



**White Paper**

**Fixing Disk Latency and I/O  
Congestion to Improve Slow  
VMware Performance**

## Executive Overview

Virtualization delivers great benefits for IT shops; however, it is plagued with the same performance issues that affect physical servers. These problems arise from the fact that in virtualized systems there are multiple instances of Windows server competing for the resources on a single host and multiple hosts connected to a SAN. Windows itself is the origin of many of these I/O congestion, disk latency and throughput problems. The Windows file system is notorious for creating lots of small I/O, which increases SCSI traffic across the storage stack, clogs queues and increases physical disk I/O. VMware never touches the Windows guest system, but it must deal with the workload it produces. This paper describes how Windows creates the problems and how to address these issues at their source; it discusses how to proactively manage and resolve virtualization performance problems.

## Virtualization Performance Issues -- Finding The Cause

There is no dispute that virtualization has provided IT shops around the world with enormous benefits. Virtual systems are easier to create and maintain. You can pack more of them per square foot of datacenter space than their physical counterparts. They also provide a considerable savings in hardware expenditures and power. Along with virtualization, some nifty new hardware technology like fiber channel and SANs have sought to enhance the virtual experience. Despite these advances, virtualized systems still have performance issues. The source of these performance issues is often not readily apparent.

Zeus Kerravala, a former Senior Vice President and Distinguished Research Fellow at the Yankee Group, wrote a blog article in the August 24, 2012 online *Network World* entitled: *VMware weather forecast: Watch out for performance storms* (<http://www.networkworld.com/community/node/81283>). Mr. Kerravala cites a survey by ZK Research and Xangati that asked about performance storms. Here are some of the survey's highlights:

- 30% of the respondents said they were the hardest track down.
- 20% claimed they were transient in nature.
- 18.8% said these issues were not flagged by current monitoring tools.
- 79% said these problems took more than an hour to resolve
- 39% said these problems took more than four hours to resolve.

From these data points, Mr. Kerravala concludes that a large number of companies are experiencing performance storms that are transient in nature and undetected by current virtualization monitoring tools. These are expensive problems. In 25% of the cases, performance storms caused the organization to roll back its virtualization to a physical server.

## The Windows File System's Impact on Virtualization Performance

The source of many performance storms could be a source you would never suspect: NTFS, the Windows file system. TechNet, Microsoft's knowledgebase of all things Windows, has long documented the impact of NTFS on system performance. It is important to understand the impact NTFS has and how it affects virtualized systems.

Let's look at what happens when you create a file with Windows. For our example, let's say we are creating a 2GB file. Initially, NTFS creates a record in the Master File Table (MFT), which is the index to the volume. This record is 1KB in size and contains the file name; file ID, the Extent List and other attributes. Next, NTFS will ask the \$Bitmap file for 2GB of space.

The \$Bitmap file is a NTFS metadata file containing one bit for every cluster on the disk. It is created when the disk is formatted. The \$Bitmap file keeps track of free and used disk clusters with respect to how NTFS "sees" the disk. It is important to note that this is different from how the disk controller sees the disk. When file space is needed, the \$Bitmap file allocates the required number of free clusters to accommodate the file. In the ideal scenario, the space allocated would be a contiguous string of clusters. The space allocated is identified by the starting Logical Cluster Number (LCN), which is recorded in the Extent List in the MFT record, along with the length of the string.

As is often the case, the ideal is hard to achieve. When users create, extend and delete files, the \$Bitmap file becomes a patchwork of free and used space, and large contiguous strings of clusters are harder to come by. A more realistic view of what will happen is for \$Bitmap to offer up 2GB of free space wherever it can find it; for our example, let's say in 200 locations. The starting LCN and length of all 200 locations are recorded in files Extent List in the MFT. For Windows, this is a fragmented file and nothing has been written to the disk yet.

This is where the performance storm starts to brew. In our ideal example, the 2GB file is represented by one entry in the Extent List. A single SCSI command with 2GB of data goes across the storage stack to the storage controller. The storage controller receives the SCSI command and understands it needs to find 2GB of physical disk space. The controller software takes over and the LCN address is mapped to physical storage on the disks (SAN or otherwise). To keep the math simple, let's assume the 2GB file is written in 10 physical I/O of 200MB each. The workload for the hypervisor and controller was processing 1 SCSI command; and for the disks it was processing 10 writes.

Now look at the typical scenario. Again, we'll assume our file is in equal size pieces. Each of the 200 chunks of the file listed in the Extent List is a separate SCSI command. Now, instead of one

SCSI commands across the storage stack, there are 200 commands and each is associated with 10MB of the file. The controller processes each SCSI command. We'll say each is written to the disks in two physical I/O of 5MB each. The workload for the hypervisor and controller was processing 200 SCSI commands (200x more than the ideal) and for the disks it was processing 400 writes (40x more than the ideal). The hypervisor now needs more CPU and memory to do its job and the disk queues are full of smaller I/O.

In a virtual environment, you have this behavior running on every Windows Server. There are multiple instances of Windows Server running on each host and competing for its resources. Several hosts are probably connected to the same SAN. Users are creating, editing and deleting files. As the virtual disks fill up, the free space gets more fragmented. When NTFS can't find sufficient free space to create a file in one piece, the result is file fragmentation. It takes longer to read a fragmented file and it takes longer to write a fragmented file. As I/O contention increases, your virtual machines get slower and slower.

When the disk queues are full of relatively small I/O, the disks are slamming away to process the workload. Disk latency, the time it takes a disk to process an I/O, increases. EMC and VMware contend that if latency is 15ms, you need to keep an eye on things; at 30ms, you have a problem. VMware introduced Storage I/O Control (SIOC) with VMware 4.1. It basically monitors latency and throttles throughput when the response time gets to 30ms. While this offers a modicum of relief, is reducing throughput a real solution? VMware does not touch the Windows guest; it simply deals with the workload it presents. In the case of NTFS, this can clobber VM performance.

## **Attacking the Source of the Problem**

The only way to address this problem is to attack it at its source in the Windows guest system. When the files are contiguous and the free space is consolidated, files are read and written in the fastest possible way. However, there are other benefits as well. As we saw in our example, when the file was contiguous, there was less SCSI overhead for the hypervisor and controller, reducing the demand for CPU and memory. Contiguous files mean larger I/O to process and correspondingly less congestion in the disk queues. Larger I/O generally improves sequential I/O at the storage level, thereby reducing disk latency with a positive impact on throughput.

We tested our PerfectDisk defragmentation software for VMware to quantify the benefits of disk optimization in a virtual environment. We used two sets of identical disks where the baseline set had fragmented files and free space and the second set was optimized with

PerfectDisk. We used VMware's *vscsiStats* utility to collect the data. The metrics captured showed some dramatic improvements:

- 28% reduction in total I/O across the stack
- 1200% increase in the largest I/O (>524K)
- 49% reduction in disk latency
- 28% increase in throughput.

## Real-world Examples

While test metrics are great, how does fragmentation manifest itself in the real world? Here are three examples of different types of performance problems where PerfectDisk disk optimization helped.

- The town of Castle Rock, Colorado did weekly full system backup of 7+TB and it was taking an unacceptable 72 hours. They had 50+ thin-provisioned VMs backed by a NetApp SAN and a CommVault Enterprise Backup system. The system administrator used PerfectDisk on the target disks in the CommVault system and the next backup was 12 hours faster. The administrator then used PerfectDisk on all the VMs and shaved another 12 hours off the backup.
- NSC manages the nationwide online ordering system for a national pizza chain. It has to be able to handle 100 orders per second across its 100+ VMs. NSC noticed that as workload increased, performance dropped and the demand for CPU and memory spiked. Its database journal and log files grow with the workload. The free space was so fragmented it was taking Windows too long to allocate space. When Windows is looking for free space, nothing else happens on that machine. PerfectDisk optimized their VMs and the problem went away. PerfectDisk is now part of their routine maintenance.
- Pasavant Hospital adopted a new electronic medical records application that severely fragmented files to the detriment of its 50 VMs running on an EMC VNX backend. PerfectDisk now keeps the files contiguous and performance has improved.

Returning to Mr. Kerravala's survey respondents, let's examine how fragmentation maps to their observations.

- Thirty percent said performance issues were hard to track down. In our experience, when performance issues arise, most customers look at VM density, the network and the hardware first. In the backup case cited above, none of Castle Rock's vendors mentioned file or free space fragmentation as a potential source of the problem. Ask yourself, would you even consider NTFS fragmentation on the Windows guest as a potential source of a VM performance problem? Probably not.
- Another 20% noted that the problems were transient in nature. In the pizza chain case, server performance plummeted and CPU and memory spiked as order volume increased and the system struggled to find contiguous free space to grow the log files. At lower workloads, the problem was less noticeable. Free space fragmentation contributes to transient performance issues, especially with write intensive tasks.
- The last 18.8% of the respondents said the problems were "under the radar" of their monitoring tools. Most VM monitoring tools measure CPU, memory and I/O. However, if the I/O spikes are due to fragmentation, these tools don't identify fragmentation as the underlying issue. They just show resources are being consumed.

As Kerravala states: "performance storms are happening, they are hard to find and they are costing the company because they take so long to isolate and respondents aren't getting help from VMware."

Prior to virtualization, corporate IT departments recognized fragmentation as a source of performance issues on physical servers. It should surprise no one that this is now the case in virtualized environments. With virtualization, multiple instances of the same file system are running on a single host and sharing its resources. In addition, multiple hosts are now accessing the same SAN via queues with finite depth. As we have explained here, fragmentation increases the SCSI workload over the virtual storage stack, increasing hypervisor CPU and memory demand. The storage controller has more work to do and the number of accesses to the physical disks increases. These are the necessary ingredients for a perfect performance storm.

## Conclusion

Performance storms can originate from a variety of sources, but some common-sense analysis can identify the most likely contributors. A look at frequently-cited VMware performance issues shows that disk latency and I/O congestion are at the forefront. This makes perfect sense. In the average business environment, I/O is the predominant resource used. Most business activity is read/write activity and not particularly CPU or memory intensive. Secondly, the hard drive is the slowest component in the computing environment. Anything increasing the number of I/O the system needs to process is just adding to the workload and ultimately creating a congestion problem. Windows guest systems that are not optimized increase the I/O workload. File fragmentation exponentially increases the number of I/O a virtual machine produces. There is a corresponding workload increase for the hypervisor, the controller and the physical disks. Each instance of Windows multiplies the problem.

VMware and EMC are correct in identifying disk latency as an important metric of I/O performance affecting system response and throughput. SIOC is not an appropriate solution; the problem has to be addressed at the source in the Windows file system. In our testing, PerfectDisk produced a 49% reduction in the total number of I/O taking 15ms or longer. At the same time, it improved system throughput by 28%.

The purpose of this paper was to highlight a likely cause of the serious performance issues in virtualized systems and what can be done about it. The case examples cited show that virtual application and system performance can be improved by implementing disk optimization as part of routine system management.

For more information or a free trial of PerfectDisk contact Raxco at [www.raxco.com](http://www.raxco.com).