



Season III, Challenge I

Project : We are now settled on Mars, but need more permanent infrastructure

Focus Area : VMware vSphere infrastructure, Compute, Storage, Network, Disaster Recovery

Created by : Lubomir Zvolensky (lubomir@zvolensky.sk)

Contents

1. Executive summary	3
2. What needs to be done, what problems we might expect, what can go wrong	4
2.1 Requirements	4
2.2 Assumptions	4
2.3 Constraints and risks	5
3. Physical infrastructure design	6
3.1 Host server information and sizing	6
3.1.1. Scale-up and scale-out environment	9
3.2 Storage information and design	10
3.2.1 VSAN Storage	11
3.2.2 Backup Storage	13
3.3 vSphere vCenter design	14
3.4 vSphere networking design	14

1. Executive summary

Here down on Earth, many people are fed up with politicians and decided to run away from this planet as there is not too much real solutions. Eh, officially they escaped due some kind of zombies or bacteria whatnot, but we've heard so many stories already ...

Preliminary infrastructure has been created on Mars by pioneers, but seems like we need something more stable, mature and permanent. Power and space are very scarce resources up there, so we can't deploy our usual super-expensive mission critical four rack storage or network solutions and wonder why they provide sub-par performance (400 spinning disks have less performance than single SSD in terms of IOPS and latency). They say cooling is also an issue, I say just open the window, there is -240 outside.

As there is no Fiber-Channel complication standing in our way, we have happy smiles on our faces and no "save me, I'm dying" calls during night. Three datacenters across The Red Planet should sustain harsh beating by cosmic rays, radiation, extreme cold ... and debris from human activity. Or may not. Nobody knows. The problem also is, nobody knows how far away these datacenters are apart, if there is any communication between them other than sending a pigeon with floppy disks.

Every human loves oxygen and would be pretty angry if life sustaining Environmental System goes havoc. We need to keep it running no matter what. Other systems, for example the one taking care about food - which doesn't have name - have been voted as very preferred to keep online, too. Respective countermeasures need to be taken, unless we want somebody banging on our door.

By the outlook of it, we will deploy completely new compute resources and storage systems. All mentioned components need to be designed, installed and configured. Some networking gear is present, with description saying "sufficient number of 10Gbit Ethernet ports and some 40Gbit ports available". Not bad !

Here is your helmet, suit, seat, buckle up.

2. What needs to be done, what problems we might expect, what can go wrong

Following are requirements, assumptions and constraints associated with infrastructure :

2.1 Requirements

R01 : RTO of Environmental System 20 minutes (RPO not specified).

R02 : Food System importance, while less than EnvSyst, is very high, although not exactly specified in IT terms. RTO and RPO for this system is unknown.

R03 : in future, unspecified business critical application will be added to this infrastructure. No workload characteristics have been provided but designed solution must be fully able to host it.

R04 : all Mars inhabitants heavily prefer email and collaborative applications to communicate. These applications must be available throughout all three locations.

R05 : infrastructure must be able flexible to scale-up and scale-out due to unknown demands in future.

2.2 Assumptions

A01 : bandwidth between all datacenters is sufficient to perform data replication

A02 : latency between all datacenters is low and stable enough to enable replication of data in order to provide disaster recovery functionality between datacenters.

A03 : sufficient and stable power, space and cooling is available at each site. No data have been provided as of watts, BTUs or rack units available. No data are available as of AC power stability.

A04 : sufficient amount of Ethernet ports is available at each location.

A05 : all requisities are available at each site, such as racks, PDUs, cabling, etc. so installation of computer and storage infrastructure is physically possible.

A06 : scale-up and scale-out principles apply to infrastructure.

A07 : thin and lean infrastructure has been proposed to run email and collaborative platform.

A08 : spinning hard-drives will not work on Mars due to gravitation issues, thus only SSDs will be used.

A09 : SSD usage is preferred due to power consumption, performance and strict RTO requirements, too.

A10 : existing networking infrastructure provides adequate Layer2 and Layer3 for three datacenters.

A11 : price of IT infrastructure is not important too much in the grand scheme of things (getting to Mars)

2.3 Constraints and risks

C01 : no detailed information was provided about bandwidth or latencies between datacenters.

C02 : information about quantity of 10Gbit and 40Gbit Ethernet available at each site is unavailable

C03 : amount of available power (Watts) is unknown for each site. There might be huge differences between three locations in these terms.

C04 : capacity of cooling system and physical space available at each site is unknown. We might not be able to host proposed number of servers at each site.

C05 : roles / hierarchy of datacenters are unknown. Is there something like primary and backup?

C06 : no information is available about future business critical application that needs to run in infrastructure.

C07 : no utilization metrics are available for existing workloads in terms of vCPU, RAM, storage usage, performance profiles

C08 : no standard SLA definitions are provided except requirement to have EnvSyst available within 20 minutes after failure as it is extremely critical.

C09 : EnvSyst application and its processing are unknown (are there live sensors ? database in background ? any form of clustering possible ?)

C10 : no future growth estimates are available. No capacity planning can be performed in terms of CPU, RAM, storage space, performance demands, networking demands.

C11 : reliability of power at each site is unknown. No UPS or generators have been mentioned.

C12 : distance from Earth precludes any vendor support (phone, live tcp/ip sessions) or spare parts availability as there is *3 to 21 minutes* latency. Basically only offline communication (email) is possible.

C13 : multicast must be available between ESXi hosts for VSAN functionality.

3. Physical infrastructure design

Following chapters deal with physical infrastructure design.

3.1 Host server information and sizing

Addressed Requirements : R01,R02, R03, R05

Addressed Assumptions : A03, A04, A05, A06, A08, A09, A10, A11

Addressed Constraints : C02, C03, C04, C06, C07, C08, C09, C10, C11, C12

For compute infrastructure, blades, hyper-converged and stand-alone servers have been evaluated and their advantages and disadvantages weighed.

Generally the size of existing and future workloads is unknown. I created infrastructure with extremely generous configuration so that it is able to supply compute, storage and network resources for many workloads. I did not specifically care about life-sustaining systems mentioned in executive summary (oxygen, water, waste, food) and “some” business critical app, this was rather the general infrastructure that can be used, AMONG OTHER THINGS, for applications already mentioned.

Basic idea here is to replace typical hyper-converged infrastructure (ie. many small blades) with “hyper-converging inside the mammoth box” to achieve further energy, space and cooling efficiency. We will try to SCALE-UP and push environment to run more virtual machines on small amount of huge servers.

Intel platform is preferred to AMD. Intel has better reputation in enterprise-class environments.

SuperMicro SuperServer 2048U-RTR4 model has been chosen for compute nodes, homepage <http://www.supermicro.nl/products/system/2U/2048/SYS-2048U-RTR4.cfm>

Each server will be configured with :

- four E5-4669 v3 CPUs, providing 18 physical cores each
- 3TB DDR4 low-voltage load-reduced 1.2V DDR4 modules for energy efficiency

- 24x Samsung 1.6TB SSD for VSAN all-flash, model MZIES1T6HMJH, 12Gbit SAS, 10 DWPD
- 16GB SuperMicro SATA disk-on-module SSD-DM016-PHI for ESXi v6.0 installation
- three additional dual-port Mellanox ConnectX-3 adapters, 56Gbit/40Gbit/10Gbit Ethernet
- additional SuperMicro adapters to utilize all 24 internal drive slots and provide redundancy paths

Reasons to choose this particular model and configuration :

- **2U form factor, providing extreme amount of CPU cores and power-efficient 1.2V DDR4 RAM**
- 4 socket provides 72 physical cores total.
- 48 DDR4 slots, up to **3TB RAM enabling to achieve extreme consolidation and monster VMs**
- **24x hot-swap 2.5" drive bays, four NVMe capable.** Important for VSAN.
- **9x PCI-Express v3.0** (some used for additional NICs)
- Dedicated IPMI LAN port for remote management with virtual media and virtual KVM over LAN
- **1000W titanium level efficiency (>96%) redundant power supplies.**
- **Six** networking ports, providing 56Gbit private Mellanox Ethernet bandwidth and industry-wide 40Gbit/10Gbit Ethernet compatibility.

Where power and cooling conditions allow, we suggest to use all four CPU sockets in these servers. In case of power constraints in some remote datacenter(s), CPU is the major component to save on and configuration with two 18core CPUs is recommended. In this case, maximum usable RAM drops to 1.5TB (because 24 memory modules can be occupied instead of all 48) and number of usable PCI-Express slots drops to 6 which is still more than sufficient for our expansion needs. In this configuration, identical number of SSDs will be used, 24.

The reason for extreme amount of RAM is to achieve highest possible consolidation, to run maximum number of VMs on host and to minimize number of physical hosts while minimizing energy consumption. Experience shows that TYPICALLY in standard environments, RAM is the primary and most demanded resource today (and the major obstacle in achieving even greater consolidation ratios), generally with plenty of CPU resources still being available far beyond RAM is exhausted.

No CPU hungry application has been identified on Mars as of yet. Design of hosts' infrastructure, mainly huge amount of CPU cores, will easier accommodate "surprises" in this area than blades allow.

When minimally loaded, these servers do not achieve energy effectivity of smaller blades, but the situation changes under load where we expect the infrastructure to leverage benefits arising from aggregation of resources. For example instead of six smaller specified blades with 512GB RAM, one proposed host with 3TB RAM can be used in lesser physical footprint, consuming less power, demanding less cooling and providing less fault domains at the same time.

In 8U rack unit space, four servers configured as above provide :

- 576 logical cores, 288 physical cores, 2.1GHz each (4x 18 cores each server)
- 12TB DDR4 RAM (3TB RAM each server)
- 24x 40Gbit Ethernet ports (compatible with 10Gbit), up to 72x 40Gbit/10Gbit ports max.
- 153.6TB flash capacity (38.4TB each server, 4x 24bays total) , not counting PCI-E VSAN Cache

This is the density of “efficient logical datacenter” - the footprint to match. Identical amount of RAM and flash storage would require twelve to sixteen blades to be used depending on particular model and manufacturer. It is impossible to match number of networking ports and their speed in 8U footprint due to lack of expansion possibilities, ie. number of PCI-Express slots, with blade infrastructure.

While we have no information about bandwidth required on network and future growth, it definitely helps to have the flexibility of installing as much physical NICs. 40Gbit Ethernet NICs have been chosen with future upgrade in sight. Today, 10Gbit Ethernet infrastructure is widely available and will be used, with some 40Gbit ports for highest demands. All 40Gbit Mellanox adapters can operate at 10Gbit setting because they are backwards compatible.

Furthermore, in pure Mellanox environment (Mellanox switches), these adapters can run Ethernet topology at non-standard 56Gbit speed, enhancing throughput and decreasing latencies furthermore compared to industry-wide 40Gbit standards. By default, 56Gbit data rate is available for InfiniBand technology only which we will not use due to lack of switches. Information about existing network has not been provided, so we don't know if there are Mellanox switches or not.

Fiber optic QSFP cables are recommended due to lower energy consumption and much lower latencies compared to metal counterparts.

Fully understanding extreme power concerns at remote planet, our calculation shows it is more effective in terms of consumed power to run mammoth space effective stand-alone server above than higher number of blades providing the same compute resources :

[didn't manage to create table with power calculation]

Advantages :

- Monster VMs. If there are 1TB/1TB+ RAM monster VMs necessary in future, this platform can host it. Typical blades are configurable with 1TB RAM max – with potential requirement for 0.5TB / 1TB virtual machine(s), we would not achieve any significant consolidation, which will lead into power/space/cooling inefficiency because more physical blades will be necessary to

host our workloads. Potentially, there might be VMs blades are not able to run at all, such as VM with 1.5TB RAM requirement while blade would top at 1TB physical RAM.

- These hosts are easily capable of hosting extremely CPU hungry VMs as up to 72 physical cores is available in each. Blades usually provide 16 to 24 physical cores max.
- When working with blades or hyper-converged infrastructure, it is much easier to run into hardware configuration limits in terms of maximum number of CPU cores, available RAM, number of available SSD drives for VSAN purposes and networking throughput (limited amount of PCI-Express slots, limited amount of additional NICs affecting available bandwidth, redundancy etc).
- **Quite often, blade enclosure/chassis might represent single point of failure. Although chassis failures are extremely rare and unlikely, we recommend against their usage at planet as remote as Mars with no real vendor support, no 4hrs contracts to meet.** Enclosures represent additional layer of complexity and failure domain.
- **Power saving : typically, 4node blade enclosures with similar number of CPU cores, RAM and storage resources use 1280W+ power supplies, quite often 1600W ; designed solution uses 1000W PSUs only !!**
- Due to the complexity of integration, this product is sold as completely assembled system only. SuperMicro sales representative must be contacted.

Hosts will be tremendously expensive. Because this is Mars, nothing will be cheap : hundreds of millions must be spent to get a box of matches to Mars, so price of IT infrastructure is not important that much as it is down on Earth.

3.1.1. Scale-up and scale-out environment

Extremely high density and powerful configuration (CPU cores, amount of RAM, number of SSDs) will be used from beginning in attempt to achieve extreme consolidation rates so there is not too much scale-up possible. 72 cores and 3TB RAM are achieved in 2U footprint together with storage.

Scale-out possibilities are limited by existing networking structure, available power, space and cooling at each datacenter. It is easy to add identical servers, spread the load between them, extend VSAN storage capacity etc. Because these specifics are completely unknown at this time, we only can state “yes, scaling out is easily possible with our infrastructure depending on given local conditions”.

Due to VSAN technology we use, at least three physical servers are necessary per each site with four nodes recommended.

Four servers will be used in datacenters where power infrastructure allows. VSAN 6.0 with VirstoFS (v2) will be set to re-protect data in case of failure of one host, returning to fail-safe status as quickly as possible.

3.2 Storage information and design

Addressed Requirements : R01, R02, R03, R05

Addressed Assumptions : A02, A03, A04, A05, A08, A09, A10, A11

Addressed Constraints : C02, C03, C04, C06, C07, C08, C10, C11, C12, C13

Due to power and space constraints, no external disk storage array will be used. Instead, VMware’s built-in VSAN functionality will provide primary storage for workloads, with single big external 4U storage unit serving for backups. Replication of data for business continuity and disaster recovery will be performed by Zerto platform.

VSAN has been chosen due to space and energy constraints – no external SAN technology can be afforded. Energy balance is much better as we have no external big box to feed ; energy consumption overhead due to additional SSD drives in compute infrastructure is very little (any form of external storage will have to have SSDs nonetheless, so there is no “penalty” of this solution).

External SANs typically require skilled administrator that will be little hard to find on Mars. Distance from Earth and round-trip times preclude any remote support, no phone calls are possible, no interactive sessions will work, no spare parts are available for months to come.

External chassis for backup is chosen in order to separate production storage from cold backup data. We can’t rely on single technology to store and protect our data. FreeNAS solution from iX-Systems has been chosen due to flexibility, synchronous and asynchronous replication possibilities, block and file-sharing via network (iscsi, nfs, cifs protocols), ZFS file-system used internally, protection against failures of up to three drives in any logical entity.

VCG, VMware Compatibility Guide must be followed at all costs for VSAN. Only certified controllers and SSDs are used.

3.2.1 VSAN Storage

Virtual SAN 6.0 supports up to 200 Virtual machines per each ESXi Host in version 6.0, with a maximum of 6,400 Virtual machines per cluster.

There is a maximum of 5 disk groups (flash cache device + capacity devices) on an ESXi host participating in a Virtual SAN cluster. We use all-flash configuration, which has maximum of 7 flash devices per disk group for the flash capacity layer and maximum 1 flash device for cache per disk group.

The largest “component size” on VSAN is 255GB. VSAN automatically divides storage objects greater than 255GB into multiple components, for example 2TB VMDK will be split into eight components in the same RAID. Configuration maximums must be considered. Maximum VMDK Size is 62TB so there is no practical limit.

Each physical host will use twenty-four SSDs to provide VSAN capacity split into four disk groups, each containing six drives. Samsung SM1683, 1.6TB model MZIES1T6HMHJH has been chosen due to enterprise-class performance, reliability and features :

- 12Gbit SAS, dual-port active-active interfaces
- 10 DWPD = 28500TBW guaranteed, it will sustain much more due to 4x nm 3D V-NAND tech
- 200k random 4K read ; 50k random 4K write IOPS performance ; 1.4GB/s read/write bandwidth
- 4.7W power consumption compared to 8W-12W of 15k rpm disks.
- power capacitors to protect data-in-flight in case of power outage

Full specification available here :

http://www.samsung.com/global/business/semiconductor/file/product/Datasheet_SM1635_v05.pdf

This configuration provides 38.4TB raw storage in each server with 2U footprint with no single point of failure for storage part. With performance figures above and 24 disks in each server, we are talking 3.6 million random 4K IOPS and 33GB/s theoretical throughput, of course this will be severely limited by controllers, throughput of PCI-Express bus, etc. Due to extreme amount of SSDs, their huge capacity and spreading the load, write endurance becomes much less of an issue as it is with small one or two SSD hybrid-array configurations.

VSAN Health Check plugin will be used to monitor conditions of VSAN environment and report operational issues.

With four hosts at site, following VSAN settings will be used :

NumberOfFailuresToTolerate: 1

NumberOfDiskStripesPerObject: 6

ForceProvisioning : no

ObjectSpaceReservation : 100%

FlashReadCacheReservation : 0% (not used with all-flash VSANs)

“Number of Failures to Tolerate” represents amount of additional copies of each storage object. We believe one single additional copy is adequate as setting of two will lead to much less effective resources usage (wasted storage space, which in turn might increase demands for number of physical servers necessary to cover particular storage capacity which is hard to justify in our power/space/cooling situation) and will require at least 5 physical servers per site. We count with 4 servers at site, additional ones added only in capacity constraint situation, so usage of NFT of 2 is impossible in initial configuration.

“Number of disk stripes per object” defines how many physical disks should be used to write single storage object, effectively raising performance.

“Force provisioning” setting allows VSAN to violate NumberOfFailuresToTolerate, NumberOfDiskStripesPerObject and FlashReadCacheReservation policy settings during the initial deployment of a virtual machine. Forcing provisioning in these conditions might easily lead to capacity issues and reduced availability of storage object in case of Maintenance Mode of particular host. Given the risks, we will NOT use forced provisioning.

Object Space Reservation is configured at 100% so there will be easy and straight-forward capacity planning. We need to play it safe here so no kind of thin provisioning will be used.

Only 10Gbit or faster NICs can be used with all-flash VSAN configuration due to the potential for an increased volume of network traffic. Adapters used for VSAN replication can be shared with other traffic types, for example with vMotion, but we designed sufficient number of ports to have dedicated NICs.

Every SSD used for capacity purposes will be marked as such.

Jumbo Frames of 9000 bytes can reduce CPU Utilization and improve throughput. Because this setting needs to be configured end-to-end throughout infrastructure and we don't have networking information available, we rather play this safe again with default 1500 byte setting. Because vSphere uses TCP Segmentation Offload (TSO) and Large Receive Offload (LRO) to deliver performance benefits and combine interrupts to save CPU processing, Jumbo frames might not have too much performance benefit at all – on the risk side, with improper settings, communication will fail completely.

Another special consideration is multicast which is a network requirement for VSAN. Multicast is used to discover ESXi hosts participating in the cluster as well as to keep track of changes within the cluster. It is mandatory to ensure that multicast traffic is allowed between all nodes participating in a VSAN cluster.

3.2.2 Backup Storage

We need to separate backup data from main storage. Stand-alone backup storage will be implemented in the form of 4U chassis with 72 Sandisk 2TB SSDs. Following configuration will be used:

- 4U SuperMicro chassis (www.supermicro.nl/products/system/4U/6047/SSG-6047R-E1R72L.cfm)
- 36x SAS/SATA hot-swap bays, two disks in each tray
- Total raw capacity : 144TB.
- Usable capacity : 102TB, three 24disk RAID-Z2 disk groups will be created.
- CIFS/Samba presentation will be used for backup purposes. Veeam backup software running in virtual machine will connect to this storage box over network via SMB/CIFS.
- CPU choice is not crucial : 6core E5-2418 v2 has been chosen due to 50W TDP
- 16GB RAM
- Two dual-port Mellanox ConnectX-3 cards, providing total four 40Gbit ports
- iX-systems FreeNAS operating system has been chosen due to proven-and-tested stability
- ZFS filesystem protects data against bit-rot issues.
- Sandisk SSDs provide 1 DWPD write endurance which is assumed to be sufficient for our needs due to large number of SSDs working in parallel, spreading out the load. Assumption is we are not going to back up (compressed) 72TB a day on average.

Storage space will be used to provide retention for backups. Amount of data to back up is unknown at this stage. By not using iSCSI, we saved two physical network ports in each server because they don't have to be dedicated to iSCSI.

[lack of time, didn't manage to finish write down configuration of this].

3.3 vSphere vCenter design

- Windows 2012 domain Mars
 - SQL and vCenter running in separated VMs (not too much overhead, gives a little flexibility although
- ... **unfinished**

3.4 vSphere networking design

Addressed Requirements : R01, R02, R03, R05

Addressed Assumptions : A03, A04, A05, A10, A11

Addressed Constraints : C02, C06, C07, C08, C09, C10, C11, C12, C13

Unknown networking infrastructure is already installed and operated in premises. The only information we have is that it supports Ethernet only, with “plenty of 10Gbit ports available” (no shortage has been mentioned site-wise, we expect to have adequate number of ports available in each site) and “some 40Gbit ports”.

Distributed switches will be used on VMware level for all types of connections except VMKernel host definition (for this, standard switch will be used). Following types of traffic have been identified :

Management traffic

Fault Tolerance Traffic

VSAN1, VSAN2

virtual machine networks 1, virtual machine networks 2 for redundancy

vMotion

Management and FT traffic will be shared on single physical port or adapter – while FT can create lot of traffic, management is not bandwidth hungry (except when deploying new VM...). We have 6 physical 10Gbit ports in each host which perfectly blends with requirements for NW traffic. Although not critically necessary, we will use Network IO Control to set “priorities” (shares) of particular network traffic types sharing each physical adapter.

Following combinations have been created :

NIC1 Port1: Management traffic && Fault-Tolerance

NIC1 Port2 : VSAN1

NIC2 Port1 : virtual machines traffic1 && vMotion

NIC2 Port2 : Fault-Tolerance && Management traffic (redundant connection)

NIC3 Port1 : VSAN2 (redundant connection)

NIC3 Port2 : virtual machines traffic2 && vMotion (redundant connection)

This way, no single physical adapter will host primary and redundant connection for identical network traffic, rendering it unusable in case of failure. All physical adapters will be connected to separate physical switches.