

Data Deduplication Explained

By: Stephen J. Bigelow, Features Writer, SearchStorage.com

Data deduplication eases storage requirements and enhances retention

Data is flooding the enterprise. Storage administrators are struggling to handle a spiraling volume of documents, audio, video and images, along with an alarming proliferation of large email attachments. More storage is often not the best answer—storage costs money and the sheer number of files eventually burdens the company's backup and disaster recovery (DR) plans. Rather than finding ways to store more data, companies are turning to data reduction technologies that can store less data. Data deduplication has recently emerged as an important part of any data reduction scheme. This article explains the basic principles and implementation issues of data deduplication and covers some examples of the technology at work today.

Understanding data deduplication

Data deduplication is basically a means of reducing storage space. It works by eliminating redundant data and ensuring that only one unique instance of the data is actually retained on storage media, such as disk or tape. Redundant data is replaced with a pointer to the unique data copy. Data deduplication, sometimes called intelligent compression or single-instance storage, is often used in conjunction with other forms of data reduction. Traditional compression has been around for about three decades, applying mathematical algorithms to data in order to simplify large or repetitious parts of a file—effectively making a file smaller. Similarly, delta differencing reduces the total volume of stored data by comparing the new and old iteration of a file and saving only the data that had changed. Taken together, these techniques can be very effective at optimizing the use of storage space.

When properly implemented, data deduplication lowers the amount of storage space required so it saves money on disk expenditures. A more efficient use of disk space also allows for longer disk retention periods, which offers better recovery time objective (RTO) for a longer time and reduces the need for tape backups. Data deduplication also reduces the data that must be sent across a WAN for remote backups, replication and disaster recovery.

Data deduplication primarily operates at the file, block and even the bit level. File deduplication is relatively easy to understand—if two files are exactly alike, one copy of the file is stored and subsequent iterations

receive pointers to the saved file. However, file deduplication is not very efficient because the change of even a single bit results in a totally different copy of the entire file being stored. By comparison, block and bit deduplication looks within a file and saves unique iterations of each block. If a file is updated, only the changed data is saved. This behavior makes block and bit deduplication far more efficient. "It's an order of magnitude difference in terms of the amount of storage that it [block deduplication] saves in a typical environment," says W. Curtis Preston, vice president of data protection at GlassHouse Technologies Inc. Other analysts note that deduplication can achieve compression ratios from 10-to-1 to 50-to-1. However, block and bit deduplication take more processing power and use a much larger index to track the individual blocks.

Data deduplication platforms must contend with the issue of "hash collisions." Each chunk of data is processed using a hash algorithm, such as MD5 or SHA-1, generating a unique number for each piece. The resulting hash number is then compared with an index of the existing hash numbers. If that hash number is already in the index, the piece of data is a duplicate and does not need to be stored again. Otherwise, the new hash number is added to the index and the new data is stored. In rare cases, the hash algorithm may produce the same hash number for two different chunks of

Find out more at.

Storage Decisions

**San Francisco
December 4-6, 2007**

**The Best 3 Days a Storage Pro Can Spend Out
of the Office**

**Apply today – Seating is extremely limited at
this FREE conference**

www.storagedecisions.com/sanfran/

data. When such a hash collision" occurs, the system fails to store the new data because it sees that hash number already. This is called a false positive and can result in data loss. Some vendors combine hash algorithms to reduce the possibility of a hash collision. Some vendors are also examining metadata to identify data and prevent collisions.

Implementing data deduplication

The data deduplication process can typically be implemented in hardware within the actual storage system, but it is also appearing in backup software. Hardware-based implementations are often easier to deploy and are primarily interested in reducing storage at the disk level within the appliance or storage system. Software-based implementations also reduce data, but the reduction is performed at the backup server. This minimizes the bandwidth between the backup server and backup system—particularly handy if the backup system is located remotely. "Users get 'end-to-end' benefits when deduplicating data at the source—less data traverses the WAN, LAN and SAN," says Lauren Whitehouse, analyst at the Enterprise Strategy Group. However, deploying deduplication in a new backup application is a bit more disruptive because it involves installing lightweight agents on the systems that must be backed up—in addition to installing the new backup engine.

There is no universal approach to data deduplication, and your re-

data tends to not dedupe very well," he says. Test a prospective data deduplication platform with various types of backups and restores, and see how it functions under actual circumstances.

Scalability is another issue that has attracted significant attention, especially in terms of performance as the data deduplication system grows. Performance might have been an issue as early hash indexes grew large and additional time was needed to look up each block, but Preston calls that FUD (fouled up data) marketing now. "All of the vendors that I am aware of that are currently shipping or about to ship have addressed this [scaling issue] in one way or another," he says. Still, he recommends that you discuss the issue with your data deduplication vendor and see how it dealt with scaling concerns.

From a management perspective, data deduplication should not present any noticeable increase in overhead. "It [management] shouldn't be any more or less than just a standard VTL [virtual tape library]." When multiple deduplication devices are needed, however, there could be an incremental increase in management effort.

The impact of data deduplication

The Appalachian and coastal areas South Carolina are enticing attractions to tourists and regional industry. Advertising, communication and literature have emerged as key assets to the Department of Parks, Recreation and Tourism—the agency responsible for promoting tourism as an industry and maintaining an extensive park system throughout the state. The agency originally had an EMC Corp. storage area network (SAN) hosting a total of 4 terabytes (TB), of which 1.2 TB comprised the actual working data set of databases and files, while 2 TB was allocated for disk backups before being relegated to DLT. Like many IT organizations, the agency sought ways to mitigate the increasing storage demands of its media and other data.

After investigating numerous data deduplication vendors, the agency settled on Data Domain Inc.'s 430 appliance for disk backup tasks. With 2 TB of onboard storage, the 430 replaced the 2 TB that had previously been set aside on the SAN. The reduction in space was dramatic with bit level deduplication. "With the compression and deduplication, I think we're using about 900 MB," says Bernie Robichau, the agency's systems administrator and security officer. The space reduction was a welcome cost savings, but it also allowed much longer backup retention on disk. "If someone had requested a two-week old file, I would have never been able to get that from a disk-based backup because I couldn't keep two sets of backups on our allocated 2 TB of hard drive [SAN] storage," Robichau notes. "Now someone can request a file from three weeks ago or six weeks ago, and it's immediately available."

Robichau says that installation of the data deduplication platform was relatively quick and easy, requiring only about four hours of onsite engineering work and minimal configuration. Its current CommVault System Inc. backup infrastructure proved to be fully



Hear it Here from the Storage Experts

Storage Soup is the SearchStorage.com Blog—featuring posts daily from the best storage experts in the world.

www.storage.blogs.techtarget.com/

sults can vary dramatically depending on the environment and selected product. It's important to note that data deduplication only makes sense when long-term retention is involved—usually for backup and archive tasks. Short-term retention sees little benefit because there is nothing to deduplicate against. Preston cautions against the misinformation circulating between deduplication vendors and suggests focusing on the key issues of performance, capacity and cost. Due diligence can identify potential performance and compression issues in your specific environment. "Let's say you're backing up seismic data or medical imaging data—this

compatible—backup agents were simply pointed to the new appliance rather than the EMC SAN. "The backups worked just as they always did, but we're consuming far less disk space and much more retention than we ever did before," he says. While the deduplication appliance requires almost no management time, Robichau notes as much as 75% labor savings in tape overhead, such as cartridge rotation, cleaning and storage. The only remaining tape effort involves full backups on weekends and systematic cartridge rotation to an offsite location.

Although there are no immediate plans to upgrade storage on the 430 appliance, the attention is clearly focused on disaster recovery. Previous considerations of complex disaster recovery plans were put on hold due to complexity. However, the 430 supports replication easily and Robichau expects to replicate the 430 to a duplicate appliance and eliminate backup tapes entirely sometime in the next fiscal year or beyond. "There's no planning beyond synchronizing an identical appliance on site and putting it in one of our remote locations."

Denver-based IT hosting provider, Data393 Holdings LLC was drowning in customer data. Its challenge: to keep its data protection business running smoothly, along with other services, like managed server hosting, managed firewalls and load balancing. However, its backup environment was formidable; handling 20,000 backups per month with each customer protecting 20 GB to 100 GB. Even with 4.5 TB of protected storage, Data393 could only keep two weeks of retention. To make matters even more challenging, its StorageTek L700 and L11000 tape libraries were managed by an outsourced provider, requiring a full time engineer at Data393.

But, it was ongoing restoration problems that really forced Data393 into action. "Our success rate from backups, at the lowest point, was roughly 70%," says Steve Merkel, senior systems engineer. "And far too often, we couldn't hit [restore] the exact day they wanted." Poor performance of the tape backup process also plagued the organization, with full backup windows often exceeding 18 hours. These problems also translated into significant customer support costs. It became clear to Merkel that disk storage was the key to beating reliability and performance woes, and data deduplication would be essential to reduce the total volume of storage needed for full and incremental customer backups.

Data393 opted for Avamar Technologies Axion software running on a cluster of 11 Dell 2850s offering about 10 TB of total storage. The actual deployment involved a forklift upgrade, but Merkel reports that the system was up and running in just a few days after installing agents on almost 400 backup servers and migrating necessary data. The move to data deduplication brought several significant benefits. Most notable was the reduction of storage requirements. For example, it might take 350 GB to protect 100 GB of customer data without deduplication (full and incremental backups). With data deduplication, it actually takes less storage than the data it's protecting. "I'm using about 7 TB of stor-



Storage RSS Center

Click here to get storage industry news and the latest tips and blog postings via RSS feed from SearchStorage.com

www.searchstorage.com/storagerss

age to protect roughly 8 TB of data," Merkel says. "That includes anywhere from two weeks to one year of retention [daily full backups]." Backup time was also slashed; in some cases an 18-hour backup window fell to one and a half hours, while improving the backup and restoration success rate to 98% or more. The need for two full-time engineers fell 75% to one half of one full-time engineer. "We wanted to have an ROI [return on investment] of 24 months, and we hit payback at 20 months," he says.

Today, the 4.5 TB of protected data has grown to about 7.6 TB protected by data deduplication. About 2 TB of that protected data is replicated to a smaller Avamar deployment at a disaster recovery site in St. Louis. Data393 continues to use tape for long-term archival backups. Merkel expects the amount of protected data to double in the foreseeable future, though less storage will be required to handle the growth.

The future of data deduplication

In the near term, industry experts see data deduplication taking an important role in disaster recovery—saving disk storage space by replicating the data of one deduplication platform to another located off site. This alleviates the need to move tapes back and forth, which can be particularly meaningful when replicating hundreds of terabytes of data.

Other analysts note that the separate "point products," like VTL, will address backup window performance, while data deduplication addresses the issue of storage capacity. "Next-generation backup solutions fix both," Whitehouse says, "deduplicating data as it's sourced from the backup target and improving the efficiency of its transfer across a LAN/WAN to the central disk repository." Deduplication is now common in VTLs and will appear as a feature of traditional backup products.