solution, have to be managed and stored, and some are likely to be stored off site, costing your company in additional space and transportation. On top of that, if tape is the primary backup medium, you have to provide the capacity to make full backups of that data through sufficient tape backup drives and libraries within a specific window, usually overnight. If those backups run over into prime-time working hours, it can introduce a new set of problems with resource utilization, open files, and specialized applications such as email and databases.

There are specific benefits that involve reduced data center footprint for the storage and backup solutions. There is also a real cost benefit in power consumption, especially when considering certain tape libraries and the power it takes to cool the data center. Those power savings also translate into environmental savings. Some other not-so-obvious pluses are the reduced network traffic for those backups using the LAN. With that comes reduced backup windows, which become more of an issue as data continues to grow. If there are any backups occurring from remote offices over a WAN or off-site backups occurring over an Internet connection, the time-to-backup-over-network benefit dramatically increases.

Explaining Data Deduplication

Data deduplication is really a complex mathematical puzzle that essentially dehydrates data, making it smaller, but with the ability to make it what it was through a rehydration process. This works by identifying the duplicate files, bytes, or blocks of data and removing those non-unique pieces. The complicated calculations are handled behind the scenes by the software. A quick and dirty estimate using duplicate files will gain you 10 to 20 percent in storage space. These files are examples of unstructured data, as opposed to structured data, in databases. (With structured data, deduplication will happen at the database system level, not at the file level.)

Virtualization is a prime candidate for data deduplication. Think of all the snapshots of virtual machines that get backed up. All these snapshots tend to have a large amount of duplicate data. Just the operating system (OS) files alone can end up being tens of thousands of duplicate files. Since they exist in a snapshot file, there is no way to just save changed files; the snapshot will be different every time it's backed up. In these scenarios, a typical deduplication backup will offer a 4 to 1 savings and, as time goes on using a post-processing system, you can see 10 to 1 on the low end, 20 to 1 on average, and realistically up to a 26 to 1 reduction or more in storage space needed for backup.

As files are backed up in a traditional solution, those files will be compared with dates, file system tags such as the archive bit, or other journaling methods to look for changes; however, the system doesn't look for data that hasn't changed. For example, consider an Outlook PST file, where the mail stored there is mostly the same, except for more recent additions. A backup job will copy that entire file every backup, whether it's a full or incremental or differential type of backup. That's because a backup solution bases its decision on the fact that the file changed at all since the last backup. With deduplication, only those changes to data will be stored in the most recent backup. The parts of the file that have not changed will not be stored.

Many deduplication implementations are based on encryption-style methods of representing data at the bit- or block-level instead of comparing actual files. Really, it's a hash that represents a dataset and is compared with other new hashes to determine whether the data is new or has already been represented. An MD5 or SHA-1 hash could be created for a block of data, for example. When the same hash is compared and found, the system will create a pointer to that data, understand it is the same data, and save the pointer instead.

The challenge then becomes where to store all that data. Will the processing be in-line before it reaches the backup medium or afterwards? Running in-line, on-the-fly deduplication, also known as *source deduplication*, is certainly more processor intensive. In fact, the process can exist either in the backup agent software itself or through an external appliance. All the heavy lifting is processor related, so an appliance can slow down due to the amount of processing it has to complete, slowing backups. If done on the backup client itself, the process can affect the performance of the server's primary function.

The other option is to store that backup data somewhere else and then have a dedicated server or appliance dig through the data for deduplication. This setup has the advantage of quick backups, leaving the data to travel as fast as possible to a central repository, leaving the backup window small and processor free to do other things. Post-processing, also referred to as *target deduplication*, usually relies on byte-level deduplication as opposed to hashes based on a block of data. Comparing bytes with previous data bytes is more efficient and avoids the chance of a hash collision where duplicate hashes are created based on data that is not a duplicate (although that is mathematically unlikely to happen).

The downside is the investment in duplicate space. The data center must have the same amount of space to store the data for post-processing as there was to store the data in the first place. There is also the issue of delayed backup as the post-processing occurs. While the system crunches the numbers, the data is not at its final resting place. You also have to deal with added complexity in the environment, where in-line methods are more straightforward and familiar, especially when compared with normal backups.

The speed in which this kind of processing needs to occur and the large amounts of space required have given way to a new class of storage known as *near-line*. Near line disk systems use lower cost, slower hard drives than their server counterparts. Often the technology is the same found in desktops, based on SATA drives and using slower 7200RPM spindle speeds instead of 10,000 and 15,000RPM SAS drives used in servers and SANs. Data access speed isn't nearly as important as capacity in this case, and the cost benefit can put them close to tape. In fact, these systems are built to be virtual tape libraries (VTL). The backup software writes to the VTL disk as if it were actually a tape device. Some systems replace tape in all cases except for archival purposes, leaving faster disk as the substitute. The VTL layer over this storage is a virtualization that allows traditional backup software to address disks as if they were a tape library, keeping the advantages of the tape software cataloging system. Some systems don't virtualize the disk space as a VTL and simply address it as a traditional NAS or file server.

Data Deduplication for Backups: Breakdown of Benefits and Risks	
Benefits	Risks
<ul style="list-style-type: none"> —Backup storage savings on average 20 to 1 —Less hardware, media, and power consumption —Remote site backup uses less bandwidth —Can reduce backup time window —Becoming a standard technology in backup 	<ul style="list-style-type: none"> —CPU intensive for in-line processing —New hardware, especially disk space, is required for post-processing de-dupe —Requires new backup methodology —Restores can take longer

How to Integrate Deduplication

When first considering data deduplication solutions, you should assess what your expected ratios of savings will be. The math could be done manually, but there are tools available from various storage and backup vendors to collect information on the data types in the environment, so a realistic ratio estimate can be attained. The other consideration is time to process. If a source in-line deduplication system is used, the processing will occur at the client level. Applications with high-transfer rates will require much higher processing load, which will equate to more time to backup. Also, the type of data deduplication chosen will affect ratios. Post-processing methods tend to give better ratios and leave the client system resources free from processing the deduplication algorithms.

Another consideration is constrained backup windows, which may cause you to avoid applying data deduplication when using in-line deduplication on certain servers. This can also affect post-processing solutions because they tend to work best when using full backups all the time; the system will get to see all data every cycle, resulting in better ratio efficiencies. However, this option might not be possible due to backup windows.

When integrating data deduplication in the environment, the consideration should be whether a move to near-line storage over straight tape backup is desired. This can serve in-line processing, where it would work as a VTL, and is required for post-processing. The requirement for disk space in a post-processing world is as much as the environment backs up. It is essentially a duplicate copy of all the backed up data. This is then processed, put to more permanent backup storage in another disk system or traditional tape, and then deleted to start the process again the next day. Remember that this will not help conserve network bandwidth, so if you are trying to conserve WAN bits for a remote office, deduplication needs to be done at the server.

Restoring data can be a time issue. The complex process that breaks down the data into small chunks and sets of pointers also has to put that data back together when restoring. The time to restore will depend on the environment, as there are many variables. Note that there is a level of processing that doesn't happen in a traditional restore, when the wait is normally for the network and tape drive transfer rate. Instead, the data has to be rehydrated after it's pulled from the backup.

When implementing data deduplication, you will want to have a clearly defined timeframe in which you want your backups to be complete. Measure the performance of the servers during backup to understand the current baseline of those systems. Include CPU statistics on the servers when you are considering in-line deduplication. You should include servers from varying application and service types. Also, there likely are additional changes that need to be made to the backup agent software. Jobs need to be set up in order to push the updated software to those servers and make the appropriate setting changes. If implementing a VTL or some other deduplication appliance, ensure compatibility to your existing backup software or the backup software of choice. An appliance-based solution may not need any client agent software updates. Be sure you work out these details and know what leverage you have to make changes before moving forward.

You'll likely be surprised when running deduplication for the first time. Be ready to take statistics from your server during a first run with data deduplication enabled. When using an agent-based in-line solution, you're likely to note a significant increase in CPU utilization. This is to be expected and is a part of using the technology. The real question is whether the systems are handling the load well. Check performance and ensure no critical applications are falling outside of performance expectations. This is especially important for your critical applications and services that need to serve their purpose even during the backup window. Also compare the total time it takes now that deduplication is enabled. If, for example, you see an average of an extra 2 hours to complete backups, you might want to modify the established backup routine to include certain servers in an earlier window or possibly move full backups to a different day when a longer backup will not affect normal business.

When using a deduplication appliance and VTL solution or target deduplication, you're likely to notice that the overall backup time is quite long at first. Even though the backups from the servers will run as normal, and actually be a bit faster because they are now streaming to disk instead of tape, the deduplication process will still need to work on the data after it is loaded to the near-line storage solution. After the initial backup, the appliance will be able to find duplicate data much more consistently now that it has indexed the data. As time goes on, the backup time should decrease and ratios of duplication should increase. In fact, if full backups are allowed to run every backup cycle, the appliance will have consistent data to compare and be even more efficient, reducing the overall backup window.

When running backups over the Internet or other wide area network (WAN) connections, a central post-processing system will give no benefit to reducing network bandwidth. Some vendors offer solutions that exist at the remote location that will compress the data before transmitting it over the wire. Otherwise, deduplicate the data using an in-line method from an enabled backup agent. With in-line deduplication, there may just be some systems that will not be able to run the deduplication features because the CPU utilization is too great a stress on the server or the amount of data is too great to fit a backup into the backup window. Thus the need for close monitoring of servers of different types and functions.

Deduplication Helps with More than Daily Backups

Deduplication can also help in areas other than traditional backup. Data deduplication is being used increasingly to help with disaster recovery. This might seem strange because the purpose of disaster recovery is to duplicate your IT capabilities in another location. Many IT departments are not just relying on a set of tape backups for a recovery and restore. Tapes can be problematic, as they have a built-in lag when they are collected, stored, and recovered for restoration. The recovery time objective tends to be one of the longer available disaster recovery solutions.

Now replication is being employed, sending backups across a WAN connection to a recovery location. This can be a choice between a very expensive connection to supply the necessary bandwidth or reducing the number of systems to be recovered. Because data deduplication reduces the amount of data that needs to be transmitted, it opens up the solution to allow for many more systems that can be recovered or allow use of a less-expensive network WAN connection to the remote recovery site. This kind of recovery scenario requires using disk-based backups. It is also worth noting that post-processing methods of data deduplication can affect recovery point objectives because of the extra processing time required before the final backups are saved. Take this into account.

Deduplication Is Now an Industry Standard Technology

The simple fact is that deduplication is becoming standard fare in enterprise backup systems, and there is no reason it shouldn't move into most data centers and IT organizations. The continued acceleration of storage capacity means administrators need to adapt with the infrastructure they have today. With a simple feature upgrade, in-line data deduplication can be an easy implementation, with the proper monitoring of performance to avoid problems. Implementing an appliance solution and moving to disk-based backup opens new abilities. Better performance and the ability to easily move backups offsite over the wire makes a compelling case to integrate this technology into not only daily backups but also disaster recovery plans. The administrator needs to keep in mind that the processing required to provide that deduplication layer means time and CPU utilization, whether it's at the client or at the deduplication appliance. Rolling out this solution requires deliberate steps and a close eye on performance. Adjust expectations and service level agreements (SLAs) accordingly once you know how deduplication affects your environment.