S9670 VIRTUAL DESKTOPS BY DAY, COMPUTATIONAL WORKLOADS BY NIGHT -AN EXAMPLE INFRASTRUCTURE



Shailesh Deshmukh Senior Solution Architect

Konstantin Cvetanov Senior Solution Architect

Eric Kana Senior Solution Architect

GPU Technology Conference 2019

AGENDA

- What We Will Discuss
- Benefits of VDI
- Computation Defined and Context
- Dual-Use and Workflow Scenarios
- Operational Challenges
- Solution Options
- Reference Architecture
- Demonstration
- Summary

WHAT WE WILL DISCUSS

A practical approach to configure intervals of VDI and Computational Resources on a daily basis - in an environment primarily designed for VDI - using commonly available tools.

More about perspective than technology

BENEFITS OF VIRTUAL DESKTOP INFRASTRUCTURE

- Enable flexible workflow scenarios
- Utilize centralized, shared, and protected storage
- Enable intellectual property protection
- Provide flexibility in configuration
- Enable user/workforce mobility
- Widely supported GPU acceleration

What you planned the system to do.

COMPUTATIONAL SPECTRUM

Additive Scale of Requirements



System Complexity

•

WHY DUAL USE?

- Cost and/or space savings
- Variable usage trends/rates
- Desire for on-prem elasticity
- Unpredictable user community
- Provide more workflow options to more users
- Effective cost justification (capital/operational) *Make best use of available resources*

SCENARIO CONSIDERATIONS FOR DUAL USE

- Creative Studio Artists go home during late hours
- Architecture Firm Engineers/Designers work daylight hours
- University/College Lower utilization during summer sessions
- Financial Services Firm Lower utilization when markets are closed
- Gov't Agency Multiple programs, duplicate (idle) resources

Primary goal is user experience

WORKFLOW CONSIDERATIONS FOR DUAL USE

- Creative Studio Create during day / Render by Night
- Architecture Firm Design during day / Render-Compute by Night
- University/College Sell cycles or run experiments during Summer
- Financial Services Firm Traders by day / Numerical analysis by night
- Gov't Agency Analysis work by Day / Image processing at Night

Get creative with workflow overlap

OPERATIONAL CHALLENGES

- What to *do* with our user VMs?
- How do we best *provision* user VMs?
- How do we *monitor* utilization?
- How do we *orchestrate* user VM state, migration, and timing?
- How do we *manage* compute jobs, and be ready for user VM restart?
- How will users **be productive** in a scheduled environment?

Manage Users, balanced with Compute Productivity

VECTORS FOR SUCCESS

- User policies reboot per day or week
- Single precision math jobs
- Single GPU compute jobs
- Jobs that may be coalesced
- Excess capacity
- Stakeholder buy-in
- Skilled admin staff



COMMON VDI INFRASTRUCTURE ASSETS

- Hypervisor(s) vSphere, AHV, RHVH, XenServer
- vGPU Software
- Compute cluster of nodes (chassis)
- CPUs, GPUs, Storage, Network Assets
- Monitoring Tools
- Orchestration / Layering Tools
- Containers
- Job Schedulers

Many common building blocks available



SOLUTION VECTORS

- Shut down (all users) and swap (in all the compute)
- Shut down (some users) and swap in (some) compute
- Migrate/degrade (users) to fewer hosts, swap (in some/all) compute
- Shut down (all users) and reprovision (to bare metal) nodes
- Keep all users intact; initiate a cycle harvester
- Some mixture of the above
- Other options...

GOAL = Use common and available tools

OPTION 1: SHUT DOWN / SWAP IN



- Shut Down User Pool
- Spin up compute Pool
- Run Scheduled Jobs
- Spin down compute Pool
- Restart User Pool

(Partial Shutdown also applies)

ARCHITECTURE DIAGRAM



Control Resources

Compute Resources

SLURM WORKLOAD MANAGER

"Slurm is an open source, fault-tolerant, and highly scalable cluster management and job scheduling system for large and small Linux clusters."

Source: https://slurm.schedmd.com/overview.html

Components:

- Centralized Manager: **slurmctld** monitors resources and work
- Compute Node daemon: slurmd waits for and executes work, returns work status

In this example:

- Slurm-ctrl = cluster controller VM
- Compute[01-07] = compute VMs (nodes)

ANATOMY OF A COMPUTE VM

- Ubuntu 16.04/18.04
- Docker, nv-docker, Anaconda, Python3-pip, ipythonnotebook
- vGPU 7.1
- CUDA 10, toolkit, and samples
- SLURM
- VMware VIEW agent
- DHCP per Active Directory DNS
- Packaged as a VM template

COMPUTE PARTITION ORGANIZATION



Template Resource Partitions (SLURM)

SLURM COMPUTE PARTITION CONFIG

COMPUTE NODES Partitions
NodeName=Compute[01-07] CPUs=2 Sockets=2 CoresPerSocket=8 RealMemory=14336 State=L
PartitionName=T4x16Q Nodes=Compute[01-04] Default=YES MaxTime=48:00:00 State=UP
PartitionName=V100x32Q Nodes=Compute[05] Default=YES MaxTime=48:00:00 State=UP
PartitionName=RTXx24Q Nodes=Compute[06-07] Default=YES MaxTime=48:00:00 State=UP

/etc/slurm/slurm.conf

PARTITION	AVAIL	TIMELIMIT	NODES	STATE	NODELIST
T4x16Q	up	2-00:00:00	4	idle	Compute[01-04]
V100x32Q	up	2-00:00:00	1	idle	Compute05
RTXx24Q*	up	2-00:00:00	2	idle	Compute[06-07]
	_				-

sinfo output

Linux VM Templates mapped to Compute Partitions

OPERATIONAL TIMELINE

VDI State	npute	Compute State	npute	VDI State	
6 VMs (Linked-clon	ies)	4 x T4-16Q	e Con	6 VMs (Linked-clones))
Windows 10	Star	1 x V100-32Q	cuat	Windows 10	
Non-persistent VMs	s 🔁	2 x RTXx24Q	Evad	Non-persistent VMs	
T4-8Q vDWS Profile	Evacuate V	7 compute VMs	Start VDI /	T4-8Q vDWS Profiles	
VDI State t1	tź	Compute State	t	VDI State :3	time
6 am	Midn	ight	6 8	am	

VCENTER INTERVAL SCHEDULING

VDI Interval:

Summary Monitor	Configur	e	Permissions	Datastore	s Ne	etworks							
 Settings vApp options More 	NEW	/ SCI	HEDULED TASK	 EDIT 	RUN	REMOVE							
Alarm Definitions			Scheduled Task		Ŧ	Schedule	٣	Last run	Ŧ	Last run result	Ψ	Next run	Ŧ
Scheduled Tasks Policies	0	>	VDI-By-Day02	- Power On		Daily		02/14/2019, 6:05:00	MA	✓ Completed		2/15/2019, 6	:05:00
VMware EVC	0	>	VDI-By-Day02 -	- Shut Down	Guest OS	Daily		02/13/2019, 8:25:00	PM	✓ Completed		2/14/2019, 8	25:00
vin haldware												1 - 2 of 2 Schedu	iled Tasks

Compute Interval:

ettings vApp options	NEV	SCHED	ULED TASK - EDIT RUN	REMOVE								
Alarm Definitions		Sc	heduled Task	▼ Sched	ule	Υ.	Last run	Ŧ	Last run result	Ψ	Next run	Ŧ
Scheduled Tasks		> ~	ompute02 - Power On	Daily			02/12/2010 9:25:00 DM		Completed		2/14/2010	2.25.00
icies	0	/ 11	mpateoz - Fower on	Daily			02/13/2019, 0.33.00 PM		♥ completed		2/14/2015, 0	
d Hardware	0	> cc	ompute02 - Shut Down Guest OS	Daily			02/14/2019, 6:00:00 AM		✓ Completed		2/15/2019, 6	00:0

SHUT DOWN / SWAP IN - HARDWARE

Component	Name
GPU	Tesla T4, V100, P40, RTX
Chassis	Supermicro 4029GP, Dell R740, HPDL380 Gen9
Storage	FA-M20R2 (Pure Storage)
Network	CISCO 10G
Endpoints	Various

SHUT DOWN / SWAP IN - SOFTWARE

Component	Name					
Hypervisor	vSphere 6.7u1					
Hypervisor Manager	vCenter 6.7					
Job Scheduler	Slurm 17.11.12					
Interval Scheduler	vCenter 6.7					
VDI Guest o/s	Windows 10					
Compute Guest o/s	Ubuntu 16.04					
NVIDIA vGPU Software	vGPU 7.1					

 ✓ vSphere - View-Environment-T4 × ✓ ✓ ✓ ✓ ▲ Not secure https://sa 	NVIDIA Enterprise × + grid-vcenter.sagrid.local/ui/#?extensionId=v	sphere.core.inventory.serverObjectViews	Extension&objectId=urn:vmomi:ClusterComp	uteResource:domain-c640:6318f9e2-f571-4396	5-a3fb-67624fa Q 🛧 🚺 Q 👌	-	o X
Apps Nvidia Personal Proj	ects 📙 Studies 📒 DEMOS 🤕 Certified	Servers N 🤕 NVIDIA Virtual GPU 🤕	Content Library 🝐 Virtual GPU Present 🝐	VirtualGPU POC tra 🔞 GRID Software: Bug	 MVIDIA Enterprise Amazon Web S C O 	ervic 📄 Lecture Notes - VN	»
Image: SAGRID-TEAM	Summary Monitor Configure	ACTIONS - Permissions Hosts VMs Da	tastores Networks Updates				- 60 22 GU
 Management Quadro-RTX SAGRID-Dell-R730 SAGRID-HP-DL380G9 View-Environment-M60 View-Environment-P100 View-Environment-P4 	Total Processors: Total vMotion Migrations:	0				Used: 11.77 GHz Capa Memory Free Used: 48.12 GB Capacit Storage F Used: 15.2 TB Capacit	city: 72 GHz e: 209.55 GB /: 255.67 GB free: 8.01 TB ity: 23.21 TB
 IView-Environment-P40 View-Environment-T4 10.31.230.45 Compute01 Compute02 	Related Objects Datacenter	SAGRID-TEAM	^	Cluster Consumers			~
Compute02 Compute03 Compute04 Compu	Tags Assigned Tag	Category	Description	Attribute	Value		*
Windows10VDI-2 Windows10VDI-3 Windows10VDI-4 Windows10VDI-5	4			Edit		No items to d	splay
 Windows10VDI-6 View-Environment-V100 	Assign Remove Cluster Resources		No items to display	Update Manager Compliance Precheck Remediation State	 Remediation Status Unknown (lass 	it check n/a)	
Recent Tasks Alarms							*
Task Name v Target	✓ Status	~ Initiator	✓ Queued For	✓ Start Time ↓	 Completion Time 	~ Server	~
Update option values	31.230.45 ✓ Complete	ed SAGRID\nvadmi	n undefined	02/27/2019, 8:17:07 AM	02/27/2019, 8:17:07 AM	sagrid-vcenter.sagrid.local	
Reconfigure scheduled task	ompute04 V Complete	d SAGRID\nvadmi	n undefined	02/27/2019, 8:14:57 AM	02/27/2019, 8:14:57 AM	sagrid-vcenter.sagrid.local	
Reconfigure scheduled task	ompute04	d SAGRID\nvadmi	n undefined	02/27/2019, 8:14:41 AM	02/27/2019, 8:14:41 AM	sagrid-vcenter.sagrid.local	
Reconfigure scheduled task	Complete	a SAGRID'invadmi	n undenned	02/27/2019, 8:14:28 AM	02/2//2019, 8:14:28 AM	sagria-vcenter.sagria.local	More Tacks

All

🛃 nva	idmin@Con	npute01: ~					_		;	nvadmin@	Compute05	: ~					-		×
00:	T4-16Q (2:02.0	SM[0%]	FB[7%]	CL[00: V100I	-32Q 02:	02.0	SM[0%]	FB[7%]	CL[-MHz
<mark>PID</mark> 1766	GPU S O F	S SM% 8 0.0	FB% 0.0	MC% 0.0	EN% 0.0	DE% 0.0	PROC /usr	ESS /lib/x	org/	<mark>PID</mark> GE 17790	PUSSM% R0.0	FB% 0.0	MC% 0.0	EN% 0.0	DE% 0.0	PROCI /usr	<mark>ISS</mark> /lib/x	org/Xo	rg -c
e nva	dmin@Com	npute02: ~					_			🛃 nvadmin@	Compute06	~							_
00:	T4-16Q 0	02:02.0	SM[0%]	FB[7%]	CL[top – 11:3 Tasks: 258	9:15 up total,	1:02, 1 ru	2 users nning, 25	, load 7 sleep	avera	ge: 0.15 0 stopp	5, 0.5 ped,	7, 0.50 0 zoml) Die
<mark>PID</mark> 1823	GPU S O R	SM% 0.0	FB% 0.0	MC% 0.0	EN% 0.0	DE% 0.0	PROC /usr	ESS /lib/x	org/ <mark>H</mark>	%Cpu(s): KiB Mem : KiB Swap:	0.5 us, 16432396 998396	0.8 s total total	y, 0.2 n , 1412177 , 99839	1, 98.5 6 free, 6 free,	id, 707	0.0 wa, 096 useo 0 useo	0.0 1, 16 1. 153	hı, 0 03524 1 21856 a	.0 sı, ouff/ca avail N
e nva	dmin@Con	npute03: ~					_			PID USER 1558 root	. PR 20	NI 04	VIRT 051932 19	RES 4228 1	SHR S 9572 S	%CPU 9 0.7	MEM 1.2	TIM 0:17.	E+ COMN 52 java
00:	T4-16Q 0	2:02.0	SM[08]	FB[7용]	CL[14061 root 2573 nvad	30 min 20	10 3	666232 20 41904	5088 1 3804	9780 S 3096 R	0.7	1.2	0:12.	10 java 23 top
<mark>PID</mark> 1824	GPU S O F	5 SM% R 0.0	FB% 0.0	MC% 0.0	EN% 0.0	DE% 0.0	PROC /usr	<mark>ESS</mark> /lib/x	org/	₫ [®] nvadmin@	Compute07	·: ~							-
e nva	admin@Con	npute04: ~					_			00: RTX60	000-240 (PU S SM%)2:02.0 FB%	SM[MC%	EN%	DE %	0%] FB PROC	[Ess	7%] C	L[-1
00:	T4-16Q (02:02.0	SM[0%]	FB[7%]	CL[188 <mark>2 0</mark>	R 0.0	0.0	0.0	0.0	0.0	/usr	/lib/x	org/Xo	rg -co:
PID	GPU S	5 SM%	FB%	MC %	EN%	DE %	PROC	ESS		F1 <mark>Help</mark> H	2 <mark>Setup</mark>	F3 <mark>Sear</mark>	ch <mark>F4</mark> Fil	ter <mark>F5</mark>	Start	F6 <mark>Sto</mark>	0		F:
										Na Slur	m@slurm-ct	trl: ~/jobs						_	
										Termina Session	Sessions Servers	View X	server Tools server Tools Games Sess	Games Set C Q ions Viev	tings Ma I v Split	acros Help Y t MultiExe	c Tunneli	ng Packag	X es erver
F1 <mark>Hel</mark>	p F2Se	etup F3	Search	F4 <mark>Fil</mark>	ter F58	Start	F6 <mark>Sto</mark>	q		✓ Quic Sl	k connect urm@slur	m-ctrl	î :∼/jobs\$	≥ 5. /home/	moba	E 6. slurr	n@slurn	×	

ENVIRONMENT MONITORING

1 - 9 of 9 item:	
D/Compute Utilization	Memory Utilization
a 🔤 🥺 🛧 👀 🔄 🔍 🔯 🐟 😫 🕫 C 🗓 🟥	😻 💼 🔤 🍕 ヤ 💵 🐼 🔲 🔍 🔯 💩 🔯 🖉 🖉 🖬
Compute07-GRID RTX6000-24Q (EV) UtilizationI3D/Compute Utilization	Compute07-GRID RTX6000-24Q
H:1	•H:1
0	••••••••••••••••••••••••••••••••••••••

FUTURE NEEDS AND ASKS

- Multiple GPUs per VM limited availability today
- Dynamic vGPU assignment per Template provisioning
- Dynamic vGPU on live migration
- vGPU + GPU ECC + UVM + P2P supports relevant compute
- vGPU + GPU memory Page retirement
- VM snapshots and user sessions
- Storage optimizations
- Live migration integration exists today

IMPORTANT: VGPU VM DEPLOYMENT POLICY (VMWARE / CITRIX)

VMware vSphere Hypervisor (ESXi) by default uses a breadth-first allocation scheme for vGPU-enabled VMs; allocating new vGPU-enabled VMs on an available, least loaded physical GPU. We need to

change that ..

[root@dclvmhvesx45:~] cat /etc/vmware/config .encoding = "UTF-8" libdir = "/usr/lib/vmware" authd.proxy.nfc = "vmware-hostd:ha-nfc" authd.proxy.nfcssl = "vmware-hostd:ha-nfcssl" authd.proxy.vpxa-nfcssl = "vmware-vpxa:vpxa-nfcssl" authd.proxy.vpxa-nfc = "vmware-vpxa:vpxa-nfc" authd.fullpath = "/sbin/authd" vGPU.consolidation = "TRUE"



FINDINGS

- At least 1 vCenter VM powered on in a pool (20/80 best practice)
- Unify the storage for users and data both VDI and Linux
- Alert users when jobs don't start properly SLURM
- Care for permissions SLURM, containers, renderers, storage
- SLURM is very powerful and potentially complex understand it
- Manage user VDI logistics and operations
- Keep the UX paramount

S9670 VIRTUAL DESKTOPS BY DAY, COMPUTATIONAL WORKLOADS BY NIGHT -AN EXAMPLE INFRASTRUCTURE



Shailesh Deshmukh Senior Solution Architect

Konstantin Cvetanov Senior Solution Architect

Eric Kana Senior Solution Architect

GPU Technology Conference 2019