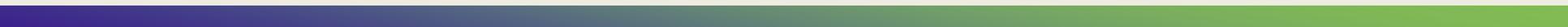




What's New with aCloud 5.8.7R1 & aCMP 5.8.6R1

Cheney Hu
Sangfor IMD



An Overview of What's New

Stretched Cluster

Sangfor's **active-active** data center solution based on aCloud, able to achieve zero data loss and near-zero downtime in the event of a site failure. The distance between 2 sites is restricted.

GPU support

NVIDIA GPU, pass-through and vGPU. For graphics-intensive and deep learning workloads

Reporting

Exporting reports of resource usage status of specified VMs and hosts

“Seed backup”

Exporting data to portable disk and transport it to secondary site to save time for first time full backup

...

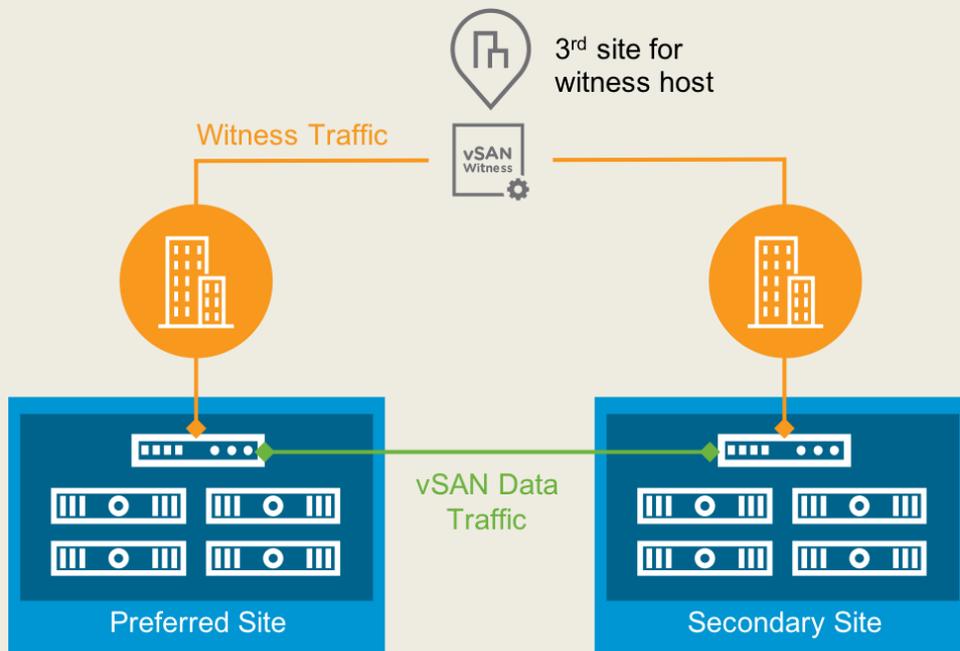


Stretched Cluster

Stretched Cluster (aCloud 5.8.7R1)

What is it?

Stretched cluster extends the aCloud cluster from a single data site to two sites for a higher level of availability and intersite load balancing. Stretched clusters are typically deployed in environments where the distance between data centers is limited, such as metropolitan or campus environments.



VMware vSAN stretched cluster

Stretched Cluster (aCloud 5.8.7R1)

Scenarios

You can use stretched clusters to manage **planned maintenance** and **avoid disaster** scenarios, because maintenance or loss of one site does not affect the overall operation of the cluster. In a stretched cluster configuration, both data sites are active sites. If either site fails, aSAN uses the storage on the other site. aSV HA restarts any VM that must be restarted on the remaining active site.

Values

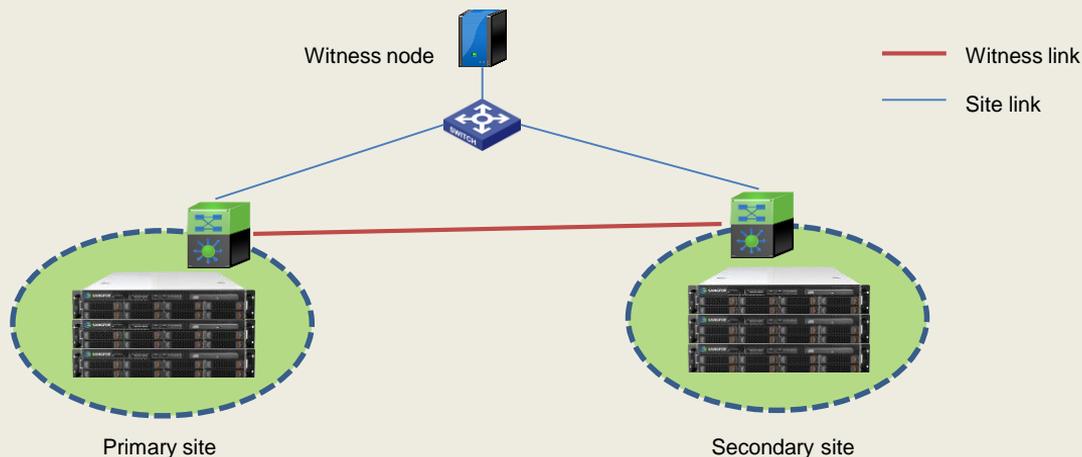
RPO = 0

RTO = minutes (HA reboot time, apps that support A-A can achieve zero downtime.)

Key words:

Zero data loss, minimal business downtime, starts small with 4 nodes

Sangfor aCloud Stretched Cluster



Site link:

- 10Gbps bare fiber is recommended, $RTT \leq 1\text{ms}$
- 1Gbps is supported for small deployment (4-6 nodes)

Witness link:

- Recommended $RTT \leq 1\text{ms}$ (fluctuation within 5ms is accepted)
- 100Mbps is recommended

Deployment scale:

- Only 1 stretched virtual volume is supported by 1 cluster
- Starts with 4 nodes, maximum 24 nodes

Supported scenarios

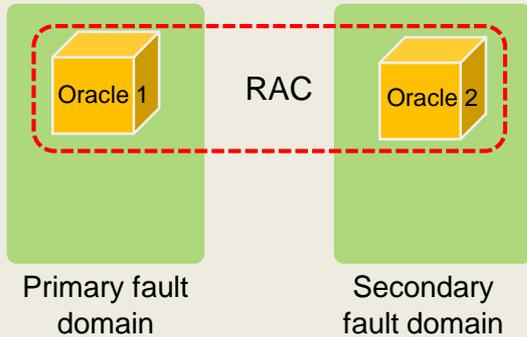
- Different floors inside the same building
- 2 adjacent buildings
- Inside the same campus

Scenarios to be supported

- 3 centers across 2 cities (unified mgmt. by aCMP, by the end of 2019)
- Metropolitan active-active (roughly 2019)

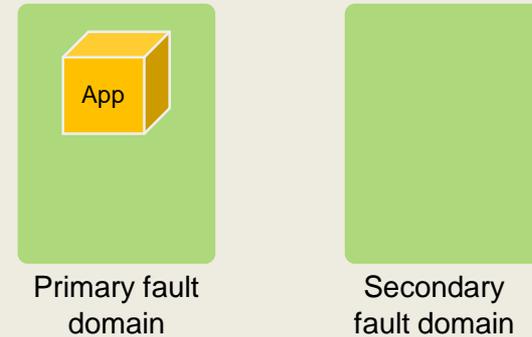
Business A-A and Data A-A

Business Active-active



- Business is running concurrently at 2 sites
- When a site fails, business is switched over to the other site by load balancer
- RPO=0, RTO≈0
- Business applications must support active-active mode, typical example is Oracle RAC
- Higher complexity

Data Active-active



- Business is running at one of the fault domains
- When primary site fails, App VM is rebooted at the other site by HA
- RPO=0, RTO=minutes
- No requirements for business apps
- Lower complexity

Licensing for Stretched Cluster

How to License

- Additional license called aSC, aSV and aSAN are prerequisite
- Licensed by the number of physical CPU sockets

Scenario 1

New construction, 1 stretched cluster with 20 physical CPU sockets, 1 traditional cluster with 40 CPU sockets, totally 60 sockets

Cluster Type	Licenses for aSV &aSAN	Licenses for aSC
Stretched cluster	20	20
Traditional cluster	40	0
Total	60	20

Scenario 2

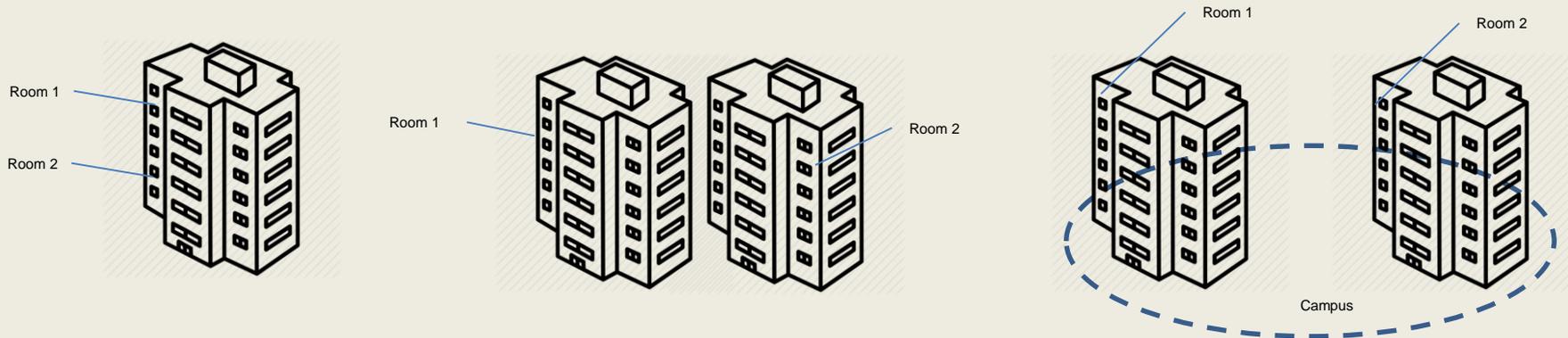
Construction in phases, primary DC at the 1st phase with 20 CPUs, add another 20 CPUs at the 2nd phase in secondary site to build a stretched cluster

Cluster Type	Licenses for aSV &aSAN	Licenses for aSC
Phase 1	20	0
Phase 2	20	40
Total	40	40

A Typical Use Case of aCloud Stretched Cluster

Customer's requirements

- Customer has 2 server rooms at floor 1 and floor 2, they expect to deploy nodes at the 2 server rooms and build a stretched cluster across to handle availability in the case of server room failure
- Normally, 2 server rooms are running different business workloads simultaneously
- The devices in the server room at floor 2 are newer and more reliable than the ones at floor1, customer expects to put important workloads in the new server room



Polling 1:

Do you have any requirement for active-active data centre solution?

- A. Yes, the server rooms are located on different floors in the same building
- B. Yes, the server rooms are located in 2 buildings side by side
- C. Yes, the server rooms are located in 2 buildings in the same campus
- D. Yes, 2 sites are located in the same city but dozens of kilometres away from each other
- E. No, never had requirement for such a high level protection



GPU Support

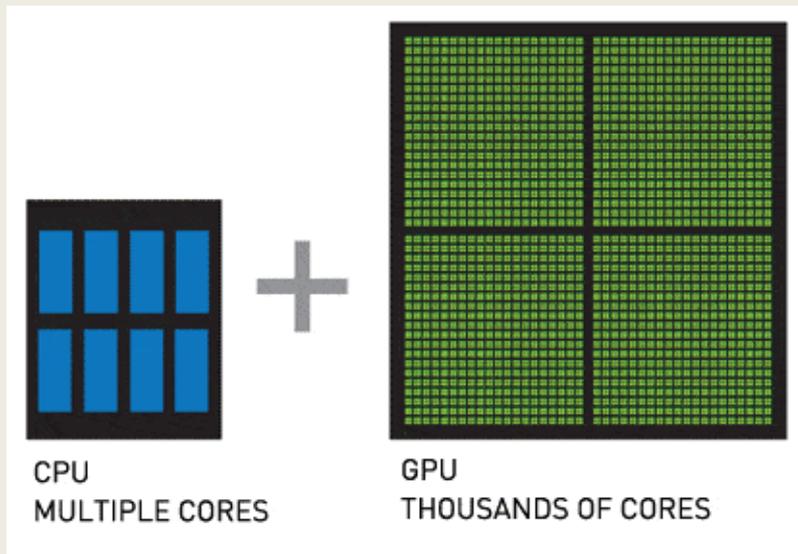
GPU Support (aCloud 5.8.7R1)

What is GPU?

A graphics processing unit (GPU) is a specialized electronic circuit designed to rapidly manipulate and alter memory to accelerate the creation of images in a frame buffer intended for output to a display device.

Characteristics of GPU Computing

High concurrency, high density, massive and parallel



GPU Support (aCloud 5.8.7R1)

Applications that require GPU

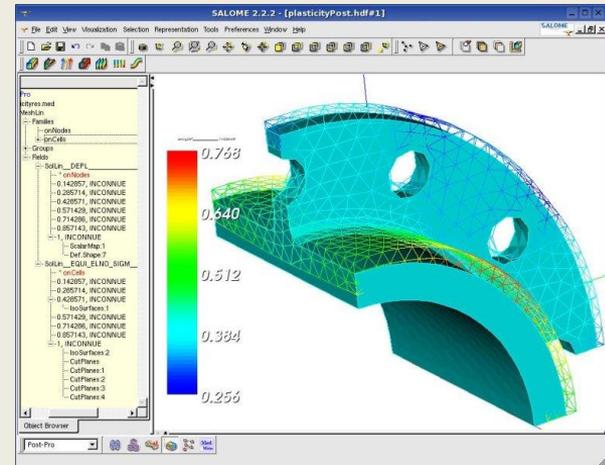
GIS (geographic information system), graphics rendering system, deep learning system, computer-aided engineering, etc

Scenarios

Video transcoding and streaming, graphics rendering, labs for deep learning and big data applications

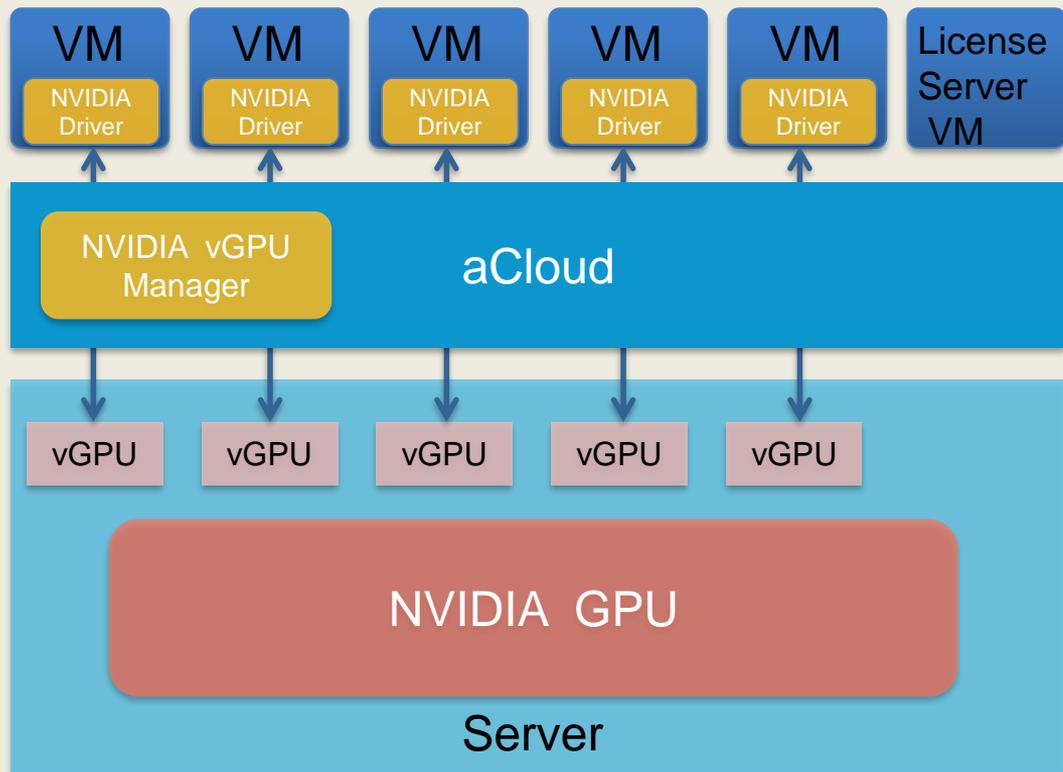


GIS



CAE

aCloud vGPU Solution



Logically partition a physical GPU into multiple vGPUs for multiple VMs to use

- Server hardware must be installed with NVIDIA GPU cards that support vGPU
- NVIDIA vGPU manager is embedded in the hypervisor, it takes charge of virtualizing physical GPU into vGPUs
- GPU drivers must be installed on VMs

GPU Support (aCloud 5.8.7R1)

Pain points

- Existing applications supported by physical GPUs need to be supported by virtualization platform in virtualization transformation
- When GPUs are deployed in a large scale, it's very difficult to match GPU resource with business requirements, thus leading to unbalanced utilization and resource waste
- GPU management in large scale environment could be very complicated

Values of aCloud vGPU solution

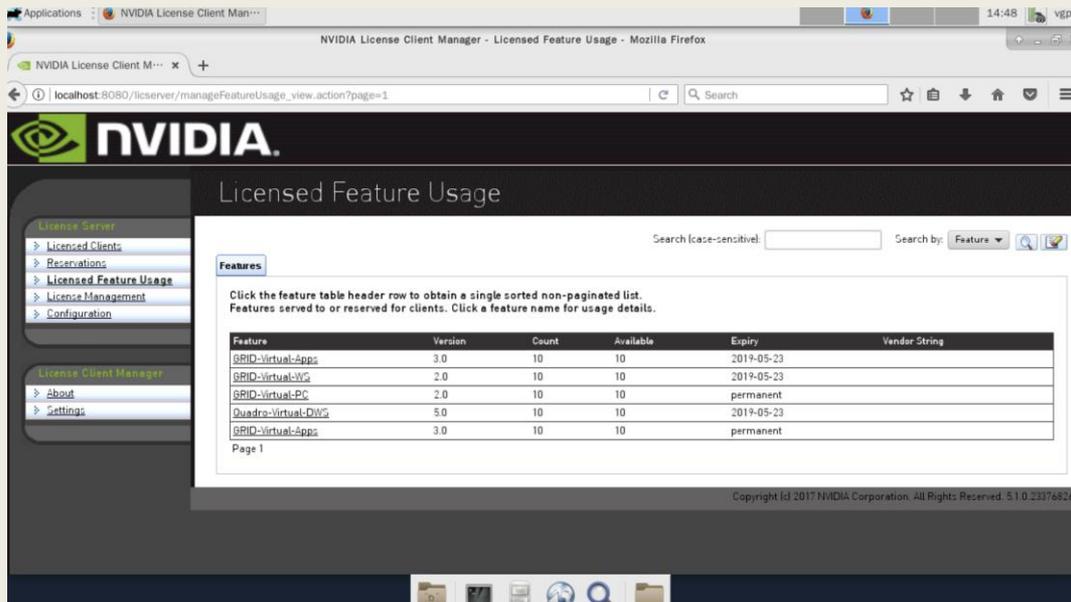
- ✓ Visualized GPU resource utilization, simplified deployment & management
- ✓ Pooling of GPU resource for multiple VMs to maximize GPU utilization

vGPU v.s. Physical GPU (In Large Scale)

	Physical servers + GPU	Virtualization + vGPU
Batch deployment	Bad	Good
Resource utilization	Low	High
Mgmt.	Per server	Unified and visualized
GPU Specs	Fixed per card	Flexible configuration of GPU memory based on needs

Licensing for vGPU

1. vGPU support requires license from NVIDIA, it's licensed by the the number of concurrent vGPU VMs
2. GPU pass-through requires no license from NVIDIA
3. No license from Sangfor for both modes



The screenshot shows the NVIDIA License Client Manager web interface. The main content area is titled "Licensed Feature Usage" and contains a table of features. The table has the following data:

Feature	Version	Count	Available	Expiry	Vendor String
SRID-Virtual-Apps	3.0	10	10	2019-05-23	
SRID-Virtual-WG	2.0	10	10	2019-05-23	
SRID-Virtual-PC	2.0	10	10	permanent	
Quadro-Virtual-DWG	5.0	10	10	2019-05-23	
SRID-Virtual-Apps	3.0	10	10	permanent	

Page 1

Licensing procedure

1. Upload licensing server template to aCloud
2. Configure IP and import license file
3. Input licensing server's IP on the GPU driver of vGPU VM

There are 3 types of vGPU licenses from NVIDIA: vPC, vDWS and EDU. vDWS has better performance than vPC, EDU is a promotion for education sector, it shares the same functionalities with vDWS, but has a lower price as vPC, only schools have the right to purchase it

Differences between vGPU Solutions on aCloud and aDesk

What aCloud has but aDesk doesn't:

1. Support GPU for Linux
2. Support launching VM console directly from the web UI
3. Support vGPU mode for NVIDIA P40
4. Optimization on UI details

Some Caveats

1. Not every server can support GPU installation, **our previous aServer models can't support GPU**. If customer has existing GPU servers, reusing those servers is a good option; Dedicated GPU aServer models (aServer-G350 or G650) are needed for new construction scenarios.
2. Only 1 graphics card (passthrough or vGPU) can be attached to a VM, this is limited by NVIDIA technology

Polling 2:

Are you looking for GPU solutions and for what workloads?

- A. Yes, graphic rendering
- B. Yes, deep learning and data analysis
- C. Yes, video transcoding
- D. No, never had such requirement

Seed Backup (aCMP 5.8.6R1)

Introduction

Seed backup is a complementary feature for off-site DR scenario. User can use this feature to make a seed file based on selected VM's backup and export it to external storage (USB flash drive or portable disk), the storage device is then transported to the secondary site and VM will have related backup on the secondary site after the seed file is imported. In a word, seed backup is a way to replace network-based data transfer with physical transportation.

Targeted customer

Customer has requirement for DR, but the link bandwidth between the 2 sites is very narrow (due to limitations on budget or physical environment).

It's more suitable for scenarios with new secondary site construction because first-time full backup can be very time-consuming.



Targeted customer

Reporting & Others

Reporting (aCMP 5.8.6R1)

Reports can be generated and exported based on selected metrics (CPU, RAM, IOPS, latency and throughput) of specified VMs and hosts

Support Importing VMs in OVF, MF and VMDK formats (aCloud 5.8.7R1)

Italian keyboard support (aCloud 5.8.7R1)

VM scheduling: affinity & anti-affinity (aCloud 5.8.7R1)

Affinity and anti-affinity rules allow you to spread a group of virtual machines across different hosts or keep a group of virtual machines on a particular host

Support CPU and memory automated hot add for Windows Server 2012, 2012R1 and 2016 VMs (aCloud 5.8.7R1)

...

Q&A



Thank You

