**Paper 4656-2020**

# Architecting SAS® Analytics Edge-to-Cloud Solutions on Hewlett Packard Enterprise Infrastructure

Mark Barnum   Hewlett Packard Enterprise

Kannan Mani    Hewlett Packard Enterprise

## ABSTRACT

We live in a world where everything computes. Where technology, apps, and data are driving digital transformation, reshaping markets, and disrupting every industry. IT needs to reach beyond the traditional data center and the public cloud to form and manage a hybrid connected system stretching from the edge to the cloud. Listen to experts on how Hewlett Packard Enterprise (HPE) and SAS help organizations to rethink and modernize their infrastructure more comprehensively to deploy their edge-to-cloud architecture and analyze data at the edge quickly.

Join us, in this session for a deep dive and demo on:

HPE Elastic Platform for Analytics (EPA) architecture for SAS and use cases
SAS Event Stream Processing (ESP) and SAS Visual Analytics (VA) deployment at edge Cloudera Data Platform and NiFi

## INTRODUCTION

Today, more data is being generated, captured, and analyzed than ever before. Data is being captured from devices across industries at the intersection of people, places, and things. When we have the infrastructure to act upon the data where it is generated, we call this the Intelligent Edge. This has been made possible, in part, by a massive proliferation of high speed sensors that can detect and transmit data from devices in real time that are coupled with powerful edge computer systems that perform historical and predictive analysis.

With the large amount of data being generated at the edge, we need to determine how that data is used to create the largest benefit. Some data will be used immediately, right at the edge to monitor, make changes, and alert based on that data in real time. Other data may not be as urgent, and needs to be sent back to the core data center (DC) or cloud, stored in a longer term repository from which batch analytics can be used to build case sets. Those case sets are then used to develop machine learning/deep learning (ML/DL) models that can be pushed back to the edge to refine real-time and predictive analysis.

Because ML/DL model accuracy depends directly upon the volume and freshness of the data, the feedback loop becomes nearly continuous. Figure 1 depicts the edge-to-core data feedback loop.
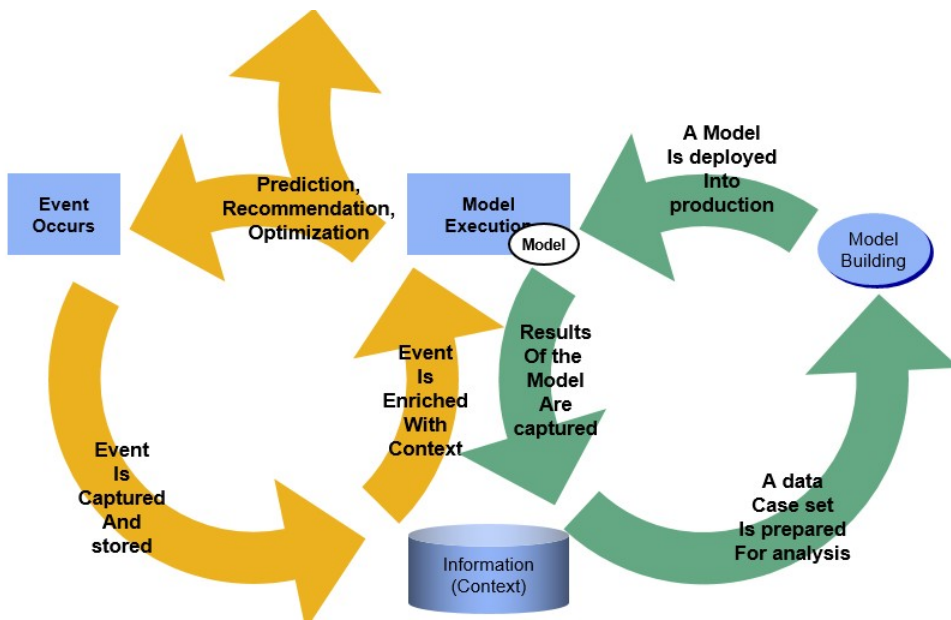
**Figure 1.** Edge-to-core data feedback loop

In order enable this continuous feedback loop, we need to create a data pipeline to deliver the right data to the right place, storing the information for easy retrieval by all participating members of the cluster.

To facilitate a design to comprise all facets of the data pipeline, we need to be able to effectively and efficiently move data from one layer of the model to the next. In effect, we need to build an end-to-end data pipeline using best of breed hardware, software and utilities to move and convert data from the source format to the target format. At the edge, data is in a number of different formats and those formats need to be transformed to a format the analysis engine understands.

HPE, SAS, and Cloudera have the components to enable all of the capture, movement, and storage analysis of that data from the edge to the core data center/cloud and back to the edge. HPE has a set of workload optimized systems that, when coupled with software from SAS and Cloudera, allows you to build a data pipeline that's high performance, scales and brings together the myriads of data sources in your environment to derive value and intelligence from the data. From the HPE Edgeline systems on the edge to a range of workload-optimized platforms in the core as part of the HPE Elastic Platform for Analytics architecture, customers can create data pipelines that can be quickly composed and recomposed to drive data intelligence and action.

This paper will provide detail about how the proof of concept (POC) was conducted and how to create the data pipeline that enables the constant feedback loop needed for up to date inference and control at the edge for electric utility use case.

In order to demonstrate the importance of a full end-to-end data pipeline, this paper will highlight the use case of a Smart Grid. A Smart Grid allows electric utilities to reduce costs, enhance the safety of its workers and customers and enhance its customer's experience.

Today, electricity generation turbines are taken offline for preventative maintenance based on time and usage parameters. This maintenance can lead to waste in that maintenance may be scheduled but not necessary at the moment. However not having a regular maintenance schedule may lead to unexpected failure of the turbine, which may end up more impactful and expensive than pulling the turbines offline on a periodic basis.

Consider how the business model would change if sensors were located on the turbine and predictive maintenance was utilized. Sensors on the turbine could be monitored every 60th of a second, less or more often based on the specific requirement needed to determine trends in the turbine's operating characteristics. Once a trend toward an error condition was identified, the turbine could be brought down immediately for maintenance or the maintenance event could be scheduled, depending on the severity of the trend.

The data from the turbine can be aggregated and sent to the data center or cloud where it can be blended with other data, such as historical equipment performance information derived during maintenance operations. That data can then be operated upon by artificial intelligence or machine learning platforms to refine the analysis model. The model(s) then gets pushed back out to the edge for use regular monitoring activity.

Another example is the electric utility industry. Transformers can be large expense items for electricity generation, transmission and distribution operations. Failure of transformers can have great impact on availability of power in the grid. New transformers are being manufactured with sensors built into the units. However, there are a large number of legacy transformers that are still performing well, but don't have sensors built into them. Electrical utilities are barred, by regulation, from modifying any of their existing transformers, (e.g. embedding sensors into the transformer). Non-invasive additions can be allowed such as adding an infrared camera to monitor transformer temperature.

In this manner, it is possible to determine, in advance, when a transformer is trending toward failure. Ambient temperature, combined with humidity can be factored with transformer bushing[1] temperature to predict failure. As additional data, such as reliability models for new transformers, are added to the model, the inference model can be refined making it more accurate over time. While the determination of an upcoming transformer failure should probably be made at the edge, the generation of new reliability models need to be generated in the data center or cloud to ensure all data is included in the new model.

Another example of the use of edge analytics in the electricity generation industry involves power line breaks. Edge analytics within the transformer vault can detect a precipitous drop in demand within a sector when a line break occurs. Collection of immediate data from smart meters can further refine the sector down to a specific location, sending an alert to the data center or cloud.

At that point data center analytics need to be brought to bear, to determine if there is a way to redirect transmission to keep electricity flowing to the affected customers. Additionally, data center analytics can determine where current work crews are deployed and the number and    importance of the customers without power to direct crews to the next issue with the greatest efficiency and impact.

Finally, in all of the scenarios where maintenance is required, the data center applications can push information out to the worker crews to make their jobs safer and easier to understand. Augmented reality can be enabled within the lineman's truck or the generation station's maintenance department so that workers can not only understand what's required, but also see the repair in graphical detail. This type of technology can aid in the determination of whether to send a worker to a transformer vault or whether to call for the fire department.

In all of the above scenarios, transmitting all of the sensor data gathered in the field to the data center for analysis may not be practical or meet real-time needs for analysis. At the same time, an aggregated or summary data stream from each collection point could be transmitted to the data center where it could be blended with data from other, related collection points to gain insight that might benefit the overall business. This consolidated data and analysis apart from any single collection point, in the data center could help recognize overall trends to be remediated for all like or similar assets.

In general, there are seven considerations why all data generated at the edge does not necessarily need to be transmitted to the DC[2]. They are:

1. Latency - The large amount of data can cause a delay between when the data was collected and when it arrives in the DC.

2. Bandwidth - The larger the amount of data to be transferred, the more bandwidth is required to transmit it in a given period of time. This problem is exacerbated if the data needs to travel outside of an enterprise's LAN.

---

[1] A bushing is an insulated device that allows an electrical conductor to pass safely through a grounded conducting barrier such as the case of a transformer or circuit breaker. Wikipedia:    https://en.wikipedia.org/wiki/Bushing_(electrical)

[2] Seven Reasons to Compute at the Edge: https://www.engineering.com/IOT/ArticleID/15540/Seven-Reasons-to-Compute-at-the-Edge.aspx

3. Cost - Expanding bandwidth is expensive. This is especially true if a WAN is required.

4. Compliance - Certain types of companies have regulations they must follow when transmitting data. For instance, health care companies dealing with patient data must comply with HIPPAA regulations.

5. Security - Any time data is placed on a wire, it is subject to unauthorized access. This is especially true if the data is transmitted to remote locations.

6. Duplication of data - The data is duplicated on each server to which it's copied. This increases the overall storage space required.

7. Reliability - The farther data is carried from the edge, the more single points of failure are introduced. Performing as much processing as possible at the edge ensures reliability remains high.

Each industry to which edge computing is applied has unique requirements and will deploy specialized clusters of computers, networking and storage devices to take advantage of those different demands.

A few examples of industries' use cases that are shaping edge computing include, but certainly not limited to:

1. Power generation and distribution - Smart Grid

2. Oil and gas - Refining, exploration and recovery

3. Medical - Magnetic resonance scanners

4. Manufacturing - Industrial controller and supervisory control and data acquisition (SCADA)

5. Motor vehicle - self-driving and assisted driving cars

6. Smart City - Traffic lights and geographical scope traffic monitoring and flow control

7. Home Automation - Thermostats, home appliances and lights

8. Telecommunication - Cell phones and cellular networks

Each of the above industries can find advantage by processing data at the edge, and transmitting a summary or exception subset of that data to a core data center or cloud computing storage environment, for analysis.

In the first instance, an electric utility needs to understand when a transformer is being pushed beyond its rated limits. The rated limit can be found on the transformer data plate. Transformers wear based on their time at a specific rating. A transformer that will last 40 years when kept to a 100% or less of their rating, will last a lesser amount of time when it's consistently asked to convert energy at 110% or even 150%.

As stated previously, there are regulations that keep power companies from modifying existing transformers without having to recertify them after any modification is made to them. For this reason, HPE and SAS are promoting the possibility of monitoring an older transformer using an infrared camera. The infrared camera shows the temperature of the bushings on the transformer and this can be used to infer the usage of the transformer. The following, figure 2 shows a demonstration environment where we have an infrared camera focused on a model of a transformer and SAS® Event Stream Processing (ESP) is monitoring the temperature of the bushings.
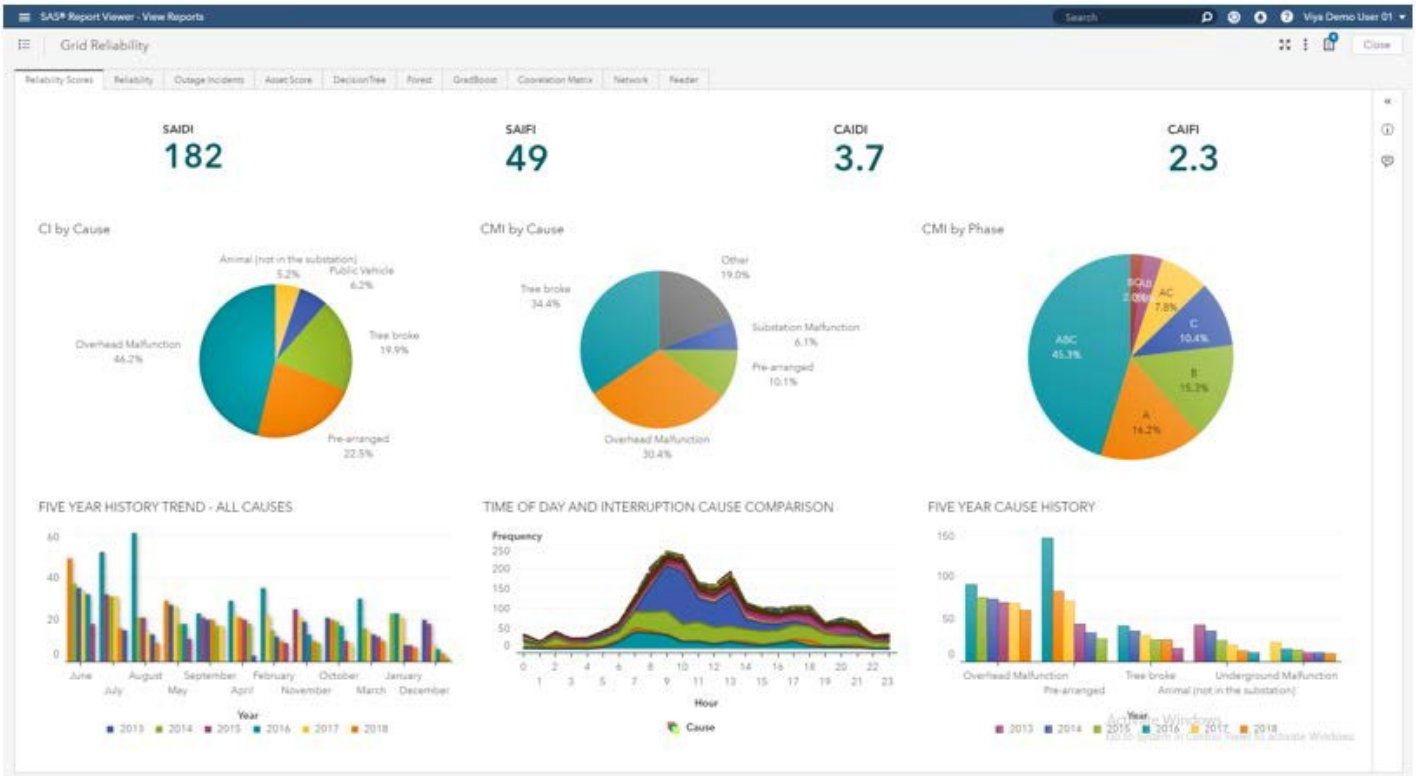
**Figure 2.** SAS Event Stream Processing monitoring the temperature of bushings on a transformer to determine at what rate the transformer is being used

The SAS Event Stream Processing model can include other factors to determine if or when it should alert based on temperature. For example, a bushing temperature of 110° F may not be an issue if the ambient temperature is 80° F, but might be an issue if the ambient temperature is 40°F.

Newer transformers have sensors built in to them that allow for more closely monitored power conditions.

The number of events created and monitored per time interval is determined by how fast a sensor can report conditions. It is better to have more events monitored over a given period of time, because that allows for swifter reactions to possible issues that can arise over the life of the transformer. However, it would be prohibitive to create a network from the edge to the data center that would be able to transmit all of the   events. For this reason, data is typically condensed, before it's sent. For instance, you may monitor a transformer 60 times a second, but you may also transmit only one synopsized record to the data center every second or even every minute.

Once the data arrives in the data center it can be combined with data from other transformers of similar or even dissimilar models. That data can then be used to determine the health of the overall grid. The following, figure 3 shows SAS Event Stream Processing monitoring the health of the grid.

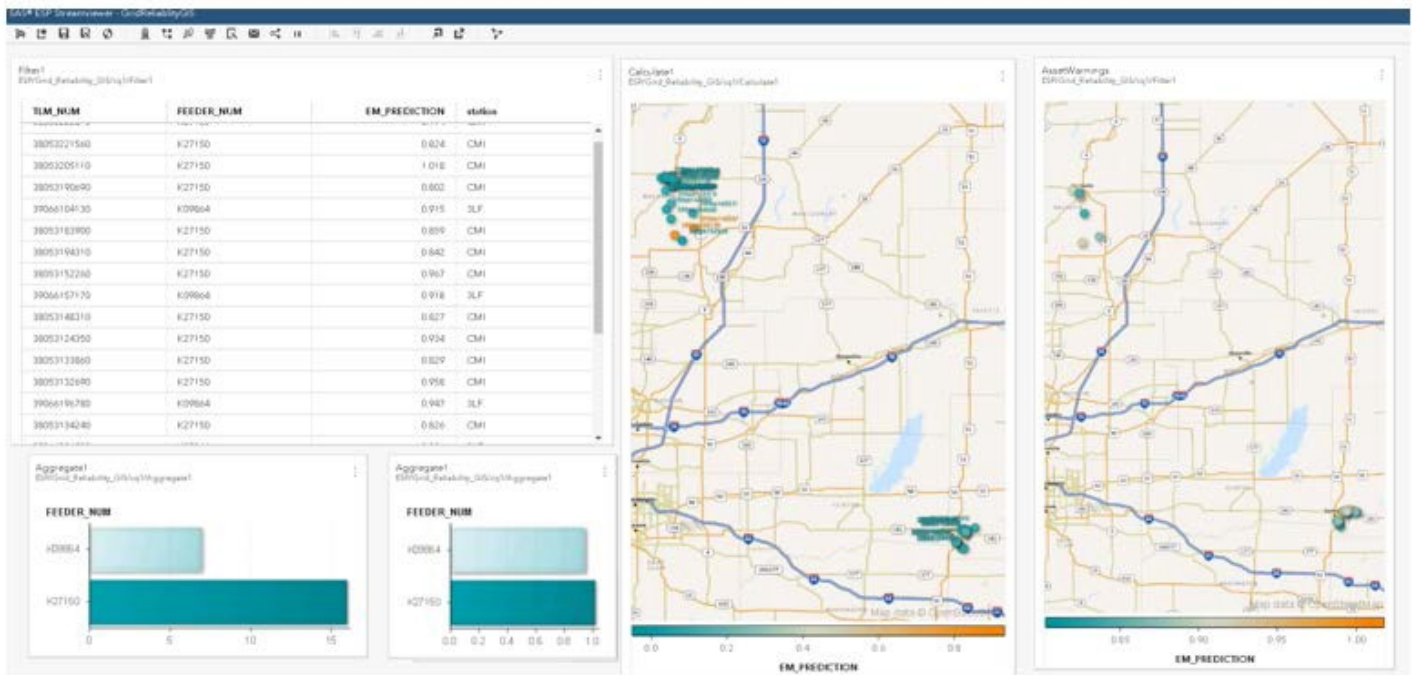**Figure 3.** SAS Event Stream Processing scoring the health of the grid over time

The metrics SAIDI, SAIFI, CAIDI and CAIFI are key performance metrics for an electric utility over time. A brief definition of each is:

1. SAIDI – The System Average Interruption Duration Index is the average outage duration for each customer served.

2. SAIFI – The System Average Interruption Frequency Index is the average number of interruptions that a customer would experience.

3. CAIDI – Customer Average Interruption Duration Index is related to SAIDI and SAIFI and can be viewed as the average restoration time

4. CAIFI – Customer Average Interruption Frequency Index is designed to show trends in customers interrupted and shows the number of customers affected out of the whole customer base.

When blended with historical information, such as interruptions by month over a 5 year history, time of day when the interruption happened and the five year causation histogram, inference can be drawn, based on known factors that allows for the projection that the cause of a power   outage is caused by weather, an animal, vehicle, tree or some other malfunction.

We can then combine the power issue data with geolocation data to show both present issues as well as predict future issues. A forecast for bad weather with the manifestation of high winds can be combined with data regarding the number of trees surrounding or close to power lines to predict when and where we believe a tree will fall on a power line.

6

We can also use the data to determine how to deploy repair crews to address issues in a timelier manner and with greater efficiency. The following, Figure 4, shows the combination of data from a power loss with geolocation information to determine the number of customers affected by the power loss.
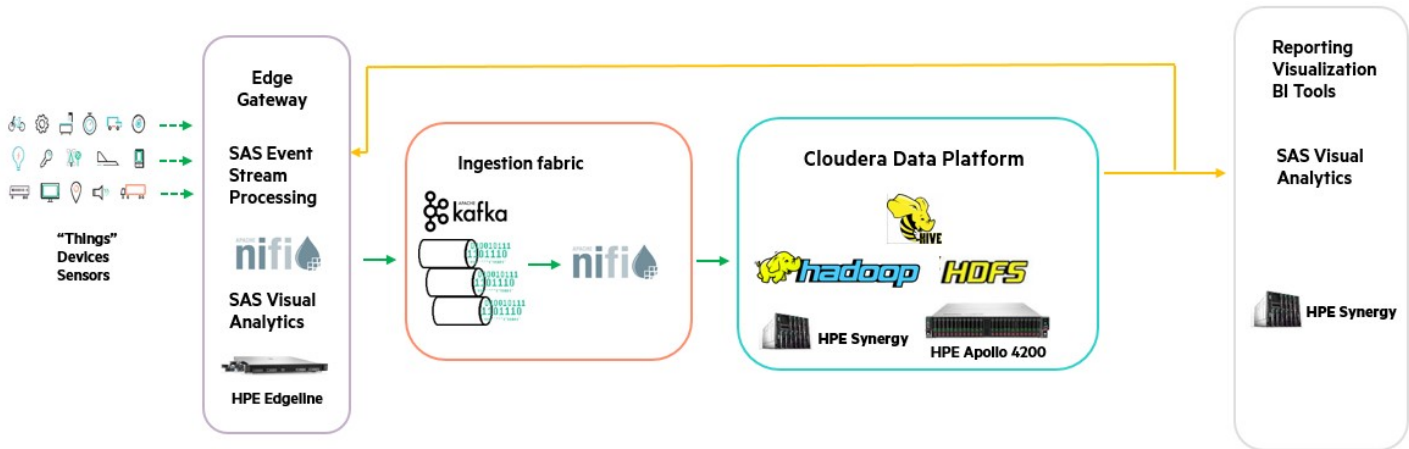


**Figure 4.** SAS Event Stream Processing displaying power loss substations and where they're located on a map

Knowing where a power loss is located so that repair crews can be dispatched is key. But we want to ensure the safety of those crews that have been dispatched. Using augmented reality in the repair vehicle can allow for a crew to see the expected cause of the power outage. They can then be presented with videos regarding how to repair the damage.

All of this lends itself to creating a Smart Grid that is more self-aware and has the ability to respond more rapidly to issues, perhaps addressing them before they become actual issues.

# SOLUTION OVERVIEW

The following, Figure 5, is the edge-to-core data pipeline with the hardware and software components that were used during the POC.



**Figure 5.** The SAS and Elastic Platform for Analytics data pipeline from edge-to-core and back

Hewlett Packard Enterprise and SAS have solutions focused on each level and layer in the solution.

## HPE Edgeline Converged Edge systems

At the edge, Hewlett Packard Enterprise has HPE Edgeline Converged Edge systems to connect to, analyze, and store data. They range from HPE Intelligent gateways, to connect, convert, and perform simple analyses to HPE Edgeline Converged Edge systems that can provide data center-like compute that can be deployed in the harsh environments of the edge. These systems are architected to address increasing performance requirements allowing a customer to choose the specific model to address the unique requirements at each and every data collection point in the process, ultimately providing the ability to run data center apps at the edge.

At the edge, there is SAS Event Stream Processing for inference and alerting. SAS Event Stream Processing can be used for everything from initial collection with a lighter level of analytics and alerting to deeper, more thorough analytics including, but not limited to, control of components on the factory floor, traffic lights and flow control in a Smart City, etc.

Additionally, if required, there is SAS® Visual Analytics which can be deployed at the edge. Visual Analytics allows for visualization of the processes as they are executing and also for some machine learning and deep learning to be executed on the edge, closer to the data collection points. Having AI performed closer to the edge devices themselves helps speed up the feedback loop and allows models used in the data collection and control layer to be cascaded to the control points more rapidly.

As data is moved from the edge to the DC, Hewlett Packard Enterprise provides the Elastic Platform for Analytics (EPA) architecture to address each customer's scaling, performance and storage requirements. Traditional big data environments, historically, have been a one size fits all scenario. This has required customers to over provision storage because they need processing power, or vice versa. The HPE EPA architecture is designed to address that specific issue. No longer is storage directly attached to and scaled along with compute. Instead a customer can scale compute without scaling storage and scale storage without having to scale compute. HPE has recommended configurations within the EPA architecture to address all of a customer's requirements and allows the customer to create an architecture that directly addresses those requirements without ever having to over provision.

The following, figure 6, is a graphical depiction of the EPA building blocks.
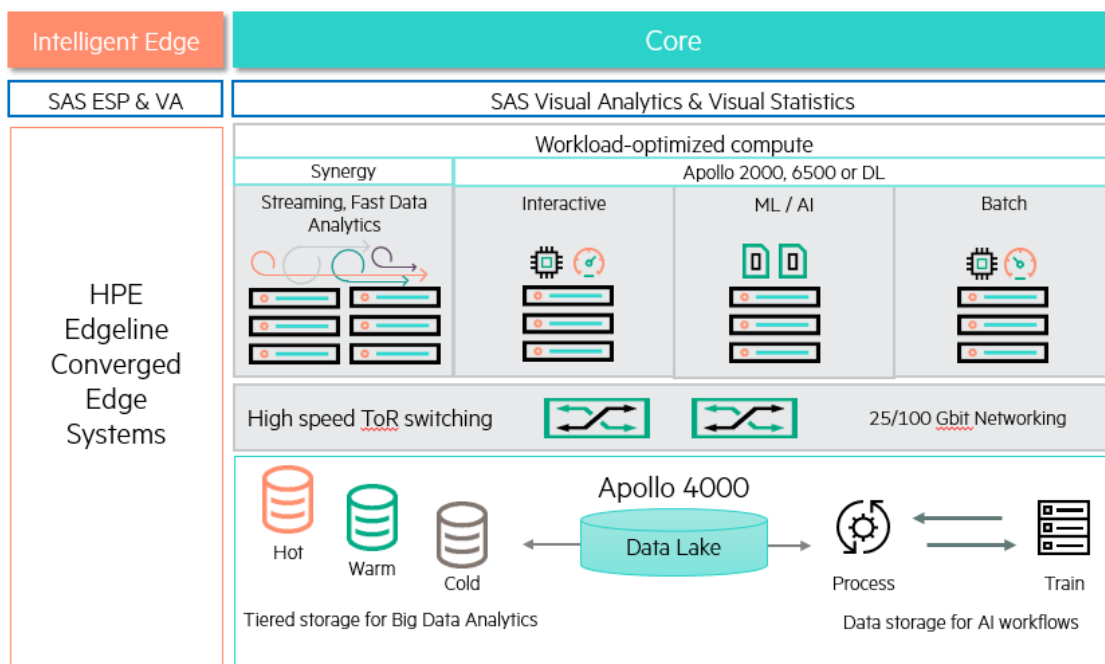


**Figure 6.** Edge-to-core HPE's Elastic Platform for Analytics and SAS software conceptual design

## HPE Synergy Composable Infrastructure for SAS Analytics in the data center/cloud

There are two primary components of the Hewlett Packard Enterprise EPA architecture as they apply to the deployment of SAS software in the  data center/cloud segment of an edge-to-core/cloud scenario. Those elements are:

HPE Synergy is a powerful software-defined solution that allows for the creation of fluid pools of physical and virtual compute storage and fabric   resources that can provisioned and reprovisioned rapidly as resource needs shift. In the HPE EPA architecture, HPE Synergy comprises one of the workload optimized compute alternatives.

HPE Apollo servers are high density rack mount servers that are purpose built to be compute, artificial intelligence or dense storage servers.

For our POC we chose HPE Synergy as the SAS Visual Analytics platform. We did this because of composability, the ability to morph the   environment over time to suit the needs of the application software. A customer doesn't have to plan for what they may need in 1-2-3 years  when they purchase for their initial requirement. They can buy only what they need now and be confident with the ability to add to the configuration at any time in the future.

We also could have chosen an Apollo server as a data center component for SAS Visual Analytics in the data center. While the Apollo 4200 doesn't allow for quite the level of configuration flexibility over time, it does allow for enough elasticity that a customer can be confident that what   they purchase today will not limit their expansion requirements in the future.

## HPE Elastic Platform for Analytics (EPA)

HPE has developed a design for Hadoop infrastructure that, we believe, allows a customer more adjustability when deploying a Hadoop cluster.  HPE calls this design the HPE Elastic Platform for Analytics.

One of the key tenets of the HPE EPA architecture is to place the right component at the right place within the Hadoop cluster such that no   segment needs to be over configured just because another component needs to be enhanced.

For example, if a customer finds they need more compute within a Hadoop cluster, they need not purchase more storage. And the reverse is true also. Each unit within the architecture is aimed at the specific need that it is intended to satisfy.

If, on the other hand a dual requirement exists, it is certainly possible to collocate compute and storage on the same node.

## HPE Elastic Platform for Analytics (EPA) Analytics block for Big Data

The Analytics block is used for compute intensive processing. For example, if a requirement exists to collocate SAS on the Hadoop cluster, then an analytics block should be purchased.

The default recommended configuration for the Analytics block is the HPE Synergy 480 Gen10 compute module. The HPE Synergy 480 Gen10 is a two processor system. For customers requiring additional compute throughput the HPE Synergy 660 Gen10 compute module is also an option. The HPE Synergy 660 Gen10 compute module is a four processor system.

The variability within HPE Synergy compute elements is in the number of and clock speed of the processors, the number of cores per processor, and the amount of memory per node. Hewlett Packard Enterprise recommends configuring memory in increments of 6 DIMMs per processor of the Intel® Scalable family of processors. This is to maximize memory access speeds.

Local storage for the Analytics block is provided by the HPE Synergy D3940 Storage Module. Each storage module has 40 drive slots, with the ability to insert both SSD drives as well as spinning media.

Because of the ability to compose the solution to fit the need, we can install as many HPE Synergy D3940s in each individual frame as needed in order to satisfy the requirement. We can also vary the number of and size of the drives to be placed in the Storage Module. The HPE Synergy D3940 has a capacity of 40 drives. Those drives may be attached to any compute module installed in the same frame as the HPE Synergy D3940. This composability benefit allows the Synergy environment to be changed if the application framework requirements change and additional drives need to be presented to compute nodes. Removing the constraints of a certain number of drive slots per server allows for a much more fluid architecture. If it is found that more storage is needed on the server, reconfiguration is as simple as reassigning a profile and then rebooting.

The operating system is provided to each of the Compute Modules by HPE Image Streamer. This allows all of the storage, both in the Compute Modules, as well as in the Storage Modules to be used for application deployment purposes, because none of the space required by the OS is consumed on the storage.

## HPE Density-Optimized Storage block

A storage block within the Hadoop cluster is used to store information that's been collected by SAS Event Stream Processing and forwarded to the Hadoop cluster to be stored in Hive tables for later access.

When a large amount of storage is required, a customer can opt for one or more Density-Optimized Storage blocks. This block is based on the HPE Apollo 4200 Gen10 server, typically configured with 28 SATA large form factor drives.

Variability in this block is in the number of processors, typically two (a customer may opt for only one processor), the clock speed of the processor, and number of cores per processor. This block starts out with 256GB of memory, but this is also customizable by the customer from 128GB to 768GB. The number and size of the hard disk drives is a choice left up to the customer.

## SAS software in the data center

In the DC, SAS Visual Analytics and SAS Visual Statistics are available for deeper AI analysis, SAS® Visual Data Mining, and Machine Learning. Because there is a great deal more processing power available in the DC, SAS software is able to more thoroughly analyze data coming from the edge. This enables machine learning and deep learning algorithms to become more accurate over time.

# SOLUTION COMPONENTS

## HARDWARE AT THE EDGE

### SAS Event Stream Processing and SAS Visual Analytics at the edge

At the edge, we placed an HPE Edgeline EL4000 Converged Edge system. The HPE Edgeline EL4000 can be configured with 1 to 4 server cartridges. Each server cartridge runs an industry standard operating system, such as Red Hat Enterprise Linux® or Microsoft Windows Server®. Figure 7 is a picture of the HPE Edgeline EL4000.

The HPE Edgeline EL4000 Converged Edge system comes in 4 chassis types. They are:

- HPE Edgeline EL4000 10GbE Switch System
- HPE Edgeline EL4000 10GbE 2xSFP+ v2 Switch System
- HPE Edgeline EL4000 4x10GbE 2xQSFP+ v2 Pass Thru System
- HPE Edgeline EL4000 10GbE 2xSFP+ Switch PXIe System

The HPE Edgeline EL4000 has the following specialty card options:

HPE Networking:

- HPE Ethernet 10/25Gb 2-port 640SFP28 Adapter
- HPE InfiniBand FDR/Ethernet 10Gb/40Gb 2-port 544_QSFP Adapter
- HPE InfiniBand EDR/Ethernet 2-port 841QSFP28 Adapter

HPE GPU Accelerator:

- HPE NVIDIA Tesla P4 8GB Computational Accelerator
- HPE AMS Radeon Pro WX4100 Graphics Accelerator

The server cartridges that are available in the HPE Edgeline EL4000 are the m510 and the m710x.



**Figure 7.** HPE Edgeline EL4000 Converged Edge system

### HPE ProLiant m510 server cartridge

The HPE ProLiant m510 server cartridge, as shown in figure 8, is available with either an eight core or a sixteen core Intel Xeon processor and up to 128GB of memory. When paired with the available Nvidia Tesla P4 GPU, this server cartridge can perform inference operations at the edge, allowing rapid enhancement of the models being executed for capture, control and alerting. SAS Event Stream Processing and SAS Visual Analytics have been validated to work on the m510 Server cartridge.



**Figure 8.** HPE ProLiant m510

## HARDWARE IN THE CORE

### Hadoop cluster

Once the data has been initially reviewed at the edge, coalesced and transferred to the core data center it was placed into a Hadoop cluster. The Hadoop cluster that we used was built using the HPE Elastic Platform for Analytics design.

The HPE Elastic Platform for Analytics is a design specification that a customer can use to build a custom Hadoop cluster specifically focused on the special needs of SAS Visual Analytics.

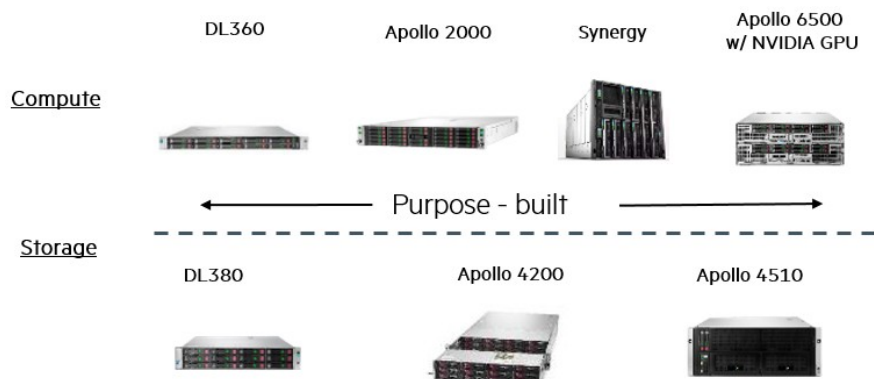Figure 9 depicts building block choices used to build the custom deployment model.



**Figure 9.** HPE Building blocks for assembling a custom SAS core architecture

For the POC we utilized a small Hadoop cluster. The compute portion of the Hadoop cluster used HPE ProLiant DL360 Gen8 servers while the storage portion utilized Apollo 4200 servers[3].

The Hive process was run on one of the compute nodes, an HPE ProLiant DL360 Gen8 server. Kafka was run on one of the data nodes, an Apollo 4200 server.

---

[3] A table of systems and the software components that were run on them is contained in Appendix A.

For more information on the HPE Elastic Platform for Analytics see the HPE Reference Configuration for Elastic Platform for Big Data Analytics   white paper, https://h20195.www2.hpe.com/v2/getdocument.aspx?docname=4aa6-8931enw. For information on sizing a Big Data Hadoop   cluster using the HPE EPA architecture, please see the HPE Sizer for the Elastic Platform for Big Data Analytics, http://www.hpe.com/info/sizers.

## SAS Visual Analytics

To support the requirements of SAS Visual Analytics the compute block should support powerful processors and fast DDR4 memory based on  computational need. For optimal performance, it is recommended datasets to be analyzed fit entirely into memory. This requires a compute   building block that is optimized for in-memory analytics and has a larger memory footprint then a compute block for general purpose workloads.   For this requirement, the recommendation is to utilize an EPA compute block based on either the Apollo 2000 or the HPE Synergy server line.

The decision between the HPE Apollo 2000 and the HPE Synergy is one of the specific use case. HPE Synergy has the unique capability to   morph itself so that it can be configured specifically to be used with SAS Visual Analytics processing during a specific time period and then have   an entirely different personality and configuration allowing it to be peaked for use during another time period. HPE calls this capability composability.

For instance, many SAS customers use SAS® 9.4 to derive raw data that is then visualized using with SAS Visual Analytics. In the past these  customers had to purchase two servers, one to run SAS 9.4 and one to run SAS Visual Analytics. By using the composability aspect of HPE  Synergy, a single compute module could be used for both needs.

The HPE Synergy compute modules works with server profiles. A server profile tells the compute module what storage and networking it's going   to use. All that's needed is to create 2 separate profiles, one for SAS 9.4 and one for SAS Visual Analytics. The major difference between the two is the storage to which each profile is directed. However, the amount of network bandwidth required may comprise an ancillary difference. The different storage then has a different OS instance and different SAS software installed. Conversion from one profile to another takes   approximately 10 minutes. In this manner HPE Synergy could use a SAS Visual Analytics profile during the day for user interactive processing   and a SAS 9.4 profile during the evening when the datasets for the next day's Visual Analytics processing are being generated.

Because the amount of data collected at the edge is very large, even after being condensed and transferred to the data center, when combined   with all of the edge data and stored for a period of time, the storage requirements will be quite large. Many customers will opt for storing IoT data   in the core on a Hadoop cluster of servers.

SAS Visual Analytics and SAS 9.4 have the ability to work in concert with all of the major distributions of Hadoop and use it as both a data   repository and also a data source.

## SOFTWARE

The following, table 1 is a list of software used in our environment.

**Table 1.** Software

| Software | Version | Where used |
|---|---|---|
| SAS Event Stream Processing | 5.2 | Edge |
| SAS Visual Analytics | 8.2 | Edge and data center |
| SAS Viya | 3.3 | Edge and data center |
| Cloudera Data Platform | 2.6.4 | Data center |
| Nifi | 1.9.1 | Data center |
| Kafka | 1.1.1 | Data center |
| Red Hat Enterprise Linux | 7.5 | Edge and data center |

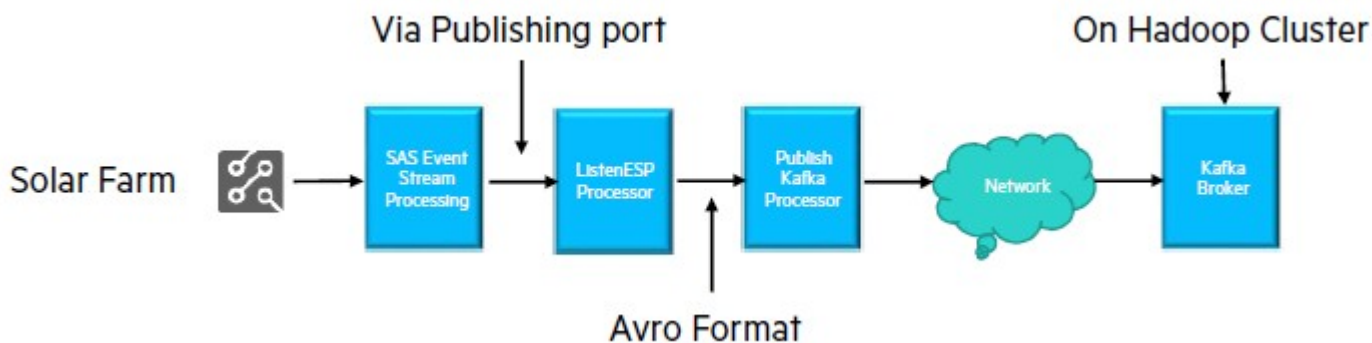## PUTTING IT ALL TOGETHER

### At the edge

Into our HPE Edgeline EL4000 we placed two m510 server cartridges. One cartridge was configured one Intel D-1548 8-core processor running at 2.0GHz, 128GB of memory and two 1TB NVMe drives. The other cartridge was configured with one D-1547 16-core processor running at 1.7GHz, 128GB of memory and two 1TB NVMe drives.

On one of the cartridges, we installed SAS Viya Visual Analytics version 8.3.1 and SAS Event Stream Processing (ESP). We installed both Viya Visual Analytics and Event Stream Processing because we wanted the ability, provided by Visual Analytics, to analyze the data generated by Event Stream Processing at the edge as soon as it was being generated.

On the other cartridge we installed Event Stream Processing only.

We also installed NiFi on the edge servers. SAS Event Stream Processing creates data windows by publishing network ports that can be monitored by other processes. We used the SAS NiFi integration processor named ListenESP[4] to watch the data window being published by Event Stream Processing as it monitored energy being generated by a solar farm. The ListenESP processor moves the data in Avro format[5]. We converted the data to a Kafka format using the PublishKafka processor and pushed it to a Kafka broker located on the Hadoop cluster.

Once on the Hadoop cluster the data was persisted in a Hive table for use by the SAS Visual Analytics deployment on the HPE Synergy 480 Gen10 serverThe following figure 10 depicts the data flow from Event Stream Processing to the Kafka Broker located on the Hadoop cluster.



**Figure 10.** Data flow from Solar Farm to Event Stream Processing to NiFi and then via the network to the Kafka broker located on the Hadoop cluster in the data center

### In the data center

We placed our Hadoop cluster in the data center (DC). As shown in Appendix A, the Hadoop cluster was made up of three HPE ProLiant DL360 and three HPE Apollo 4200 servers. The HPE ProLiant DL360 servers were used as the management server, the Name Node and the Secondary Name Node. The Apollo 4200 Gen9 servers were used as data nodes within the cluster.

We utilized Cloudera Data Platform (HDP), NiFi and Kafka on the Hadoop cluster.

Once the data has been posted to the Kafka Broker, which we located on one of the data nodes within the Hadoop cluster, we used NiFi to extract the data and post it to Hive.
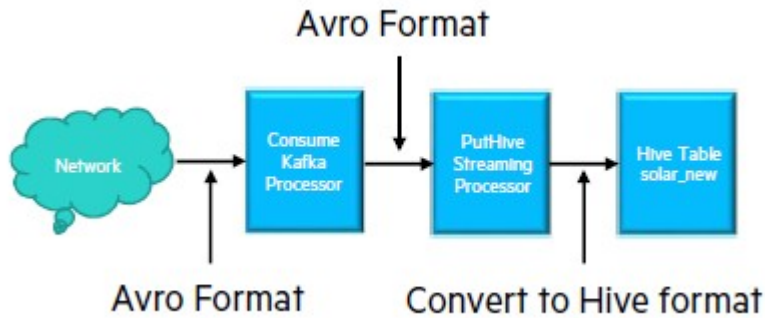
---

[4] All NiFi development palettes are located in Appendix B

[5] Avro is an open source Apache project. Avro relies on schemas written in JSON format which defines the fields and their types. The data is then stored in a binary format which makes it compact and efficient. For more information on the Avro format see http://aseigneurin.github.io/2016/03/04/kafka-spark-avro-producing-and-consuming-avro-messages.html

To pull the data from Kafka, we used the ConsumeKafka processor. The ConsumeKafka processor took the data that was posted to Kafka with the PublishKafka processor and pushed it to the PutHiveStreaming processor. It was the PutHiveStreaming processor that posted the data to the Hive table.

In order for the PutHiveStreaming processor to be able to post the data to a Hive table, that table must have been created. In our environment we used the default database and created the table within that database[6].
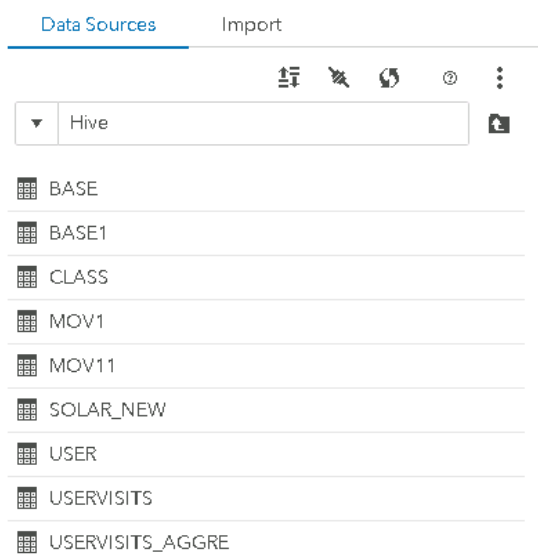
The following, figure 11 is a graphical depiction of the data flow from the network through to the Hive table within the default database.



**Figure 11.** Data flow from the network to the Hive table solar_new via the ConsumeKafka and PutHiveStreaming processors

Once we have the data collocated in a Hive table, it's a simple task of accessing that table with SAS Visual Analytics running in the data center. To do this we use the SAS/ACCESS® Interface to Hadoop. Once we've configured the SAS/ACCESS to Hadoop, we are provided with a list of Hive tables to which SAS Visual Analytics has access. The following, figure 12, shows the table listing when deciding which table to load into memory within SAS Visual Analytics.



**Figure 12.** Listing of available tables that can be loaded in to memory within SAS Visual Analytics

---

[6] Table definition script can be found in Appendix C.

In the above listing, the table SOLAR_NEW is the one that we're interested in. But the other tables listed are available for upload also. Basically,      any table being stored by Hive is available, so long as the user has read access to that table. Once the table has been loaded into memory, all SAS  Visual Analytics algorithms are available to operate on that data.

## BEST PRACTICES AND CONFIGURATION GUIDANCE FOR THE SOLUTION

The HPE Edgeline EL4000 that was used during the testing had 2 x 10Gb network connections available per cartridge. We used only one of the  network ports and as you will see in the performance section, only a small amount of the network bandwidth was consumed. Customers may opt  for a 1Gb network uplink, which the 10Gb network connections can use. However care is needed to ensure the network is not a constriction point  in a deployment scenario.

## CAPACITY AND SIZING

### WORKLOAD DESCRIPTION

The following, figure 13 is a graphical depiction of the environment that was tested. We used the HPE Synergy as the driver for the tests. We   utilized a file to simulate sensor monitoring. The file had 241 million events that were published from the HPE Synergy 480 Gen10 system to a  SAS Event Stream Processing engine running on the HPE m510 Converged Edge server cartridge. We were able to scale the number of events   per second from 10K per second up to in excess of 1.26M events per second. This was done by running multiple instances of SAS Event Stream  Processing on the HPE m510 Converged Edge Server cartridge.

Once the event was received on the first HPE m510 server cartridge, it was formatted and transmitted to SAS Event Stream Processing running   on the second HPE m510 server cartridge. It was these instances of SAS Event Stream Processing that performed most of the work.

Once the instance of SAS Event Stream Processing running on the second HPE m510 server cartridge had performed the analysis and alerting   based on the data, it published that data to a processing window that was monitored on back on the HPE Synergy 480 Gen10 system.

It should be noted that during this test we did not utilize the NiFi to Kafka data flow. The test was meant to demonstrate the scalability of SAS   Event Stream Processing and the impact that software made on the m510 cartridges. We did, however, test the proof of concept and validated   the total edge to core/cloud data pipeline outside of this performance testing.
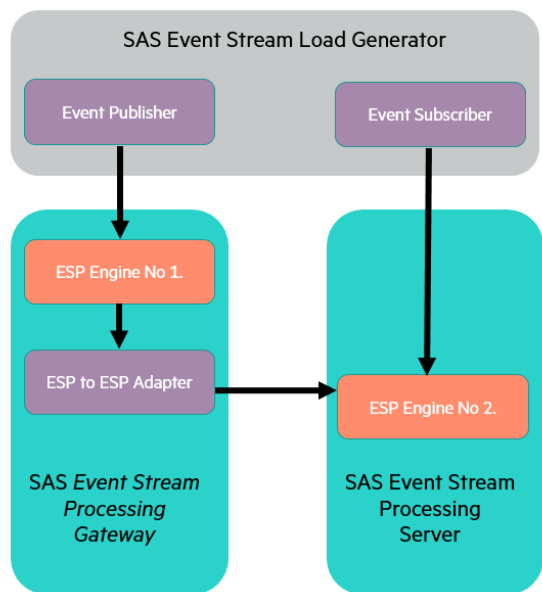


**Figure 13.** Graphical depiction of the SAS Event Stream Processing test environment

The test that we ran is intended to emulate phasor Measurement Unit (PMU) monitoring in a Smart Grid environment.

A phasor measurement unit (PMU) is a device used to estimate the magnitude and phase angle of an electrical phasor quantity like voltage or current in the electricity grid using a common time source for synchronization. Time synchronization is usually provided by GPS and allows synchronized real-time measurements of multiple remote measurement points on the grid. PMUs are capable of capturing samples from a waveform in quick succession and reconstruct the phasor quantity. The resulting measurement is known as a synchrophasor. These devices can also be used to measure the frequency in the power grid. A typical commercial PMU can report measurements with very high temporal resolution in the order of 30-60 measurements per second. [7]

Since a typical PMU will only send measurements 30 to 60 times every second, scaling up the test provides information about how a single instance of SAS Event Stream Processing running on HPE Converged Edge Systems could monitor multiple PMUs simultaneously.

All of the tests being reported consumed the entire 241 million events and ran for approximately 40 minutes.

The following, figure 14 shows the peak network bandwidth consumed by system number 1. We measured peak network bandwidth because we wanted to show that at no point was there a network bottleneck. Network bandwidth is an important metric because, if the bandwidth is exceeded, we face the possibility of artificially limiting the amount of work the system can perform. In order to perform its work, a system needs to get data to and from the processors. If that data is coming via the network, as it was during our tests, then a network bottleneck would limit the amount of data flowing in to and out of a system, which in turn would limit the amount of work able to be performed.

During the test, network utilization was sampled at 2 second, 10 second and 40 second intervals. We chose the 2 second interval, because it showed the highest network utilization for all of the tests. We also show the number of events being processed per second so you can see the scaling of SAS Event Stream Processing and the first HPE m510 server cartridge.

It should be noted that the values are given in megabits per second. We have two 10Gb network ports, of which only one was used. So the available bandwidth is 10,000 megabits per second.
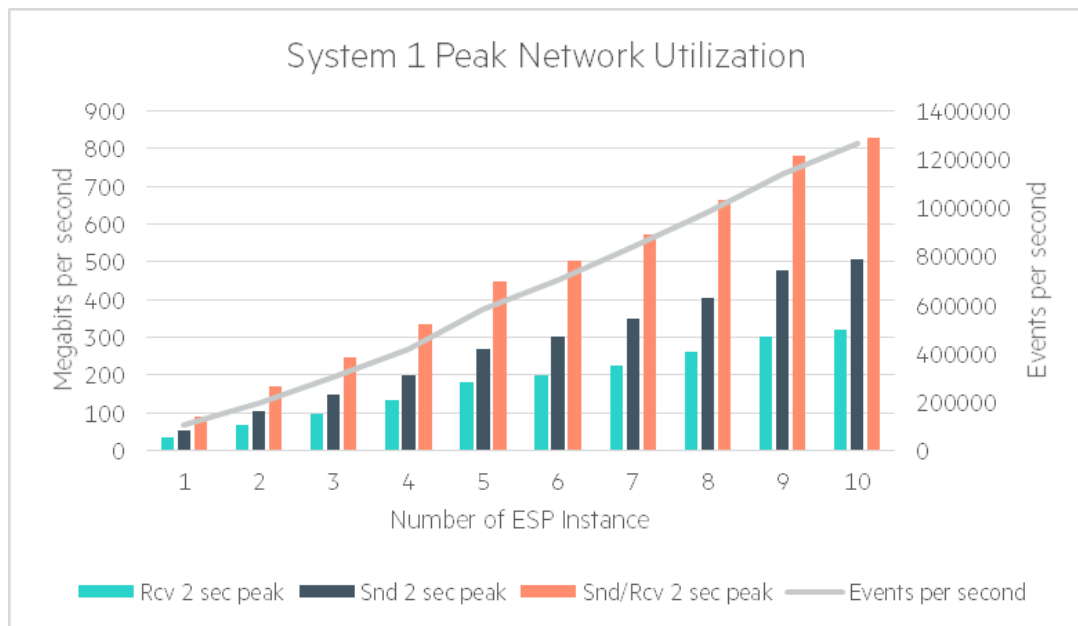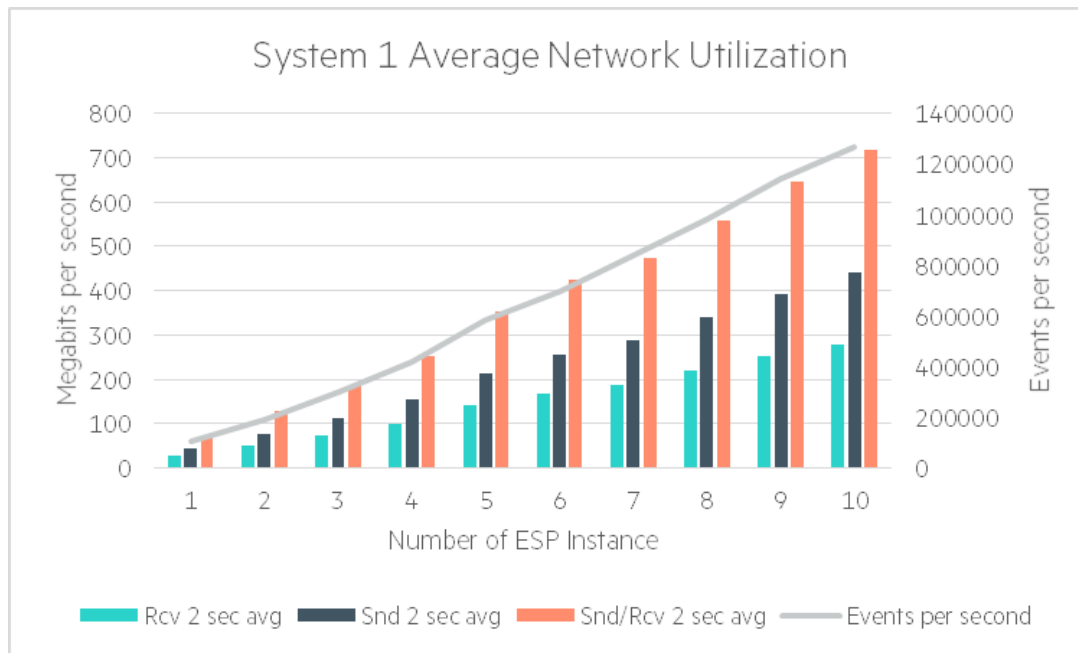


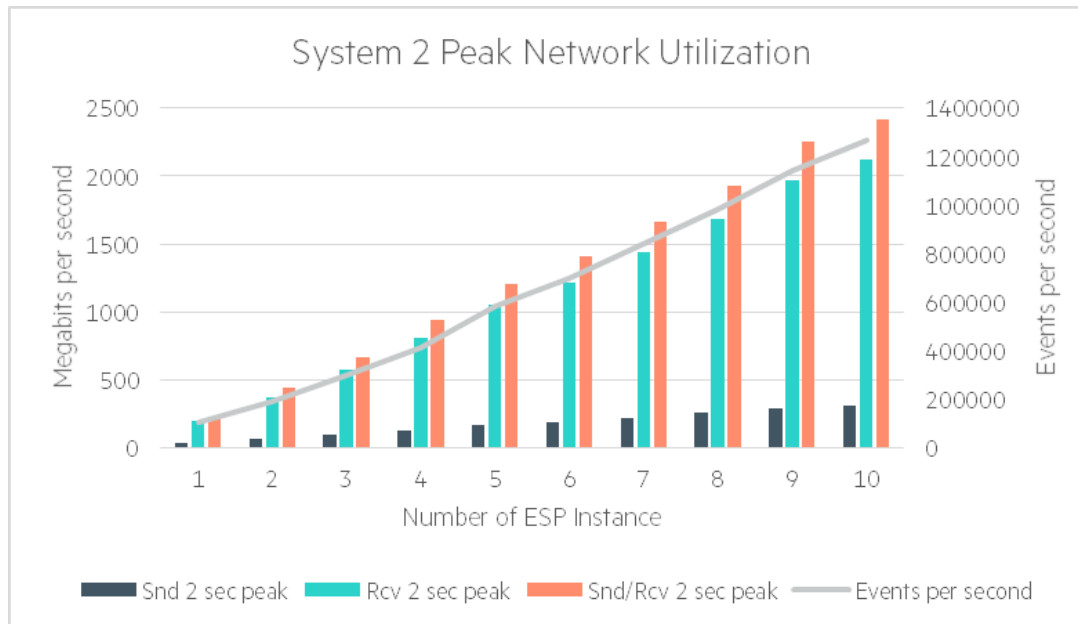**Figure 14.** Peak network bandwidth consumed on system 1 during the tests

[7] Reprinted from Wikipedia Phasor measurement unit article: https://en.wikipedia.org/wiki/Phasor_measurement_unit

The following, figure 15 is the same type of graph, however, this one shows the average network utilization when running the tests.



**Figure 15.** Average network bandwidth consumed during the tests

The following, figure 16 shows the peak network utilization on system number 2. Note that system 2 did most of the work and its network consumption reflects that. Even though we consumed almost 2.5Gb of peak network bandwidth, there remains more than 7.5Gb of available bandwidth available, should it be needed.



**Figure 16.** Peak network bandwidth consumed on system number 2 during the tests

And finally, figure 17, shows the average network bandwidth consumed during the tests.
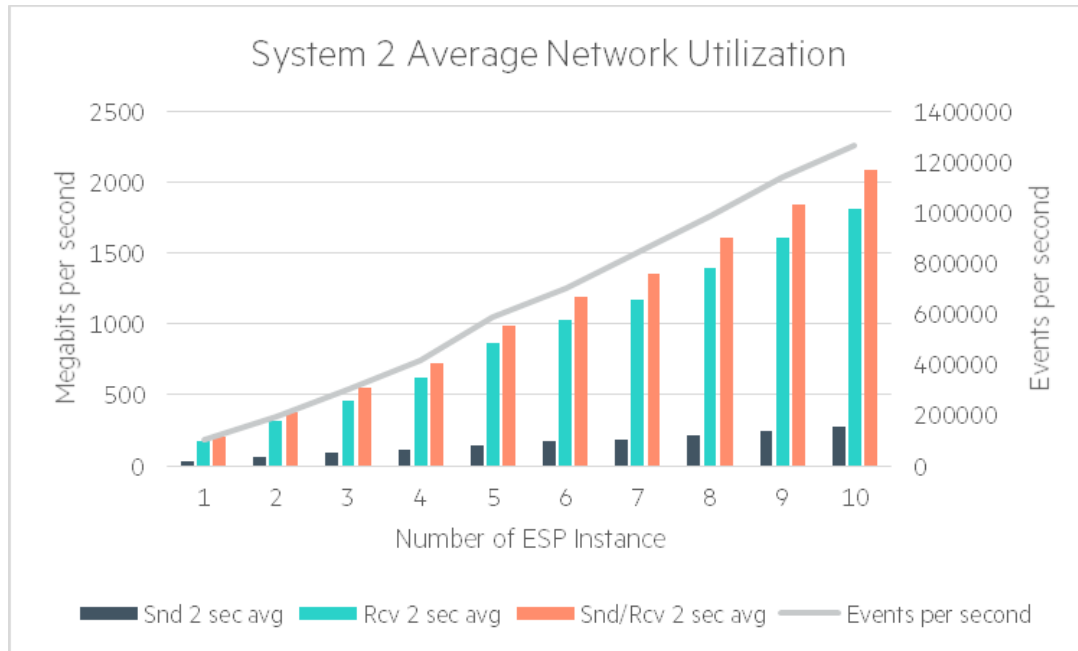


**Figure 17.** Average network bandwidth consumed during the tests

Let's look at the system CPU utilization during these various tests. This is important to understand to determine how much further the system could be stressed by placing more load. Figure 18 shows the system utilization levels during the various tests.
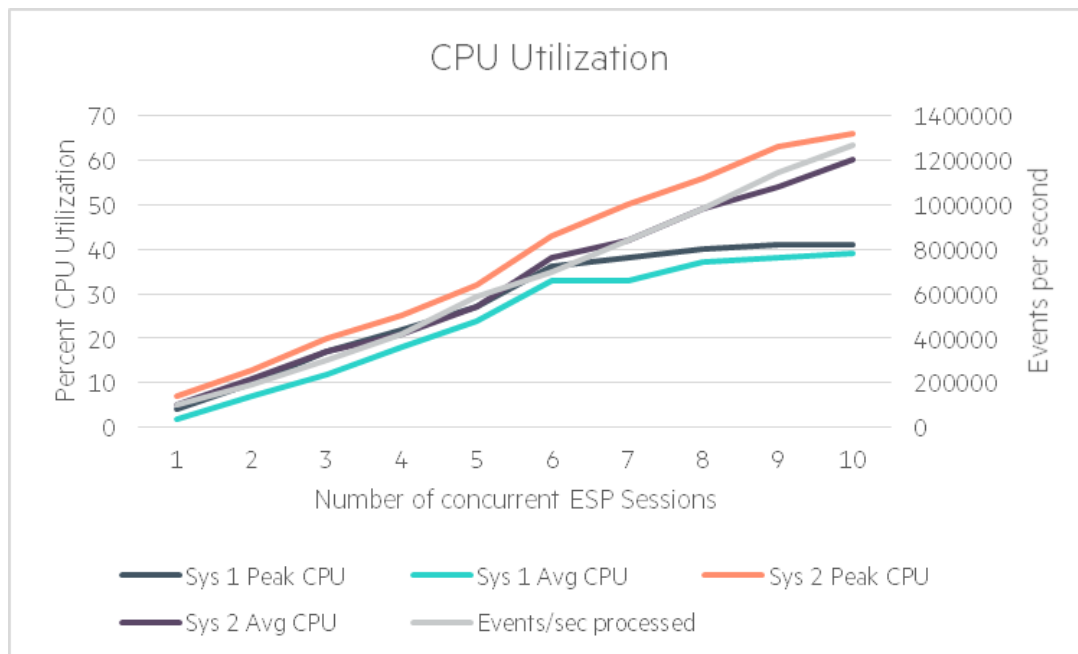


**Figure 18.** Peak and average system utilization during the 10K, 50K, 100K and 200K event per second tests

## ANALYSIS AND RECOMMENDATIONS

Our testing has shown the ability of both HPE and SAS to scale from a small, single device monitoring solution up to a full scale data acquisition and analytics server. The number of events per second scaled in a nearly linear fashion, which indicates the ability of both SAS Event Stream Processing and the HPE Edgeline Converged Edge systems ability to scale. While we terminated testing with 10 simultaneous instances of SAS Event Stream Processing, it should be noted that there was still available network bandwidth, as well as CPU available. Based on the linearity of the CPU utilization and the number of events per second being processed, we estimate this configuration would have been able to support almost 2 million events per second when fully consumed.

If a customer wishes to deploy fewer resources at the edge, they may choose the EL1000 with a single cartridge. And since the EL1000 utilizes the same family of cartridges as the HPE Edgeline EL4000, they can tailor the solution for their specific requirements.

During our testing, based on the above graphs, we have shown that the server cartridges have the ability to scale beyond the number of events per second that we tested. There is still quite a bit of available network bandwidth remaining. In our tests, we never exceeded 2500Mb/sec (250MB/sec) on a link that is capable of providing 10Gb/sec (1GB/sec). Additionally, our system came with 8 x 10Gb QSFP+ pass-through networking ports that are divided equally among the four cartridge slots effectively doubling the 10Gb/sec bandwidth.

While we utilized 2 of the available 4 cartridge slots during our testing, if further compute resources are required an additional 2 cartridges can be placed in the HPE Edgeline EL4000 frame.

Scaling SAS' Event Stream Processing software at the edge is as easy as deploying a second, third, fourth, etc. instance of Event Stream Processing on the same cartridge. Additionally SAS Event Stream Processing can be placed on any numbers of servers within a customer's environment.

In the data center, scaling SAS Viya and SAS Visual Analytics on the HPE Synergy 480 Gen10 can take two forms.

We have the HPE Synergy 660 Gen10 which is a 4 processor system which effectively doubles the computing horsepower for SAS Viya.

SAS also provides SAS Viya and SAS Visual Analytics in a MPP model where we deploy multiple servers in a tightly coupled, clustered environment. In this scenario, one of the servers becomes the controlling node. The rest of the servers are worker nodes. When a query is requested, the controller node distributes that query among the worker nodes based on the number of worker nodes and the distribution of the data within the query set. As the worker nodes complete their processing they pass their result set back to the controller node, whose job is to combine the data from all of the worker nodes and present the final model to the user requesting the query.

Because HPE Synergy is built from the ground up to be flexible and composable, all that's required to add a server into the equation is to locate an empty slot in the HPE Synergy Frame and install the new server. Then create a profile where fabric bandwidth and local storage is defined along with an OS image and bring up the server. At this point SAS software can be installed and the cluster can be configured to accept the new member.

Finally, HPE's Elastic Architecture for Analytics allows for ease of scaling the Hadoop cluster. If additional storage capacity is required, acquire as many storage blocks as required and include them in the Hadoop cluster. If additional compute capacity is desired, purchase additional compute capacity. Using HPE's EPA, a customer never has to pay for one resource, for example storage or compute to get the other.

## SUMMARY

At the edge, HPE has a large portfolio that allows the selection of hardware to fit each unique requirement. The product portfolio starts with lightweight single processor, dual core systems and the product portfolio scales all the way up to 64-cores in a single chassis. Additionally, the larger HPE Edgeline Converged Edge systems allow for analytical processing engines (referred to as GPUs) to be added allowing for inference modeling to be performed right at the edge.

This breadth of offering allows a customer to acquire a server that fits both budget and workload profile.

When combined with SAS software both the edge and in the data center, the total offering provides the ability to capture and react to massive   amounts of edge data, the ability to move that data to the data center and then analyze the data in order to provide enterprises with the type of   information required to make both reactive as well as planned business decisions.

The ability to scale systems and software using HPE's Elastic Platform for Analytics allows customers to place the right resource in the right place   at the right time.

**IMPLEMENTING A PROOF-OF-CONCEPT**

As a matter of best practice for all deployments, Hewlett Packard Enterprise recommends implementing a proof-of-concept using a test  environment that matches as closely as possible the planned production environment. In this way, appropriate performance and scalability   characterizations can be obtained. For help with a proof-of-concept, contact an HPE Services representative hpe.com/us/en/services/consulting.html) or your HPE partner.


## RESOURCES AND ADDITIONAL LINKS

HPE Reference Architectures, hpe.com/info/ra

HPE Servers, hpe.com/servers

HPE Storage, hpe.com/storage

HPE Networking, hpe.com/networking

HPE Technology Consulting Services, hpe.com/us/en/services/consulting.html

HPE Reference Configuration deploying a composable SAS® infrastructure with HPE Synergy and HPE Image Streamer, https://h20195.www2.hpe.com/V2/GetDocument.aspx?docname=a00061618enw

HPE Reference Architecture for SAS® 9.4 on HPE Synergy and HPE 3PAR 8400, https://h20195.www2.hpe.com/V2/GetDocument.aspx?docname=a00044785enw

HPE Reference Architecture for SAS® 9.4 on HPE 3PAR 8400 and HPE Synergy, https://h20195.www2.hpe.com/V2/GetDocument.aspx?docname=a00044784enw

To help us improve our documents, please provide feedback at hpe.com/contact/feedback.


## CONTACT INFORMATION

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

Mark Barnum
Hewlett Packard Enterprise
585.236.6510
Mark.barnum@hpe.com

Kannan Mani
Hewlett Packard Enterprise
650.258.1678
Kannan.mani@hpe.com

# APPENDIX A: BILL OF MATERIALS

**Table 1a.** Hardware used in the POC and software that ran on that hardware

| System | Components | Software |
|---|---|---|
| **HPE Edgeline EL4000** | | |
| HPE ProLiant m510 server cartridge | 1 x Intel Xeon D-1548 8-core 2.0 Ghz Proc | SAS Viya |
| | 128GB Memory | SAS Event Stream Processing |
| | 1 x 120GB SATA m.2 SSD | SAS Visual Analytics |
| | 2 x 1TB Read Intensive NVMe drives | |
| HPE ProLiant m510 server cartridge | 1 x Intel Xeon D-1587 16-core 1.7 Ghz Proc | SAS Event Stream Processing |
| | 128GB Memory | |
| | 1 x 120GB SATA m.2 SSD | |
| | 2 x 1TB Read Intensive NVMe drives | |

| System | Components | Software |
|---|---|---|
| **HPE Synergy** | | |
| HPE Synergy 480 Gen10 | 2 x Intel Xeon Gold 6144 8-core 3.5GHz Procs | SAS Viya |
| | 768GB Memory | SAS Visual Analytics |
| | OS via Image Streamer | |
| HPE Synergy D3940 | 1 x RAID-5 2.4TB logical drive | SAS Software |
| | 4 x 800GB drives RAID-5 (3+1) | SAS XML and data files |
| **Hadoop Cluster** | | |
| HPE ProLiant DL360 | 2 x E5-2670 8-core 2.6GHz Procs | Hadoop Ambari console |
| (management node) | 128GB Memory | |
| | 1 x 1.2TB Logical drive | |
| HPE ProLiant DL360 | 2 x E5-2670 8-core 2.6GHz Procs | Hadoop Name Node |
| | 128GB Memory | Hive |
| | 1 x 1.2TB Logical drive | |
| HPE ProLiant DL360 | 2 x E5-2670 8-core 2.6GHz Procs | Hadoop Secondary Name Node |
| | 128GB Memory | |
| | 1 x 1.2TB Logical drive | |
| Apollo 4200 | 2 x E5-2660 v4 14-core 2GHz Procs | Kafka |
| | 128GB Memory | Hadoop data node |
| | 27 x 4TB drives | |
| | 1 x Dual 128 VU m.2 drives (boot) | |
| Apollo 4200 | 2 x E5-2690 v3 12-core 2.6 GHz Procs | Hadoop data node |
| | 256GB Memory | |
| | 48 x 600GB drives | |
| Apollo 4200 | 2 x E5-2690 v3 12-core 2.6 GHz Procs | Hadoop data node |
| | 256GB Memory | |
| | 45 x 600GB drives | |

# APPENDIX B: NIFI DATA FLOWS

## NIFI DATA FLOW DEVELOPMENT PALETTE ON THE HPE EDGELINE EL4000

The following, figure 19 is the NiFi development palette after the data flow was established between the Edgeline m510 server cartridge and the Hadoop cluster.
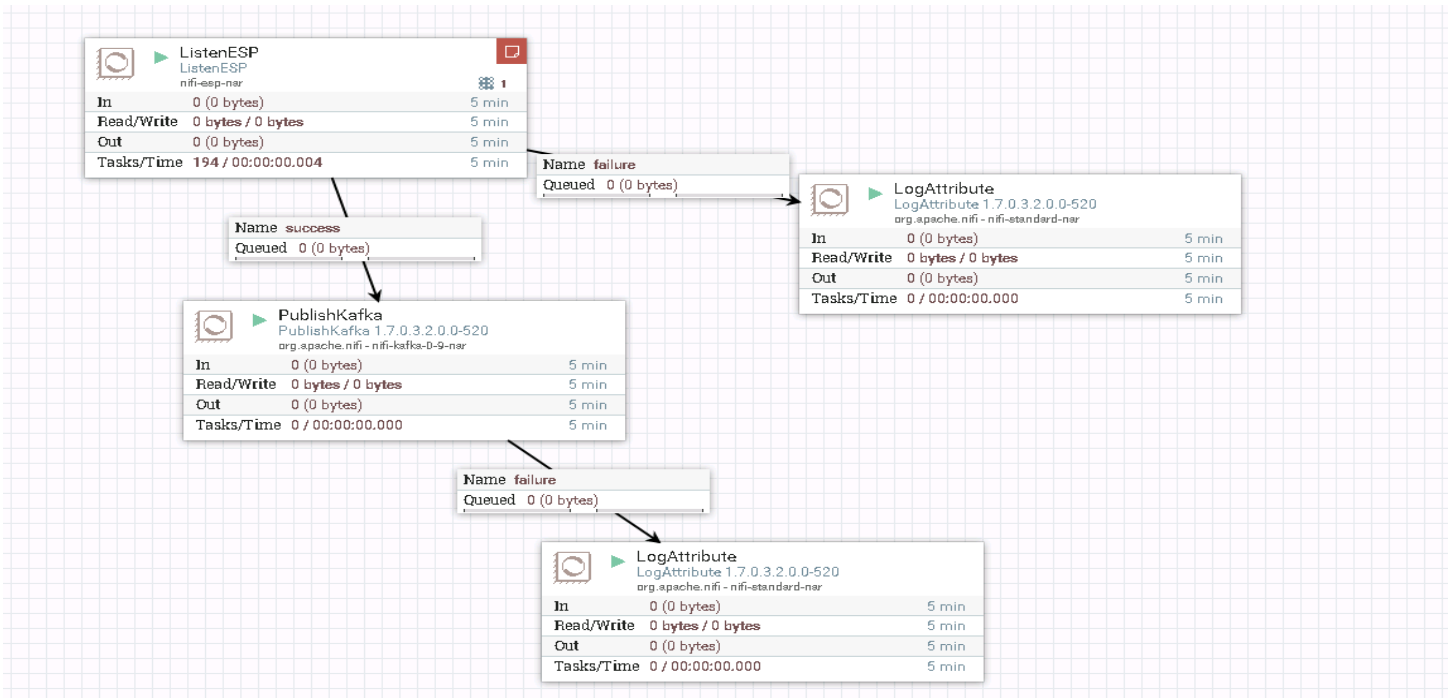


**Figure 19.** The development template for the data flow taking data from SAS Event Stream Processing and posting that data to the Kafka broker running on the Hadoop cluster

As you can see, as we listen to Event Stream Processing, if we happen to have a failure, we log that failure. Success moves the data to the next processor, PublishKafka. There, if we happen to have a failure, we log that failure, otherwise for either processor, no additional steps are taken if we are successful in our operation.

### NIFI DATA FLOW DEVELOPMENT PALETTE ON THE HADOOP CLUSTER

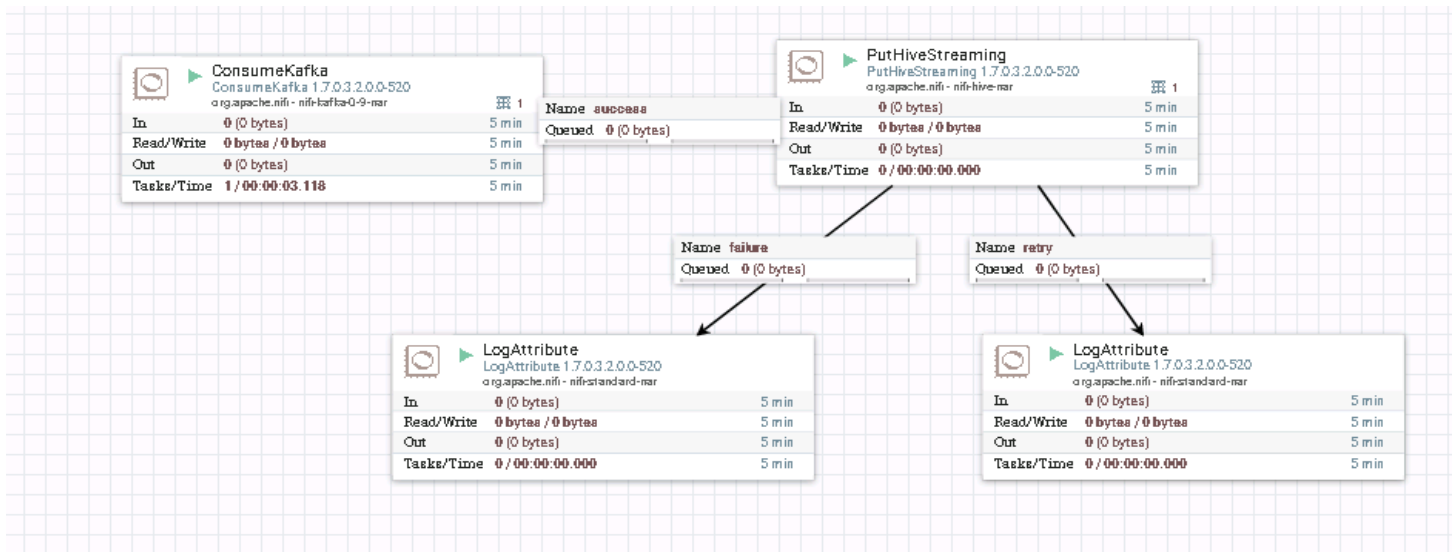The following, figure 20 is the NiFi development palatte on the Hadoop cluster.



**Figure 20.** The development template for the data flow taking data in via the Kafka broker and writing that data to the Hive table

## APPENDIX C: HIVE TABLE CREATION SCRIPT

```
create table solar_new (opcode   string,
                        solarfarm      string,
                        season         string,
                        timedate       string,
                        avgIntkWh      double,
                        avgIntkW       double)
                        clustered by (timedate) into 10 buckets
                        stored as orc
                        tblproperties('transactional'='true');
```