



VDI & Storage: Deep Impact



PQR

Eenvoud in ICT

Author(s) : Herco van Brug
Version: 1.25
Date: September 2, 2011

© 2011 PQR, all rights reserved.

All rights reserved. Specifications are subject to change without notice. PQR, the PQR logo and its tagline "Eenvoud in ICT" are trademarks or registered trademarks of PQR in the Netherlands and/or other countries. All other brands or products mentioned in this document are trademarks or registered trademarks of their respective holders and should be treated as such.

REFERENCES

Reference	Title
http://www.vmware.com/files/pdf/resources/vmware-view-reference-architecture.pdf	VMware View Reference Architecture
http://h71028.www7.hp.com/ERC/downloads/4AA2-3017ENW.pdf	HP: Implementing a VDI infrastructure
http://support.citrix.com/servlet/KbServlet/download/19754-102-376939/XD%20-%20Design%20Handbook.pdf	Citrix XenDesktop Design Guidelines
http://technet.microsoft.com/en-us/library/cc748650.aspx	Address Space Load Randomization
http://technet.microsoft.com/en-us/library/aa995867(EXCHG.65).aspx	Windows disk alignment
http://download.microsoft.com/download/C/E/7/CE7DA506-CEDF-43DB-8179-D73DA13668C5/DiskPartitionAlignment.docx	
http://en.wikipedia.org/wiki/Redundant_Array_of_Independent_Disks	Wikipedia: RAID levels
http://en.wikipedia.org/wiki/Scsi	Wikipedia: SCSI
http://en.wikipedia.org/wiki/Integrated_Drive_Electronics	Wikipedia: ATA
http://en.wikipedia.org/wiki/Flash_memory	NOR and NAND Flash memory
http://virtuall.eu/blog/creating-a-vdi-template	Create a VDI image that minimizes IOps
http://virtuall.eu/blog/hbr/ita-s-sample-time-with-vscsistats	Sample blocksizes on VDI Luns

HISTORY

Version	Date	Author(s)	Remarks
1.1	29-03-2009	Herc0 van Brug	2.1: added boot and logon storms 4.3: Corrected a few calculations 6.2: Extended SSD explanation
1.2	5-10-2009	Herc0 van Brug	1.1 added 3.2.4 removed 5. removed DP references 6.3 added
1.25	2-9-2011	Herc0 van Brug	7 Introduction to blocksize

CONTENTS

1.	Introduction.....	1
1.1	Intended audience	1
1.2	Update 1.2.....	1
2.	The client IO.....	2
2.1	Boot and Logon storms.....	2
3.	The storage IO.....	4
3.1	SCSI vs ATA.....	4
3.2	RAID level	4
3.3	Disk alignment	5
3.4	Prefetching and Defragging	6
4.	The mathS.....	7
4.1	Processor.....	7
4.2	Memory	7
4.3	Disks	8
4.4	Practical numbers.....	8
5.	Summary.....	9
6.	Alternatives.....	10
6.1	cache	10
6.2	SSD.....	10
6.3	Serializing random writes	11
7.	Block size	12
7.1	VDI storage profile	12
7.2	The IOPS multiplication	12
8.	In conclusion	13
9.	About	14
9.1	About PQR.....	14
9.2	About the author.....	14

1. INTRODUCTION

Virtual Desktop Infrastructure, or VDI, is hot. It's cool, secure, centrally managed, flexible - it's an IT manager's dream.

VDI comes in two flavors; Server Hosted VDI (Centralized, single-user remote vDesktop solution) and Client-Side VDI (local, single-user vDesktop solution).

The advantages of a VDI infrastructure are that virtual desktops are hardware independent and can be accessed from any common OS. It is also much easier to deploy virtual desktops and to facilitate the freedom that the users require of them. And because of the single-user OS, application compatibility is much less of an issue than it is with terminal servers.

However, when implementing a VDI infrastructure certain points need to be addressed. First of all, the TCO/ROI calculation may not be as rosy as some people suggest. Secondly, the performance impact on applications, specifically multimedia and 3D applications, needs to be investigated. And finally, don't forget to check licensing aspects, as this can be a very significant factor in VDI infrastructure.

While centralized desktop computing provides important advantages, all resources come together in the datacenter. That means that the CPU resources, memory resources, networking and disk resources all need to be facilitated from a single point - the virtual infrastructure.

The advantage of a central infrastructure is that, when sized properly, it is more flexible in terms of resource consumption than decentralized computing. It is also more capable of handling a certain amount of peak loads, as these only occur once in a while on a small number of systems in an average datacenter.

But what if the peak loads are sustained and the averages are so high that the cost of facilitating them is disproportionate to that of decentralized computing?

As it turns out, there is a hidden danger to VDI. There's a killer named "IOPS".

1.1 INTENDED AUDIENCE

This document is intended to warn people not to oversimplify storage design, especially for VDI projects. It does not try to be a fully accurate storage sizing guide. Some information is aggregated to a simple number or rule of thumb to allow people who are not in-depth storage experts, to get the general message and still know what's going on.

When I talk about IOPS per disk, the table simply shows the number of IOPS that a disk will be able to handle. It is however more complicated than that. The SAN setup, disk latency, seek time, etc. all influence the net IOPS a disk can handle. But for the sake of argument, an aggregated number is used that is extracted from practice and several vendor storage sizing calculators. The numbers used here are correct for general setups but some storage vendors will argue that they can handle higher IOPS per disk, which they probably can. But now you'll have a good starting point for the discussion.

1.2 UPDATE 1.2

Despite several people from NetApp reading this whitepaper before it was released, the RAID DP information was only partly correct. In <http://media.netapp.com/documents/tr-3298.pdf> they state that RAID DP is comparable to RAID 4 for performance. But there's more to it than that. The typical Read/Write ratio of VDI highlights another feature of NetApp arrays. The way their WAFL works allows random IOs to be consolidated and effectively written to disk sequentially. Therefore, writing to a WAFL with 100% random blocks is faster than reading from it. That means that the write penalty for RAID DP is not comparable to RAID 1 at all. Therefore the RAID DP chapter has been replaced by a section about serializing random IO (chapter 6.3).

2. THE CLIENT IO

A Windows client that is running on local hardware has a local disk. This is usually an IDE or SATA disk rotating at 5,400 or 7,200 RPM. At that rate it can deliver about 40 to 50 IOPS.

When a Windows client starts, it loads both the basic OS and a number of services. Many of those services provide functionality that may or may not be needed for a physical system and make life easier for the user. But when the client is a virtual one, a lot of those services are unnecessary or even counter-productive. Indexing services, hardware services (wireless LAN), prefetching and other services all produce many IOPS in trying to optimize loading speed, which works well on physical clients but loses all effectiveness on virtual clients.

The reason for this is that Windows tries to optimize disk IO by making reads and writes contiguous. That means that reading from a disk in a constant stream where the disk's heads move about as little as possible is faster than when the head needs to move all over the disk to read blocks for random reads. In other words, random IOs are much slower than contiguous ones.

The amount of IOPS a client produces is greatly dependent on the services it's running, but even more so on the applications a user is running. Even the way applications are provisioned to the user impacts the IOPS they require.

For light users the amount of IOPS for a running system amounts to about three to four. Medium users show around eight to ten IOPS and heavy users use an average of 14 to 20 IOPS. The exact number of IOPS for an environment can be determined by running an assessment in your current environment with tools like Liquidware Labs or other capacity planning tools.

Now the most surprising fact; those IOPS are mostly WRITES. A great many researchers have tested the IOPS in labs and in controlled environments using fixed test scripts. The read/write ratio turned out to be as high as 90/10 as a percentage. But in reality users run dozens or even hundreds of different applications, whether virtualized or installed. In practice, the R/W ratio turns out to be 50/50 percent at best! In most cases the ratio is more like 30/70, often even 20/80 and sometimes as bad as 10/90 percent.

But why is that important? Most vendors don't even mention IOPS or differentiate between reads and writes in their reference designs.

2.1 BOOT AND LOGON STORMS

When implementing a VDI infrastructure, it should be optimized for IOPS as that is one of the main bottlenecks. Powering on a VM at the time a user actually wants to use it, would put unnecessary strain on the infrastructure. The user would have to wait relatively long before he can log on and booting a machine creates IOPS that can also be read before working hours start. So for those two reasons, it makes sense to power on the amount of machines you will use at a certain day before users actually start using them. On the other hand, 90 to 95 percent of the IOPS during boot are reads. If you have a large enough cache or a cache that dedicatedly captures OS reads, the strain on the storage would be minimal and hosts would suffer more from impact on memory and CPU than on storage. All in all, when designed correctly, the impact of boot storms on a VDI infrastructure can be minimized or even avoided completely during production hours.

Logon storms are something else though. The impact on IOPS when a user logs on is very dependent on the way the profiles and policies are set up and also on how application are delivered. The way to deal with this, is to not only optimize the VM image but also optimize user environment management. When done correctly, the impact of logon storms can be reduced to a manageable factor. In practice, the read/write ratio during login turns out to be 80/20 percent to even 95/5 percent. While the amount of write IOPS during logon can be more than double the amount of average production use, by far the most amount of IOPS are reads

which are much more manageable than writes. Nevertheless, if all users start working at the same time, you need a substantial larger amount of IOPS (and storage) to accommodate them than when the logons are dissipated over several hours. The only way to properly design for this is to know when users log in. For example, if a logon takes 30 seconds and the maximum amount of simultaneous logons at a peak moment is 10%, with double the amounts of write IOPS and 10 times the amount of read IOPS, the storage would need to facilitate an additional 36% of IOPS compared to regular production use. But if 3% of the users log in simultaneously, the storage only needs to facilitate an additional 11% of IOPS.

Besides logging in, there's a 3rd factor to take into account; application's first run. The first couple of minutes after a user logs on, applications will start for the first time. As it turns out, the read/write ratio during this stage is about 50/50 percent. When the first runs are completed, the reads drop by a factor of about 5 on average but the amount of writes stay the same! That means that a few minutes after a user logs in, the R/W ratio goes down from 50/50 percent to 20/80 percent which is what people see in almost all production environments. It's important to notice that the reads go down but the writes don't. And it's the writes that cause trouble.

3. THE STORAGE IO

When all IOs from a client need to come from a shared storage (attached directly to the virtualization host or through a Storage Area Network) and many clients read and write simultaneously, the IOs are, from the storage point of view, 100 percent random IOs.

3.1 SCSI vs ATA

There are two main forms of disks - SCSI and ATA. Both have a parallel version (regular SCSI vs IDE or PATA) and serial version (SAS vs SATA).

The main differences between the architecture of the SCSI and ATA disks are rotation speed and protocol. To start with the protocol, the SCSI protocol is highly efficient with multiple devices on the same bus, and it also supports command queuing. ATA devices have to wait on each other, making them slower when grouped together.

The higher rotation speed means that when the head needs to move to a different location, it does not need to wait as long for the data to pass beneath it. So a SCSI disk can produce more IOPS than an ATA disk. The faster a disk rotates, the less time the head needs to wait before data passes beneath it and the sooner it can move to the next position, ergo the more IOs it can handle per second.

To give some idea of the numbers involved; a 15,000 RPM disk can handle about 180 random IOPS, a 5,400 RPM disk about 50. These are gross figures and the number of IOPS that are available to the hosts depend very much on the way they are configured together and on the overhead of the storage system. In an average SAN, the net IOPS from 15,000 RPM disks is 30 percent less than the gross IOPS.

3.2 RAID LEVEL

There are several ways to get disks to work together as a group. Some of these are designed for speed, others for redundancy or anything in between.

3.2.1 RAID5

The way a traditional RAID5 system works is that it writes the data across a set of hard disks, calculates the parity for that data and writes that parity to one of the hard disks in the set. This parity block is written to a different disk in the set for every further block of data.

To write to a RAID5 system, the affected blocks are first read, the changed data is inputted, the new parity is calculated and the blocks are then written back. On systems with large RAID5 sets this means a write IO is many times slower than a read IO. Some storage systems, like HP's EVA, have a fixed set of four blocks for which parity is calculated, no matter how many disks are in a group. This increases overhead on a RAID5 group because every set of four disks needs a fifth one, but it does speed things up. Also, on most storage systems, write operations are written to cache. This means that writes are acknowledged back to the writing system with very low latency. The actual write to disk process takes place in the background. This makes incidental write operations very speedy, but large write streams will still need to go directly to disk.

With 15,000 RPM disks the amount of read IOPS are somewhere in the 150-160 range while write IOPS are closer to the 35-45 range.

3.2.2 RAID1

A RAID1 set is also called a mirror. Every block of data is written to two disks and read from either one. For a write IO to occur, the data doesn't need to be read first because it does not change part of a parity set of blocks but rather just writes that single block of data. This means that writing to a RAID1 is much faster than to a RAID5.

With RAID1 the data is read from one of the two disks in a set and written to both. So for 15,000 RPM disks, the figures for a RAID1 set are still 150-160 IOPS for reads, but 70-80 for writes.

3.2.3 RAID0

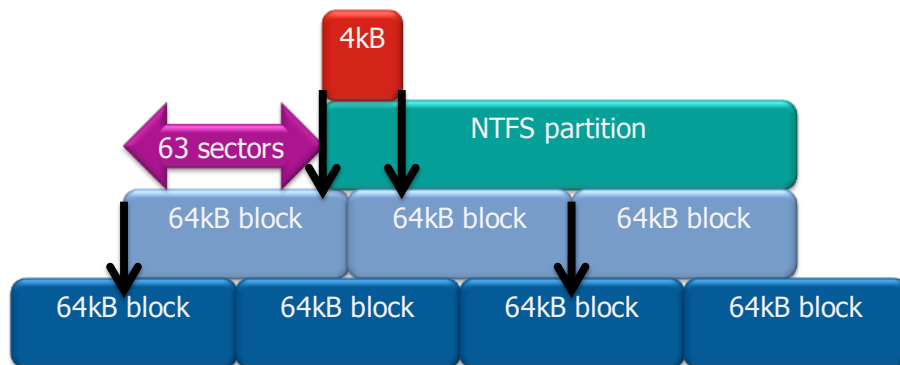
RAID0 is also called striping. Blocks of data are written in sequence to all disks in a RAID0 set but only to one at the time. So if one disk in the set fails, all data from the set of disks is lost. But because there is no overhead in a RAID0 set, it is the fastest way of reading and writing data. In practice this can only be used for volatile data like temporary files and temporary caches, and also perhaps for pagefiles.

If used, the amount of IOPS a RAID0 set can provide with 15,000 RPM disks is 150-160 for reads and 140-150 for writes.

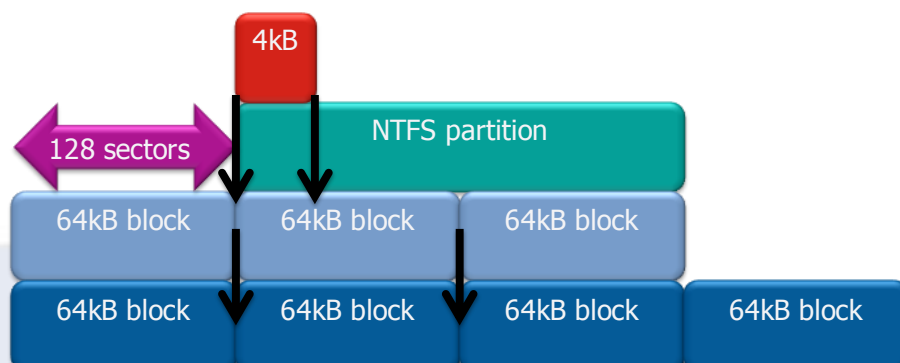
3.3 DISK ALIGNMENT

Because we want to minimize the amount of IOPS from the storage we want every IO to be as efficient as possible. Disk alignment is an important factor in this.

Not every byte is read separately from the storage. From a storage perspective, the data is split into blocks of 32 kB, 64 kB or 128 kB, depending on the vendors. If the filesystem on top of those blocks is not perfectly aligned with the blocks, an IO from the filesystem will result in 2 IOs from the storage system. If that filesystem is on a virtual disk and that virtual disk sits on a filesystem that is misaligned, the single IO from the client can result in three IOs from the storage. This means it is of utmost importance that all levels of filesystems are aligned to the storage.



Unfortunately, Windows XP and 2003 setup process misalign their partition by default by creating a signature on the first part of the disk and starting the actual partition at the last few sectors of the first block, misaligning the partition completely. To set this up correctly, create a partition manually using 'diskpart' or a Linux 'fdisk' and put the start of the partition at sector 128. A sector is 512 bytes, putting the first sector of the partition precisely at the 64 kB marker. Once the partition is aligned, every IO from the partition results in a single IO from the storage.



The same goes for a VMFS. When created through the ESX Service Console it will, by default, be misaligned. Use fdisk and expert mode to align the VMFS partition or create the partition through VMware vCenter which will perform the alignment automatically.

Windows Vista and later Windows versions try to properly align the disk. By default they align their partition at 1 MB, but it's always a good idea to check if this actually is the case¹.

The gain from aligning disks can be 3-5 percent for large files or streams up to 30-50 percent for small (random) IOs. And because a VDI IO is an almost completely random IO, the performance gain from aligning the disks properly can be substantial.

3.4 PREFETCHING AND DEFRAGGING

The NTFS filesystem on a Windows client uses 4 kB blocks by default. Luckily, Windows tries to optimize disk requests to some extent by grouping block requests together if, from a file perspective, they are contiguous. That means it is important that files are defragged. However, when a client is running applications, it turns out that files are for the most part written. If defragging is enabled during production hours the gain is practically zero, while the process itself adds to the IOs. Therefore it is best practice to disable defragging completely once the master image is complete.

The same goes for prefetching. Prefetching is a process that puts all files read more frequently in a special cache directory in Windows, so that the reading of these files becomes one contiguous reading stream, minimizing IO and maximizing throughput. But because IOs from a large number of clients makes it totally random from a storage point of view, prefetching files no longer matters and the prefetching process only adds to the IOs once again. So prefetching should also be completely disabled.

If the storage is de-duplicating the disks, moving files around inside those disks will greatly disturb the effectiveness of de-duplication. That is yet another reason to disable features like prefetching and defragging.

¹ A quick way to check if a partition is aligned is by typing "wmic partition get BlockSize, StartingOffset, Name, Index" in a command shell. If the number isn't a multiple of 65536 (64 kB) or 1048576 (1 MB) the partition is unaligned.

4. THE MATHS

So much for the theory. How do we use this knowledge to properly size the infrastructure?

4.1 PROCESSOR

On average, a VDI client can share a processor core with six to nine others. Of course, everything depends on what applications are being used, but let's take an average of 7 VMs per core. With a dual socket, quad core CPU system that means we can house $7 \times 2 \times 4 = 56$ clients. However, the Intel Nehalem architecture is very efficient with hyper-threading and allows 50-80 percent more clients. That means that when it comes to the CPU, we can host $150\% \times 56 = 84$ VMs.

4.2 MEMORY

The amount of memory the host must have depends primarily on the applications the users require and the OS they use. On average a Windows XP client needs 400-500 MB of RAM for basic operations and a standard set of applications. Add some caching and the memory usage should stay below 700 MB.

The Windows OS starts paging when 75 percent of its memory is allocated. It will always try to keep at least 25 percent free. But paging in virtual environments is a performance-killer. So instead of giving it the recommended (in physical systems) amount of 1.5 to 2 times the amount of memory in swap space, we limit the pagefile size to a fixed amount of 200 to perhaps 500 MB. If that is not enough, just add more RAM to the client, rather than extending the pagefile.

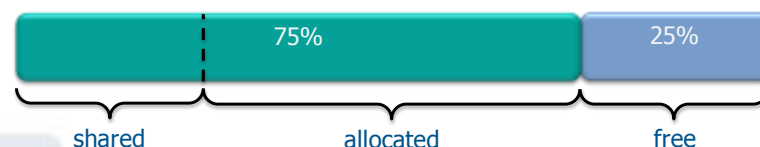
This also means we aim for at least 25 percent free RAM space with most applications running. Additionally, about half of the used memory contains the same blocks in all clients (Windows DLLs, same applications, etc). This is lower on Windows 7 clients because of ASLR (Address Space Load Randomization), which means that the amount of memory shared between clients is 25% (empty space) + $75\% / 2 = 62.5\%$.



So when running Windows XP on ESX servers, if 60 percent of memory per client is actually being used, 50 percent of which is shared between clients, we need $1 \text{ GB} \times 60\% \times 50\% = 300$ MB per client. Every VM needs about 5 percent more than the amount allocated as overhead from the host. So you need an additional 50 MB (5 percent of 1 GB) per client.

We have seen from the CPU calculation that we can host 84 clients, so a host would need 4 GB (for the host itself) + $350 \text{ MB} \times 84 =$ at least 34 GB of RAM.

However, if 75 percent of memory is used and only a third of that can be shared, every client needs $1 \text{ GB} \times 75\% \times 67\% = 512$ MB of dedicated host memory. So for 84 clients the host needs $4 \text{ GB} + (512 + 50) \text{ MB} \times 84 = 52$ GB of RAM.



Of course if you run on a host that doesn't support transparent page sharing, the amount of memory needed is $4 \text{ GB} + 84 \times (1024 + 50) \text{ MB} = 96$ GB of RAM.

For Windows 7 clients the numbers are $(2\text{ GB} + 100\text{ MB}) \times 60\% \times 50\% = 660\text{ MB}$ per client on average, $4\text{ GB} + 660\text{ MB} \times 84 = 60\text{ GB}$ of minimum host memory and $4\text{ GB} + 84 \times (2\text{ GB} + 100\text{ MB}) = 188\text{ GB}$ per host if the host doesn't support memory over-commitment.

4.3 DISKS

The amount of IOPS a client produces is very much dependent on the users and their applications. But on average, the IOPS required amount to eight to ten per client in a read/write ratio of between 40/60 percent and 20/80 percent. For XP the average is closer to eight, for Windows 7 it is closer to ten, assuming the base image is optimized to do as little as possible by itself and all IOs come from the applications, not the OS.

When placing 84 clients on a host, the amount of IOPS required would be 840, of which 670 are writes and 170 are reads. To save on disk space, the disks are normally put in a RAID5 set up. But to deliver those numbers, we need $670 / 45 + 170 / 160$ (see 'RAID5' section earlier in this document) = 16 disks per host. Whether or not this is put in a central storage system or as locally attached storage, we will still require 16 disks for 84 VMs. If we used RAID1, the number changes to $670 / 80 + 170 / 160 = 10$ disks. That means, however, that using 144 GB disks, the net amount of storage drops from $16 \times 144\text{ GB} \times 0.8$ (RAID5 overhead) = 1840 GB to $10 \times 144\text{ GB} \times 0.5$ (RAID1 overhead) = 720 GB; less than 2.5 times the amount of net storage.

4.4 PRACTICAL NUMBERS

All these numbers assume that clients are well-behaved and that most of the peaks are absorbed in the large averages. But in reality you may want to add some margins to that. To be on the safe side, a more commonly used number of clients per host is 65 (about 3/4 of 84). That means that the *minimum* amount of memory for the average XP client solution would be $65 \times 350\text{ MB} + 4\text{ GB} = 27\text{ GB}$, or for Windows 7: $65 \times 660\text{ MB} + 4\text{ GB} = 47\text{ GB}$.

The amount of IOPS needed in these cases is $10\text{ IOPS} \times 65\text{ clients} = 650\text{ IOPS}$ where 80 percent (= 520) are writes and 20 percent (= 130) are reads. With RAID5 that means we need $(520 / 45) + (130 / 160) = 13$ disks for every 65 clients. Should you require 1,000 VDI desktops, you will need $(1000 / 65) \times 13 = 200$ disks. Put on RAID1, that number decreases to 112, which is quite substantial considering that it only serves nine clients per disk.

So, to be sure of the number you need to use, insist on running a pilot with the new environment where a set of users actually use the new environment in production. You can only accurately size your infrastructure once you see the numbers for those users, the applications they use and the IOPS they produce. Too much is dependent on correct sizing - especially in the storage part of the equation!

5. SUMMARY

The table below summarizes the sizing parameters:

Setting	Windows XP sizing	Windows 7 sizing
Number of VDI clients per CPU core	6-9	6-9
Hyper-threading effectiveness on Intel Nehalem systems	150 – 180%	150 – 180%
Amount of memory per VDI client	1 GB	2 GB
Amount of memory actually allocated (in actual host memory)	Minimum: 37.5% Average: 50% Not shared: 100%	Minimum: 20% Average 50% Not shared: 100%
Number of clients per host	65	65
Minimum amount of memory per host	27 GB	30 GB
Average amount of memory per host	37 GB	76 GB
Amount of memory per host (if not shared)	76 GB	144 GB
IOPS per VDI client (also depends on image optimization)	Light user: 3-4 Medium user: 6-8 Heavy user: 12-16	Light user: 4-5 Medium user: 8-10 Heavy user: 14-20

The following table summarizes the IOPS for the different RAID solutions:

RAID level	Read IOPS 15k	Write IOPS 15k	Read IOPS 10k	Write IOPs 10k
RAID 5	150-160	35-45	110-120	25-35
RAID 1	150-160	70-80	110-120	50-60
RAID 0	150-160	140-150	110-120	100-110

To illustrate the above figures, a few samples follow:

Scenario with 65 clients/host	10 IOPS		5 IOPS	
	R/W 20/80%	R/W 50/50%	R/W 20/80%	R/W 50/50%
VDI clients per disk	RAID5: 5	RAID5: 7	RAID5: 10	RAID5: 14
	RAID1: 9	RAID1: 10	RAID1: 17	RAID1: 21
	RAID0: 15	RAID0: 16	RAID0: 30	RAID0: 31
Number of disks per host	RAID5: 13	RAID5: 10	RAID5: 7	RAID5: 5
	RAID1: 8	RAID1: 6	RAID1: 4	RAID1: 3
	RAID0: 5	RAID0: 4	RAID0: 3	RAID0: 2

6. ALTERNATIVES

6.1 CACHE

There are many solutions out there that claim to speed up the storage by multiple factors. NetApp has its FlashCache (PAM), Atlantis Computing has vScaler, and that's just the tip of the iceberg. Vendors such as Citrix with its Provisioning Server and VMware with its View Composer and storage tiering, but also cloning technology from storage vendors, aid storage by single-instancing the main OS disk, making it much easier to cache it.

But in essence they are all read caches. Caching the IOPS for the 30 percent that are reads, even with an effectiveness of 60 percent, will still only cache $30\% \times 60\% = 18\%$ of all IOs. Of course they can be a big help with boot and logon storms but all write IOs still need to go to disk.

However, most storage systems also have 4 GB, 8 GB or more cache built-in. While the way it is utilized is completely different for each vendor and solution, most have a fixed percentage of the cache reserved for writes, and this write cache is generally much smaller than the read cache.

The fact is that when the number of writes remains below a certain level, most of them are handled by cache. Therefore it is fast; much faster than for reads. This cache is, however, only a temporary solution for handling the occasional write IO. If write IOs are sustained and great in number, this cache needs to constantly flush to disk, making it practically ineffective. Since, with VDI, the large part of the IOs are write IOs, we cannot assume the cache will fix the write IO problems, and we will always need the proper number of disks to handle the write IOs.

6.2 SSD

SSD disks are actually more like large memory sticks rather than disks. The advantage is that they can handle an amazing amount of IOPS; sometimes as high as 50,000 or 100,000. They have no moving parts so accessing any block of data takes mere microseconds, instead of milliseconds. Also, the power consumption of SSD disks is only a fraction of a spinning SCSI disk. An array of SSD disks consumes a only a few 100s of Watts while an array of traditional disks can consume many 1000s of Watts.

There are two types of flash memory; NOR and NAND. With NOR-based flash memory, every bit is separately addressable but writes to it are very slow. NAND-based memory is much faster and because it requires less space per memory cell, has a higher density than NOR memory. It's also much cheaper. Downside of NAND memory is that it can only address blocks of cells at once. For block based devices this works perfect and because of its speed and density, NAND memory is used for SSD disks.

NAND memory traditionally stores one bit per cell. This is called a Single-Level Cell (SLC). Newer types can store multiple levels of electrical charge allowing more than one bit of information per cell. Those are called Multi-Level Cells (MLC). MLC memory is most commonly used because it is much cheaper to make than SLC memory (currently about 4 times cheaper). But the downside is that where SLC memory has a 100,000 write cycle endurance, MLC only allow 5,000 to 10,000 write cycles. Because of the longer lifespan of SLC memory, this type is used by most enterprise vendors in their SAN solutions where cost is subordinate to reliability. Some solutions however, use MLC memory. By caching and optimizing write IOs, the amount of writing of the actual cells can be greatly reduced, thus extending the expected lifetime of MLC memory considerably.

Currently the SSD disks are four to ten times more expensive than fiber channel hard disks and most storage vendors don't guarantee a lifespan of multiple years with an IO profile like that of VDI. But better SSD cells are being developed constantly. With a more even read/write ratio, a longer lifespan, larger disks and better pricing, we may see SSD disks in a SAN become more common within a year or two.

6.3 SERIALIZING RANDOM WRITES

When sizing storage for IOPS, all formulas assume 100% random IO which is especially true in VDI environments.

But when IOs are sequential, a spinning disk can handle much more IOPS than with random IO. The exact numbers when using sequential IO vary between 300 to 400 IOPS to sometimes over 600 to 800 IOPS for a 15000 RPM fiber channel disk depending on the overhead and RAID levels involved.

If a storage system manages to cache random IO and write it to disk sequentially, it would boost the write performance of that system significantly. If the reads for that system are fully random, the impact on read performance would be small. There are several storage systems that achieve this.

NetApp with its ONTAP OS and WAFL filesystem serializes random IO by journaling all writes in NVRAM, coalescing them, and writing the data on the first available block nearest to the head of the spinning disk. This way, the effective write IOPS of a NetApp storage system is some 3 times higher than a RAID 1 system². The exact number of IOPS varies with how much free space is in the system but that's something that can be taken into account when sizing a NetApp system. Sizing for VDI is determined by IOPS which usually leaves plenty of capacity to work with. NetApp has several calculators that can help determine the exact number of disks needed when sizing a system for VDI. It's not uncommon that NetApp's storage systems need half the number of disks compared to more traditional storage arrays for the same amount of IOPS.

Another system that uses serialization is Whiptail. They have SSD storage systems that cache writes in amounts of a full row of SSD cells before writing that row. That way, cells get written far less often than when every IO would be written directly. This effectively eliminates the need for Wear Leveling and prolongs the life of SSDs considerably. That's why they can use MLCs instead of SLCs, making their systems more cost effective. Also, the IOPS such a system can handle runs in the 100.000s.

There are probably more systems out there that can serialize random IO but this should give a pretty good idea of the effects of serialization.

² Section 3 of this document gives a detailed view of how this write optimization works: <http://media.netapp.com/documents/wp-7107.pdf> and section 7.2 of this document: <http://media.netapp.com/documents/tr-3705.pdf> brings it in the context of VDI (VMware View).

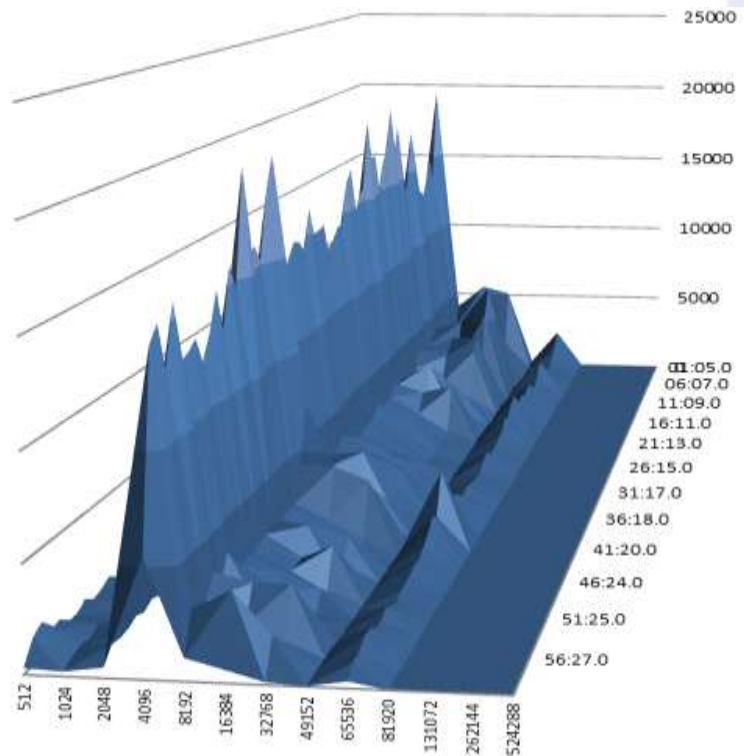
7. BLOCK SIZE

Unfortunately, there is still more to it. After sizing and implementing several dozens of VDI projects, most fit well within the layout set in this document. But some of them deviate horribly. After some investigation it turns out that there are applications out there (mostly the streaming kind) that use a block size that is totally out of scope for VDI projects.

7.1 VDI STORAGE PROFILE

Most VDI environments show a specific mixture of block sizes that we call the VDI Storage Profile. Measured over time we counted the blocks per size entering the storage system. When we draw this up in a 3D graph, it usually looks like shown in the graph to the right.

These are taken from a live VDI environment using the method described in <http://virtuall.eu/blog/hbr/ita-s-sample-time-with-vscsistats>. It uses 60 second intervals to count the amount of different sized blocks. It shows a healthy 75-80% amount of 4kB blocks and calculates down to 10-15kB block size on average.



7.2 THE IOPS MULTIPLICATION

With certain application however, block sizes can be much bigger and as high as 80kB or 100kB on average throughout the day. If your storage is built to use 4kB blocks internally, this could become a problem really fast. Every 100kB block would show as a single IO but become 25 blocks of 4kB each on the storage array. So a client producing 5 IOPS in that block size would actually require $25 * 5 = 125$ IOPS at the storage level. Luckily, these are usually the exceptions and they dissipate away in the large averages. But if those applications are more generally used, sizing the storage can become a daunting task. Also, just imagine what swapping, which usually happens in 64kB blocks, would add to this load.

It's therefore imperative to run a decent Proof of Concept, measure block sizes and know exactly how the intended storage solution handles blocks internally!

8. IN CONCLUSION

It should be obvious by now that calculating the amount of storage needed in order to properly host VDI is not to be taken lightly. The main bottleneck at the moment is the IOPS. The read/write ratio of the IOPS that we see in practice in most of the reference cases demonstrate figures of 40/60 percent, sometimes even as skewed as 10/90 percent. The fact is that they all demonstrate more writes than reads. And because writes are more costly than reads - on any storage system - the number of disks required increases accordingly, depending on the exact usage of the users and the application.

Some questions remain:

- What is the impact of application virtualization on the R/W IOPS?
- What exactly is the underlying cause of the huge difference in read/write ratios between lab tests and actual production environments?
- What if all the write IOs only need to be written to a small part of the total dataset (such as temporary files and profile data)? Could all the data, or at least most of it, be captured in a large write cache?
- How is the average block size determined reliably to predict the *actual* load on the storage and not become surprised when some applications turn out to behave differently from the average applications.

These questions will be investigated as an increasing number of VDI projects are launched.

And as a final note, it is imperative that you run a pilot. Run the actual applications with actual users in the production environment beforehand so that you know how they behave and what the read/write ratio is. If you don't size correctly, everybody will complain. All users, from IT staff to management and everybody in between, will complain and the VDI project... will FAIL.

9. ABOUT

9.1 ABOUT PQR

PQR is the professional ICT infrastructure specialist with a focus on availability of data, applications and work spaces with optimized user experience in a secure and manageable way. PQR provides its customers innovative ICT solutions that ensure the optimization of application availability and manageability, without processes getting complex. Simplicity in ICT, that's what PQR stands for.

PQR has traceable references and a wide range of expertise in the field, proven by many of our high partner statuses and certifications.

PQR's approach is based on four main pillars:

- Data & Systems Availability
- Application & Desktop Delivery
- Secure Access & Secure Networking
- Advanced IT Infrastructure & Management

PQR provides an ICT infrastructure that is stable, flexible and future proof. PQR has extensive experience in designing and implementing server & storage environments, including networking and security. Traditionally, massive storage environments have been PQR's specialty.

PQR, founded in 1990, is headquartered in De Meern and counts over 100 employees.

9.2 ABOUT THE AUTHOR

Herco van Brug was born in 1968 and studied mechanical engineering at the University of Twente in the Netherlands. Immediately after graduation he started working at Rijnhaave, later Syntegra. When Syntegra was taken over by British Telecom his position shifted to that of technical specialist, focusing mainly on specialized solutions and migrations.



At present he is a Solutions Architect at PQR, with his primary focus being business continuity solutions in the datacenter. He is the co-author of the Data & System Availability diagram and is certified for Microsoft, RedHat, Citrix and VMware, while as a VMware Authorized Consultant, he undertakes VMware branded Professional Services assignments. He has been a speaker at several national conferences and published a number of articles, all related to virtualization.



PQR B.V.
Rijnzathe 7
3454 PV De Meern
The Netherlands

Tel: +31 (0)30 6629729
Fax: +31 (0)30 6665905
E-mail: info@pqr.nl
www.PQR.com