# SAS Enterprise Miner 12.1 on IBM PureFlex System with IBM POWER7 processor-based compute nodes and IBM XIV Storage System Gen3

*Performance and scalability with 4- to 16-user workloads*

*Jerrold Heyman*
*Srirama Sharma*

*IBM Systems and Technology Group ISV Enablement*
*October 2013*

@IBMSystemsISVs

# Table of contents

*SAS Enterprise Miner 12.1 on*
*IBM PureFlex System with IBM POWER7 processor-based compute nodes and IBM XIV Storage System Gen3*

# Abstract

*This paper demonstrates the excellent performance and scalability of the IBM Flex System p460 Compute Node with the IBM AIX version 7.1 operating system when running the SAS Enterprise Miner 12.1 application. IBM and SAS engineers collaborated to run a 4-, 8-, and 16-user benchmark using small, medium, and large data tables from a retail bank's marketing initiative. The results show that the IBM POWER7 processor-based compute node scales proportionally as the number of users grow. The use of IBM XIV Storage System Gen3 completes the solution stack to meet the need for higher I/O bandwidth requirements.*

# Introduction

Turning increasing amount of raw data into useful information has become increasingly important in today's highly competitive business environment, increasing the demand for predictive, analytics data mining solutions. SAS Enterprise Miner a powerful and comprehensive data mining software help organizations explore large quantities of data to discover relationships and patterns that lead to proactive decision making.

Unmatched IBM® expertise in hardware and software technology enables the SAS Enterprise Miner for IBM AIX® 7.1 on IBM POWER7® processor-based server to deliver significant benefits. It can be deployed on an infrastructure that is designed to improve reliability, performance, and scalability. Combined with IBM PowerVM® and IBM System Storage® solutions, the solution described in this paper provides unparalleled function and performance for implementing data-mining solutions in small to large environments.

# Test scenario

The test scenario is designed to simulate a data mining environment under heavy use. The performance tests simulate 4 to 16 analysts who run predictive modeling projects using the SAS Enterprise Miner software. The goal is to understand and predict the customer behavior for a retail bank's marketing initiatives.

Real-world data volumes and structures are used for these tests. Input tables for data mining projects typically contain many rows (observations) and columns (explanatory variables) of data. It is increasingly common for data miners to use hundreds of independent variables to explain the behavior of a single dependent variable.

The users in this scenario previously used SAS Enterprise Miner to develop modeling process flow diagrams and save the flows as batch code. For these batch tests, all users start running their modeling flows simultaneously. The analysts save their modeling flows as batch code so that they can schedule the execution to occur only when the server is available to handle the processing. This represents an unusually high level of demand for the server.

Three tests are run for each of the 4-, 8-, and 16-user scenarios, each using a different input data table.

Table 1 shows the volume of data used for these tests.

| Table | Size | Number of observations | Number of variables |
|-------|------|------------------------|---------------------|
| Small | 439 MB | 500,000 | 100 (50 interval, 25 categorical, 25 binary) |
| Medium | 1.7 GB | 1,000,000 | 200 (100 interval, 50 categorical, 50 binary) |
| Large | 7.6 GB | 2,000,000 | 500 (250 interval, 125 categorical, 125 binary) |

*Table 1: Data volumes used in the test scenario*

Figure 1 shows the modeling flow diagram run by each of the analysts. All the tests involved concurrent execution of SAS Enterprise Miner modeling flows for all users. Concurrency is defined as workloads running during the same period of time.
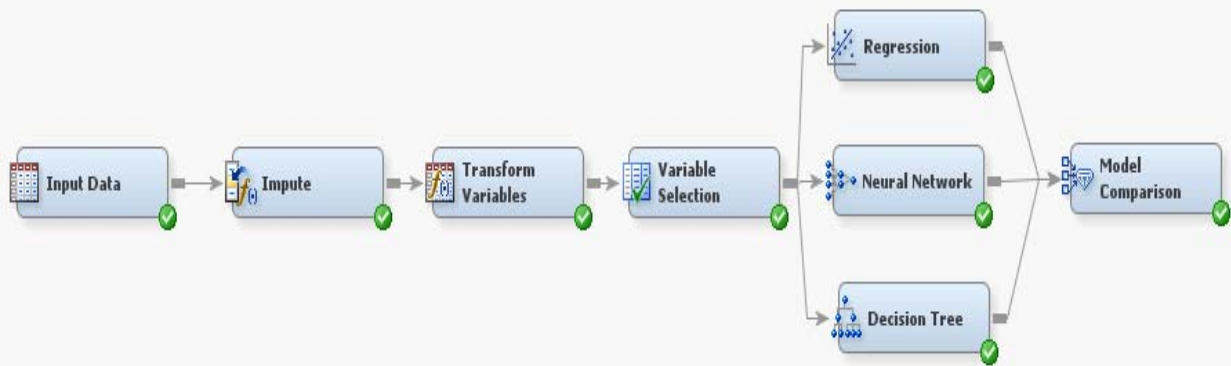


*Figure 1: Modeling flow diagram run by each user*

A workload can be a single process or multiple processes using the resources constantly or intermittently, depending on the task being run for a certain period of time. Not all workloads use the same amount of system resources. Each modeling flow runs a decision tree, regression, neural network node, and other common nodes.

# Deploying the SAS Enterprise Miner server

The 4-, 8-, and 16-user tests are conducted on a SAS environment installed on a POWER7 system with PowerVM features. For all tests, dedicated logical partitions, and virtual I/O attached storage for all storage requirements are used.

## Deploying the hardware

The SAS Enterprise Miner is deployed on a single logical partition installed with AIX 7.1. The number of active cores was varied for each test depending on the number of users running the tests. The 4-user test ran on a 4-core dedicated configuration, the 8-user test ran on an 8-core dedicated configuration, and the 16-user test ran on a 16-core dedicated configuration. In all cases, there were 2 Virtual I/O Server (VIOS) partitions, each of them configured with 2 dedicated cores.

A VIOS serves the disks that are used for the AIX system file systems and SAS file systems (/sasem_data and /sasem_work). The hardware used is an IBM PureFlex™ System with a single IBM Flex System™ p460 Compute Node, which is a POWER7 processor-based system. Figure 2 shows the basic deployment of the hardware environment.
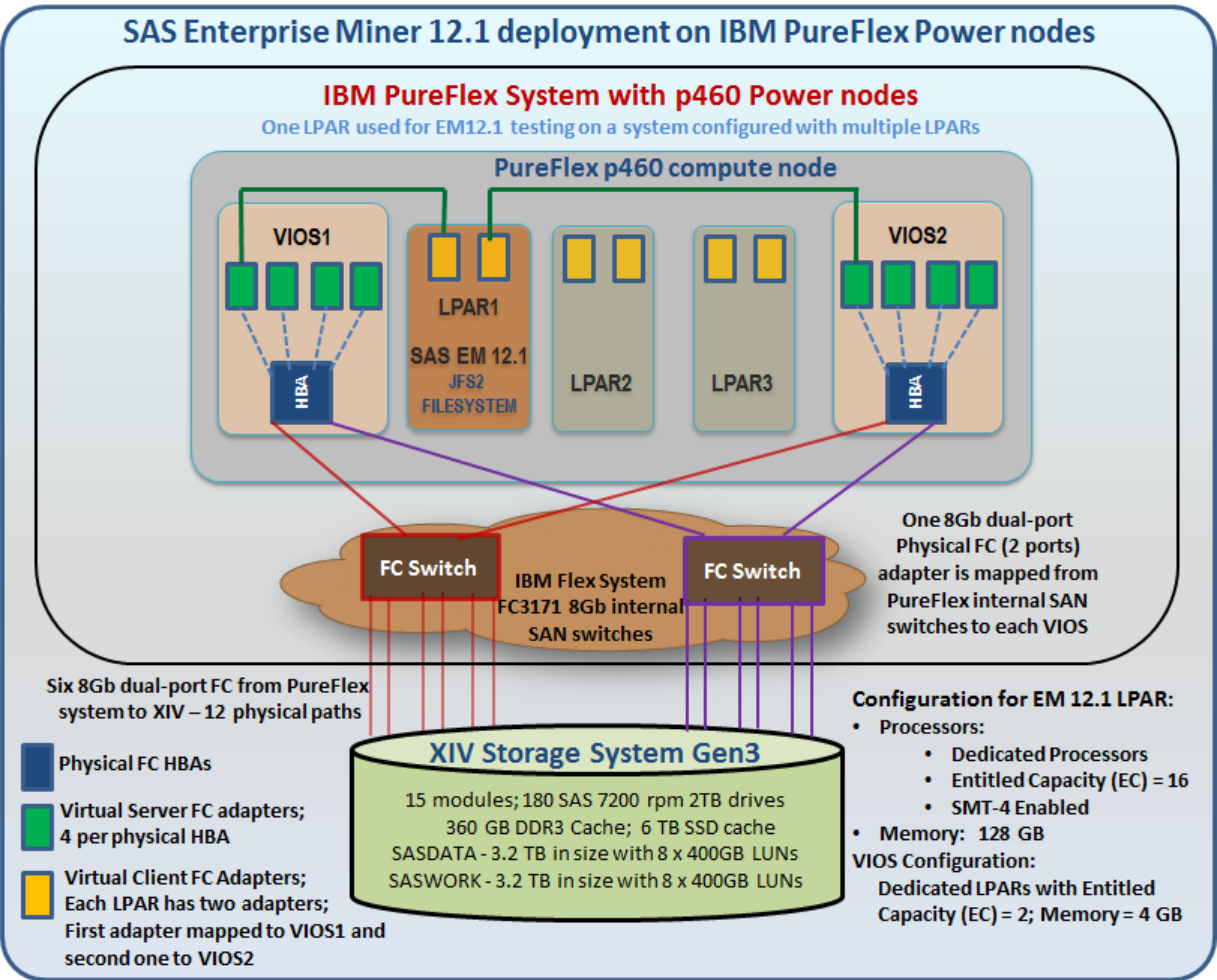


*Figure 2. High-level logical partition server configuration*

The SAS Enterprise Miner 12.1 is deployed on the AIX 7.1 logical partition (LPAR) along with the prerequisite IBM WebSphere® Application Server (version 7.0.25). This LPAR is connected to VIOS server through virtual Fibre Channel (FC) adapters using N-Port ID Virtualization (NPIV). Through this, the external storage disks from IBM XIV® Storage System Gen3 storage device are directly accessed by the client LPAR in pass-through mode. This VIOS-attached storage for this partition contains the SAS binary files, various log files, the /sasem_work temporary file system, and the /sasem_data file system that are allocated to logical volume groups.

Table 2 shows the core and memory allocation for each of the three tests. The XIV Gen3 storage solution provides the necessary storage.

| User test | Number of cores | Memory |
|---|---|---|
| Four users | 4 (dedicated) | 32 GB |
| Eight users | 8 (dedicated) | 64 GB |
| Sixteen users | 16 (dedicated) | 128 GB |

*Table 2: Core and memory allocations for the 4-, 8-, and 16-user tests*

The configuration on Figure 2 shows the use of a VIOS. The VIOS provides the interface between the physical storage area network (SAN) attached storage with the LPAR. The SAN network switch is contained within the IBM PureFlex System chassis and is attached directly to XIV Gen3 through FC.

The XIV Gen3 system consists of 180 physical disks. Best practice dictates that SAS runtime data is spread out to as many physical disks as possible. The storage configuration for the SAS Enterprise Miner LPAR is used to illustrate some main points in managing disk space that is provided by the XIV Gen3 system.

To strike a balance between disk I/O performance and economies of scale, the storage system is configured as RAID5. The minimal unit, also known as a logical unit number (LUN), on the system that a piece of storage can be created is called an *array*. Each array consists of eight physical disks. Figure 3 shows a representation of a LUN consisting of an array with eight disks.
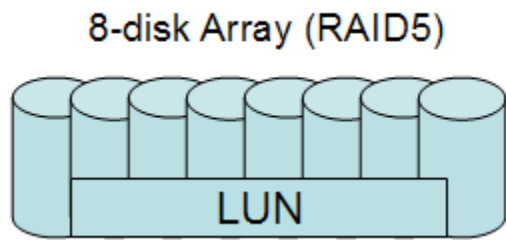


*Figure 3: Representation of a LUN that consists of an array with eight disks*

With 180 physical disks, the disks are grouped into arrays of eight disks, resulting in 16 arrays. To simplify AIX disk management, you can create a LUN from each array and then present it to AIX
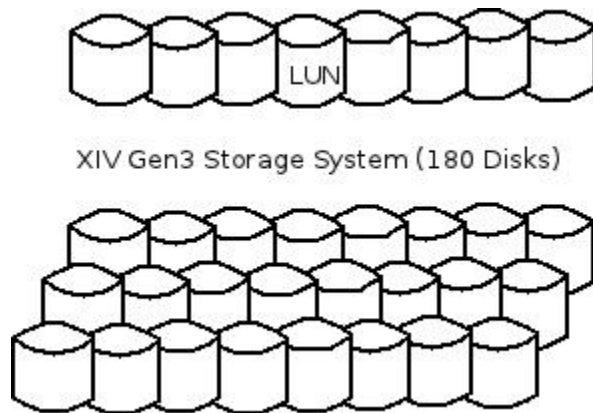(refer to Figure 4).

*Figure 4: A 180-disk subsystem configured with 16 arrays*

From the receiving AIX partition point of view, the LUNs are seen as 16 logical disk drives (hdisks) from which two volume groups (sasdatavg and sasworkvg) are created. Striped logical volumes are added, resulting in a file system layout that is evenly distributed among all the disks.

Separate file systems for /sasem_work and /sasem_data are created on the volume groups **sasworkvg** and **sasdatavg** respectively. The /sasem_data file system contains a subdirectory associated with each user. The SAN configuration to XIV Gen3 provided by the PureFlex System has two 6-port switches, with each port of the first switch allocated to a channel on XIV Gen3, and the second switch mirroring the first – providing two distinct SAN paths to a given XIV channel.

Finally, the environment was tuned by following the tuning guidelines for SAS 9.3 on AIX 7.1 available in the link: **ibm.com**/support/techdocs/atsmastr.nsf/WebIndex/WP101529.

# Benchmark results

The tests consist of running 4-, 8-, and 16-user scenarios. Each user test is run three times separately using a different input data table each time. The data used for these tests has the characteristics shown in Table 3.

| Table | Size | Number of observations | Number of variables |
|-------|------|------------------------|---------------------|
| Small | 439 MB | 500,000 | 100 (50 interval, 25 categorical, 25 binary) |
| Medium | 1.7 GB | 1,000,000 | 200 (100 interval, 50 categorical, 50 binary) |
| Large | 7.6 GB | 2,000,000 | 500 (250 interval, 125 categorical, 125 binary) |

*Table 3: Data volumes for the test scenario*

# Four-user scenario

Table 4 summarizes the maximum amount of time taken by each SAS Enterprise Miner node to run across all four users. As the data volume increases, the run time also increases. Time is displayed in the hh:mm:ss format.

| Table | Input data set | Impute | Transform | Variable selection | Neural | Regression | Decision tree | Model comparison |
|---|---|---|---|---|---|---|---|---|
| Small data | 0:00:04 | 0:00:06 | 0:00:10 | 0:00:12 | 0:04:37 | 0:00:17 | 0:00:38 | 0:00:32 |
| Medium data | 0:00:10 | 0:00:16 | 0:00:24 | 0:00:28 | 0:14:02 | 0:00:36 | 0:00:51 | 0:01:15 |
| Large data | 0:00:30 | 0:00:47 | 0:01:07 | 0:03:19 | 0:44:30 | 0:01:56 | 0:04:08 | 0:05:52 |

*Table 4: Maximum amount of time taken by each SAS Enterprise Miner node to run for the 4-user test*

Figure 5 shows the average percentage of time for which the processors are busy processing the small, medium, and large input data tables in the 4-user scenario.
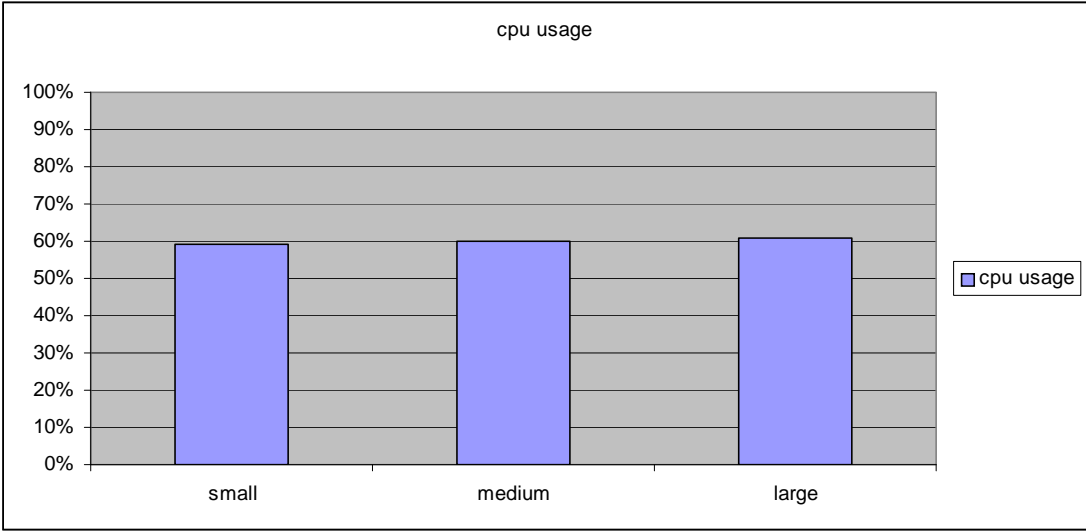


*Figure 5: Percentage of time the processors in the system are busy in the 4-user scenario*

## Eight-user scenario

Table 5 summarizes the maximum amount of time taken by each SAS Enterprise Miner node to run across all eight users. As the data volume increases, the run time also increases. The time is displayed in the hh:mm:ss format.

| Table | Input data set | Impute | Transform | Variable selection | Neural | Regression | Decision tree | Model comparison |
|-------|---------------|--------|-----------|-------------------|--------|-----------|--------------|------------------|
| Small data | 0:00:06 | 0:00:06 | 0:00:11 | 0:00:13 | 0:05:04 | 0:00:18 | 0:00:35 | 0:00:40 |
| Medium data | 0:00:04 | 0:00:13 | 0:00:20 | 0:00:38 | 0:14:05 | 0:00:40 | 0:01:10 | 0:01:32 |
| Large data | 0:00:44 | 0:00:54 | 0:01:12 | 0:05:05 | 0:50:53 | 0:02:31 | 0:04:28 | 0:05:17 |

*Table 5: Maximum amount of time taken by each SAS Enterprise Miner node to run for the 8-user test*

Figure 6 shows the average percentage of time for which the processors are busy processing the small, medium, and large input data tables in the 8-user scenario.



*Figure 6: Percentage of time the processors in the system are busy in the 8-user scenario*

# Sixteen-user scenario

Table 6 summarizes the maximum amount of time taken by each SAS Enterprise Miner node to run across all 16 users. As the data volume increases, the run time also increases. The time is displayed in the hh:mm:ss format.

| Table | Input data set | Impute | Transform | Variable selection | Neural | Regression | Decision tree | Model comparison |
|-------|--------|--------|-----------|-------------------|--------|------------|--------------|------------------|
| Small data | 0:00:06 | 0:00:07 | 0:00:12 | 0:00:15 | 0:05:48 | 0:00:23 | 0:00:35 | 0:00:46 |
| Medium data | 0:00:04 | 0:00:14 | 0:00:22 | 0:00:42 | 0:14:46 | 0:00:47 | 0:01:15 | 0:01:48 |
| Large data | 0:00:49 | 0:00:53 | 0:01:16 | 0:05:53 | 0:54:52 | 0:02:39 | 0:05:36 | 0:05:43 |

*Table 6: Maximum amount of time taken by each SAS Enterprise Miner node to run for the 16-user test*

Figure 7 shows the average percentage of time for which processors are busy processing the small, medium, and large input data tables in the 16-user scenario. The figure shows the fall in percentage of processor utilization specifically when using the large data set. In this case, the processors are not fully used because they are waiting for data from the storage array.



*Figure 7: Percentage of time the processors in the system are busy in the 16-user scenario*

Note that because of the high I/O bandwidth requirement (more than 700 MBps throughput) of the 16-user and large data test, the percentage of the processor usage is down to 60%. In this case, it is recommended to configure additional disks to the XIV Gen3 device to accommodate the I/O throughput requirement when running this particular test.

# Comparison to previous SAS Enterprise Miner benchmark

In a previous SAS Enterprise Miner benchmark as documented in *SAS Enterprise Miner Performance on IBM System p 570* at http://www.sas.com/partners/directory/ibm/SASEMonSystemp570.pdf, you can see that using later versions of both AIX and SAS Enterprise Miner results in better performance in the 4- and 8-user tests using the small- and medium-sized data sets respectively.

Table 7 shows the 4-user and small volume test comparison between the benchmark using the earlier and the current versions of SAS Enterprise Miner and AIX.

| Version | Transform | Variable selection | Neural | Decision tree | Model comparison |
|---|---|---|---|---|---|
| SAS Enterprise Miner 5.2 and AIX 5.3 | 0:00:30 | 0:00:50 | 0:09:30 | 0:05:50 | 0:0000:35 |
| SAS Enterprise Miner  6.1 and AIX 6.1 | 0:00:09 | 0:00:20 | 0:08:25 | 0:01:16 | 0:00:35 |
| SAS Enterprise Miner 12.1 and AIX 7.1 | 0:00:10 | 0:00:12 | 0:04:37 | 0:00:38 | 0:00:32 |

*Table 7: A 4-user, small data volume test*

Table 8 shows the 8-user and medium volume test comparison between the benchmark using the earlier and the current versions of SAS Enterprise Miner and AIX.

| Version | Transform | Variable selection | Neural | Decision tree | Model comparison |
|---|---|---|---|---|---|
| SAS Enterprise Miner 5.2 and AIX 5.3 | 0:01:45 | 00:13:20 | 0:35:00 | 0:15:12 | 0:01:35 |
| SAS Enterprise Miner 6.1 and AIX 6.1 | 0:00:33 | 0:00:45 | 0:29:49 | 0:05:28 | 0:01:25 |
| SAS Enterprise Miner 12.1 and AIX 7.1 | 0:00:20 | 0:00:38 | 0:14:05 | 0:01:10 | 01:32 |

*Table 8: An 8-user, medium data volume test*

# Summary

The results of this SAS Enterprise Miner performance test suite successfully demonstrate the performance, reliability, and scalability on the IBM Flex System p460 compute nodes. IBM breakthrough hardware and operating system technology, together with the robust scalability of SAS software, delivers an ultimate solution to address analytical mining demands for SAS and IBM clients.

# Acknowledgments

# Resources

The following websites provide useful references to supplement the information contained in this paper:

- Power Systems on IBM PartnerWorld®
  **ibm.com**/partnerworld/systems/p

- AIX on IBM PartnerWorld
  **ibm.com**/partnerworld/aix

- IBM Systems on IBM PartnerWorld
  **ibm.com**/partnerworld/systems/

- IBM Publications Center
  www.elink.ibmlink.ibm.com/public/applications/publications/cgibin/pbi.cgi?CTY=US

- IBM Redbooks®
  **ibm.com**/redbooks

- IBM developerWorks®
  **ibm.com**/developerworks

# About the authors

**Jerrold Heyman** is a Technical Consultant in IBM Systems and Technology Group ISV Enablement. He has more than 25 years of experience working in the IT Industry. He currently works with several software vendors to enable their enterprise applications on the latest IBM operating systems. You can reach Jerrold at jheyman@us.ibm.com.

**Srirama Sharma** is a Technical Consultant in IBM Systems and Technology Group ISV Enablement. He has 9 years of experience working in the IT Industry. He currently works with several independent software vendors (ISVs) to enable their enterprise applications on the latest IBM platform. You can reach Srirama at sriramsh@in.ibm.com

@IBMSystemsISVs

# Trademarks and special notices