

Important Performance Considerations When Moving SAS® to a Public Cloud

Margaret Crevar, SAS Institute Inc.

Updated: 29MAR2018

ABSTRACT

Any hardware infrastructure that is chosen by a SAS® customer to run their SAS® applications requires the following:

- a good understanding of all layers and components of the SAS® infrastructure
- an administrator to configure and manage the infrastructure
- the ability to meet SAS requirements—not to just run the software but to also allow it to perform as optimally as possible

This paper will talk about important performance considerations for SAS® 9 (both SAS® Foundation and SAS® Grid Manager) and SAS® Viya® when hosted in any of the available public clouds (Amazon AWS, Microsoft Azure, and Google Cloud, to name a few). It will also give guidance on how to configure the cloud infrastructure to get the best performance with SAS.

INTRODUCTION

Many SAS customers are making the decision to move their current SAS applications from their on-premises data centers to a public cloud. The hype around public clouds portrays this as a very simple task that saves SAS customers a lot of money.

It should be noted that the information discussed in this paper is based on what is available from the public clouds and SAS experience with the public clouds at the time of the writing of this paper. Public cloud offerings are constantly changing. Therefore, it is in your best interest to understand the rationale used in the selection process and to consider what was done as a point-in-time design.

However, there is a lot of planning needed, and depending on the requirements from the SAS customer, the price might not be cheaper than on-premises hosting. The chief reason is that if IO throughput is crucial to the success of your SAS applications in the public cloud, you might need to provision more cores and disk space capacity to ensure the success of SAS public cloud deployments. I will explain the reasons for this.

BEFORE YOU START

As mentioned in the introduction, a good understanding of the SAS application workload requirements, along with the required hardware infrastructure required to meet the service objectives (SLAs), specifically time to task complete, is crucial. For existing SAS customers, the following questions help guide that examination:

- Are there SAS jobs that need to execute within a certain time frame? Are you expecting your SAS jobs to execute in the same time as—or faster than—these jobs are currently finishing in their existing data center? If so, a determination of the IO throughput for each file system being used must be made. It must be determined if this same IO throughput can be achieved in the public cloud.
- Where is the source data for the SAS jobs located? Does this data already reside in the public cloud of choice? If not, the length of sourcing data to the cloud space that SAS is executing in needs to be determined. This added time will affect the SLA of the jobs that consume off-cloud data.

- Is the customer's IT staff willing to do stand-up authentication in the public cloud?
- What security is needed for the data and/or SAS code?

The results of the above questions and fact-finding need to be fully understood so that the correct hardware and storage is selected from the available public cloud offerings. There are many different hardware and storage types. Some are hardware equipped for the heavy analytical and large sequential IO that SAS 9 does. And others are better equipped for the in-memory needs of SAS® Viya. It is important to understand the workload profile needs of the customer's SAS application(s) to ensure that correct hardware and storage selections (cloud server and storage types) are made for the best performance. Please note to get the best achievable performance, the least expensive hardware and storage types from the public cloud offerings might not be suitable. The customer might also have to have a stand up server and storage instances with more physical cores than required for compute needs and/or set up more storage capacity than the initial sizes needed to acquire the maximum IO bandwidth available for their SAS application(s).

Now let's talk about what needs to be considered to ensure that you can configure the hardware infrastructure in the public cloud to perform as optimally as possible. These things include the following:

- what server instance type to use
- what storage type (for both persistent and nonpersistent storage) to use
- if deploying SAS® Grid Manager, what shared file system to use
- where to place temporary (SASWORK/UTILLOC and CAS_Disk_Cache) and permanent (SASDATA) data to be used by SAS
- where to place the SAS clients that will be used
- where to place authentication tools
- are high availability and security required?

WHAT INSTANCE TYPE TO USE

In a SAS 9 and SAS Viya infrastructure, there are several SAS server types and uses. Each has different and specific requirements for CPU, IO throughput, and memory provisioning. We will list each SAS server type and discuss its provisioning requirements. Please remember that the most public cloud instances list CPUs as virtual CPU(s). These CPUs are hyper-threaded (two threads per CPU core). SAS does not recommend the use of hyper-threads for SAS 9, so you need to understand that the cloud instance vCPU count must be divided by 2 to obtain the number of usable physical cores for SAS.

You might have to use an instance with more physical cores in it than your workload requires. This is because a higher CPU count machine might be required to obtain a dedicated network interface card (NIC) of sufficient bandwidth to maximize IO to and from off-server cloud storage. Server instance types come in set models, with set CPU counts, set NIC card installation, RAM, and so on. To get a higher bandwidth NIC card, you might have to upgrade to a larger server type (with more CPU and RAM than needed for compute). When calculating NIC card capacity to drive storage IO bandwidth, be aware that sharing a NIC card in a cloud server host with other multi-tenant applications residing on the same physical server, might result in inferior performance that occurs randomly. This is especially true for virtualized cloud host instances.

When setting up the instances, please make sure that all the instances are in the same cloud space (for example, region, area zone, and placement group). Failure to do so results in an additional WAN connection that severely impacts performance.

SAS 9 SERVERS

These are general guidelines. We advise that you have a detailed workload assessment done to determine which hardware is needed to support the usage of SAS 9 in the public cloud.

SAS® Compute Tier with SAS® Grid Node

This server needs fast CPUs for processing data, a minimum of 8 GB of RAM per physical core, and robust server host local IO storage and throughput (especially to SAS WORK and SAS UTILLOC).

- Amazon
 - I3 series – This series can provide NVMe (Non-Volatile Memory Express) SSD drives that can be striped together in a RAID0 configuration to amplify bandwidth for SAS WORK and UTILLOC file systems. This series provides the best IO bandwidth for SAS IO needs for SASWORK and UTILLOC file systems.
- MS Azure
 - DSv2series
 - GS series
 - L series (only available in limited areas)
- Google
 - N1-standard series

Shared File System Storage Required for SAS Grid

These servers need robust IO throughput to the permanent storage they access. The instances will also need a minimum of 8 GB of RAM per physical core.

- Amazon
 - R4.8xlarge – IO throughput to permanent storage is 875 MB/second per server node
 - R4.16xlarge – IO throughput to permanent storage is 1,750 MB/second per server node
- MS Azure
 - Hm series
- Google
 - N1-standard series

SAS Mid-Tier and Metadata Servers

These servers do not require compute-intensive resources and lesser IO bandwidth, but do require access to more memory than the SAS Compute Tiers. The recommendation is a minimum of 24 GB of physical RAM or 8 GB of physical RAM per physical core—whichever is larger.

- Amazon
 - R4 series (for mid-tier servers)
 - M series (for metadata servers)
 - X1e series
- MS Azure
 - DS series
- Google
 - N1-highmem series

SAS VIYA SERVERS

For this paper, we are going to list several machines to use for the most robust SAS Viya offering. Some of the SAS Viya offerings (single, non-distributed node going against small data from a source other than Hadoop) will not need all the servers listed below. These are general guidelines. We advise that you have a detailed workload assessment done to determine which hardware is needed to support the usage of SAS Viya in the public cloud.

You might have to deploy an instance with more physical cores in it than your workload requires, to upgrade to an instance type that provides a dedicated NIC (ethernet) card. Sharing a NIC card with multi-tenant applications that are not in SAS on the same physical hardware might result in inferior performance that occurs randomly.

CAS Nodes – Minimum of Three

These servers require fast CPUs for processing data, enough physical RAM to hold all the data files to be analyzed by all the concurrent SAS Viya users, and robust IO throughput (especially to CAS_Disk_Cache). If you are not sure how much data will be accessed at any given time, but know you SAS users will be accessing files in the 100s of gigabytes in size, we recommend 64 GB of RAM per physical core (more if you have large files that will be in memory).

- Amazon
 - I3 series– the primary reason is the high internal IO bandwidth from striped NVMe SSD drives for SAS WORK and UTILLOC file systems.
- MS Azure
 - Hm series
- GS series Google
 - N1-highmem series

MicroServices Node

These servers do not require high computational speed or power. To run all SAS® Visual Analytics products with SAS Viya, you will need at least 96 GB of RAM in your MicroServices Node.

- Amazon
 - R4 or X1e series
- MS Azure
 - Hm series
- Google
 - N1-highmem series

SAS® Programming Run-Time Node

This server will run SAS® 9 code per the application. This server needs fast CPUs for processing data, at least 16 GB of RAM per physical core, and robust IO throughput (especially to SAS WORK and SAS UTILLOC).

- Amazon
 - I3 series– The primary reason is the high internal IO bandwidth from striped NVMe SSD drives for SAS WORK and UTILLOC file systems.
- MS Azure
 - Hm series
- Google
 - N1-standard series

Hadoop Nodes

If you are using Hadoop file systems for data access (consume or store), then the following server instance types are preferred:

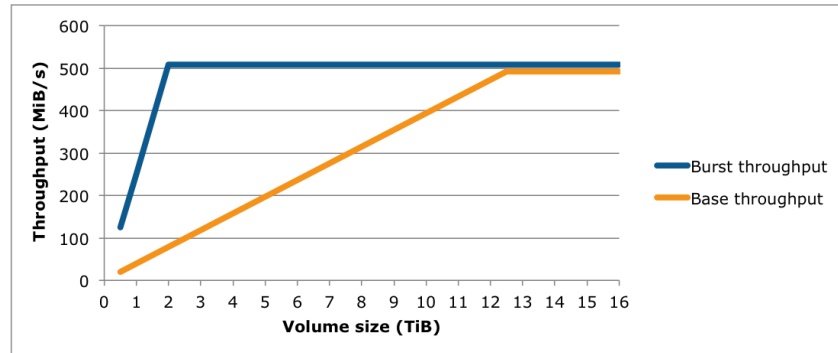
- Amazon
 - D2 instances
- MS Azure
 - DS or DSx_v2 series
- Google
 - N1-standard series

WHICH STORAGE TYPE TO USE

The public clouds have multiple types of permanent and temporary data storage offerings for use with SAS. Discussed below are findings from general field experience and lab testing for SAS in public cloud space.

Permanent SAS Data Storage Will Persist through a Reboot/Restart

- Amazon
 - Elastic Block Storage (EBS) – Our recommendation is to stripe together a minimum of 4 EBS volumes together for IO bandwidth aggregation.
 - EBS ST1 (large block IO) Storage – Designed for large block sequential IO. A 12.5TB volume can sustain 500 MB/second. (This is a published value from Amazon.) If your volume size is less than this, you will only get 500 MB/second total bandwidth during your burst window. Please see the [EBS Storage guide](#) on the Amazon website for more details.



- Other EBS storage types (like general purpose SSD (gp2), sustained IOPS SSD (io1), and cold storage (sc1) should not be used for permanent SAS 9 data files.
- Elastic File System (EFS) storage – you can use this type of storage, but you need to understand that the maximum IO throughput for each EFS file system is only 250 MB per second per node. You can have multiple EFS files systems per node to overcome the 250 MB per second per node limitation, but you need to understand that there is a limitation of the single NIC card in each AWS EC2 instance.

In addition, if you are looking to place SAS Home and SAS Configuration files on EFS, you will need to add the following statement to your sasv9.cfg file.

```
"-filelocks <HOME> none" where <HOME> is the directory containing the SAS Home and SAS Configuration files.
```

File locks need to be turned off is because EFS has a hard limitation of 87 locks per file and SAS manages locks on various files in SAS Home.
- MS Azure – P20 or P30 Premium Storage Disk Type. You should stripe at least 4 (preferably 8) volumes together per instance to get the maximal IO needed.
- Google – Use the storage with the highest IO throughput. You should stripe at least 4 preferably 8 together to get the IO throughput for each instance.

Temporary SAS Data Storage – Will Not Persist through a Reboot/Restart

Temporary storage is most commonly used for SAS WORK, SAS UTILLOC, and CAS_Disk_Cache (with SAS Viya) because this data does not need to persist through reboots and restarts. When placed on striped SSDs, SAS WORK and SAS UTILLOC should share a single file system.

- Amazon
 - Internal SSD devices striped together with RAID0.) The above I3 instances have NVMe SSD (ephemeral) devices with high bandwidth, low latency, and sequential IO, all of which are ideal for temporary SAS data.
- MS Azure

- Internal SSD devices striped together with RAID0.
- Google
 - Internal SSD devices striped together with RAID0.

IF DEPLOYING SAS GRID MANAGER, WHICH SHARED FILE SYSTEM TO USE

SAS Grid Manager requires a shared file system for the permanent files being shared by all the SAS Grid compute nodes. Below are several shared file systems that have been tested with SAS Grid in several public cloud infrastructures.

- Intel Cloud Edition for Lustre (ICEL) Software is available in Amazon and Microsoft Azure (see links below). Information about how to configure Lustre for the public cloud is discussed in this paper from 2015. EBS file systems used by Lustre should be created in the same manner as those used for permanent SAS data files (that is, striping together at least 4 EBS volumes).

For optimal performance, you will need the following as *minimums* for your Lustre shared file system.

- OssCount: 4
- OssInstanceType: C4.8xlarge
- OstVolumeCount: 4

This will give you a total of 16 OSTs (Logical Units). You should also set “EnabledOstRaid” to false. This will change the striping on directories of files (via “lfs setstripe”) to true.

To learn more about ICEL in the public cloud, please refer to the following links:

- Amazon - https://aws.amazon.com/marketplace/seller-profile/ref=dtl_pcp_sold_by?ie=UTF8&id=d1c6e336-5f6f-4234-82a8-a57463081a35
- Microsoft Azure - <https://azuremarketplace.microsoft.com/en-us/marketplace/apps/intel.intel-cloud-edition-gs>
- IBM Spectrum Scale (formerly known as GPFS) - IBM and SAS tested the IBM Spectrum Scale in AWS at the end of 2017 with positive results. A paper is currently being written to document the testing and the recommended hardware and environment configurations for achieving optimal performance. This paper (“Performance and Tuning Considerations on Amazon WebServices with SAS® 9.4 using IBM Spectrum Scale™”) will be ready by SAS Global Forum 2018. IBM is working with Amazon to have Spectrum Scale added as an offering in the Amazon Marketplace.

- EFS storage - you can use this type of storage, but you need to understand that the maximum IO throughput for each EFS file system is only 250 MB per second per node. You can have multiple EFS file systems per node to overcome the 250 MB per second per node limitation, but you need to understand that there is a limitation of the single NIC card in each AWS EC2 instance.

In addition, if you are looking to place SAS Home and SAS Configuration files on EFS, you will need to add the following statement to your sasv9.cfg file.

“-filelocks <HOME> none” where <HOME> is the directory containing the SAS Home and SAS Configuration files.

File locks need to be turned off is because EFS has a hard limitation of 87 locks per file and SAS manages locks on various files in SAS Home.

WHERE TO PLACE SAS CLIENTS

The most common placement of SAS clients like SAS® Enterprise Guide®, SAS® Data Integration Studio, and SAS® Studio (even though this is a web-based client, the issue with its placement in comparison to the SAS processing is very important) is within the public cloud infrastructure (either a Windows server or virtual desktops). The question is where to put these Windows systems—in the public cloud in the same availability zone and placement group, on a system within your current data center, or on a desktop/laptop?

The answer depends on the volume of data being transferred back to the SAS client. If there is a lot of data being transferred to populate drop-down windows or to view table, then having the clients and the back-end SAS servers on the same LAN will equate to the fastest results. For public clouds, this would mean standing up an instance with a Windows server on it and placing your SAS clients on that instance. Your SAS users will have to access their SAS clients on the Windows server inside the public cloud infrastructure.

In that case, it is important that the Windows server is located in the same region, placement group, and availability zone as the SAS deployment.

WHERE TO PLACE DATA TO BE USED BY SAS

SAS has several tools that allow sharing of SAS data files on-premises with SAS applications that are run in a public cloud and vice versa. While these tools function well, please note that the IO throughput between your data center and the public cloud might be as low as 500 KB/second. Rarely have we seen them to be more than 20 MB/sec. If the SLAs for your SAS applications in the public cloud can be met with this slow data transfer rate, then it is acceptable to keep your data in a different physical location than where SAS is running.

Our best practice is to have SAS applications and the data associated with the SAS application reside in the same data center—whether on-premises or in the same placement group in the public cloud.

WHERE TO PLACE AUTHENTICATION TOOLS

The considerations needed when deciding where to place your authentication tools are similar to those needed when deciding where to place SAS clients or data used by SAS: in the public cloud in the same availability zone and placement group or on a system within your current data center?

This depends on how often your SAS application will reach out to your authentication tools for permission to use a file. If this is a high number, it is recommended that you move your authentication tools and associated data into the public cloud. Again, these need to run on a system in the same region, placement group and availability zone.

ARE HIGH AVAILABILITY AND SECURITY REQUIRED?

Depending your need for high availability (HA), you might just need processes in place to quickly create a new cloud host instances in case one of your existing instances fail. This is more of a failover HA practice. The ability to quickly spin up a new instance is one of the benefits of running SAS in a public cloud.

For SAS Grid customers, please note that a shared file system (for example, Intel Lustre, IBM Spectrum Scale) will remain operational if one of the nodes associated with the shared file system goes down. But any data that is associated with that node will not be available until the node is restored. This is due to only one copy of the data being stored by default. However, you can enable replication services with these shared file systems so that two or three copies of your data are stored on multiple nodes of the shared file system. This does drive up the cost, especially if you have hundreds of terabytes of data.

If your definition of High Availability includes Disaster Recovery, then you will need to look at mirroring their SAS deployment, SAS files, and their data store in another region of the public cloud or a separate cloud. The SAS Application considerations for Disaster Recovery will apply to the public cloud in the same manner as on-premises infrastructures. Details on what needs to be considered in a Disaster

Recovery SAS implementation can be found in [“Do You Have a Disaster Recovery Plan for Your SAS® Infrastructure”](#) (*Proceedings of the SAS Global Forum 2016 Conference*).

CONCLUSION

This paper has been written to help SAS customers understand the public cloud instance types that best meet the needs of their SAS users. Many of the topics discussed in this paper are based on real world experiences by the first set of SAS customers standing up SAS in a public cloud offering. These customers have learned that there are very few things from a hardware and storage perspective that can be tuned, but the ones that can be are listed in this paper. The chief takeaway is that IO throughput is crucial, and unfortunately, is a limiting factor in the success of SAS public cloud deployments.

The intent of this paper is to raise concrete awareness of SAS application requirements, and how best to meet them in a public cloud. The choice of hardware resources, data stores, and application architecture placement is crucial to achieving the best performance the cloud can offer. The typical SAS minimum recommended IO throughput of 100 MB/s per core, versus most public cloud’s maximal delivery of 50 MB/s per core for off-board storage (for example, EBS) is a key consideration. You must carefully consider the performance ramifications of moving on-premises SAS applications to the public cloud, unless you can suffice with a dropped IO bandwidth performance for persistent data activity.

And as mentioned at the start of this paper, it should be noted that the information discussed in this paper is based on what is available from the public clouds and SAS experience with the public clouds at the time of the writing of this paper. Public cloud offerings are constantly changing. Therefore, it is your best interest to understand the rationale used in the selection process and consider what was done as a point-in-time design.

REFERENCES

- Amazon Cloud. “Amazon EBS-Optimized Instances.” Available at <https://docs.aws.amazon.com/AWSEC2/latest/UserGuide/EBSOptimized.html#ebs-optimization-support>.
- Amazon Cloud. “Amazon EC2 Instance Types.” Available at <https://aws.amazon.com/ec2/instance-types/>.
- Crevar, Margaret. April 2015. “Performance and Tuning Considerations for SAS Grid Manager 9.4 on Amazon Cloud using Intel Cloud Edition for Lustre File System.” Available at <http://support.sas.com/rnd/scalability/grid/SGMonAWS.pdf>.
- Google Cloud. “Machine Types.” Available at <https://cloud.google.com/compute/docs/machine-types>.
- Microsoft Azure. “High Performance Compute VM Sizes.” Available at <https://docs.microsoft.com/en-us/azure/virtual-machines/windows/sizes-hpc>.
- Red Hat. “Optimizing SAS on Red Hat Enterprise Linux (RHEL) 6 & 7.” Available at http://support.sas.com/resources/papers/proceedings11/342794_OptimizingSASonRHEL6and7.pdf.

ACKNOWLEDGMENTS

Many thanks to Ande Stelk and Ed Gaines for their help researching several of the sections of this paper. And to Tony Brown and Jim Kuell for reviewing the content of the paper.

CONTACT INFORMATION

Your comments and questions are valued and encouraged. Contact the author at:

Margaret Crevar
SAS Institute Inc
Margaret.Crevar@sas.com

Ande Stelk

SAS Institute Inc
Ande.Stelk@sas.com

SAS and all other SAS Institute Inc. product or service names are registered trademarks or trademarks of SAS Institute Inc. in the USA and other countries. ® indicates USA registration.

Other brand and product names are trademarks of their respective companies.