

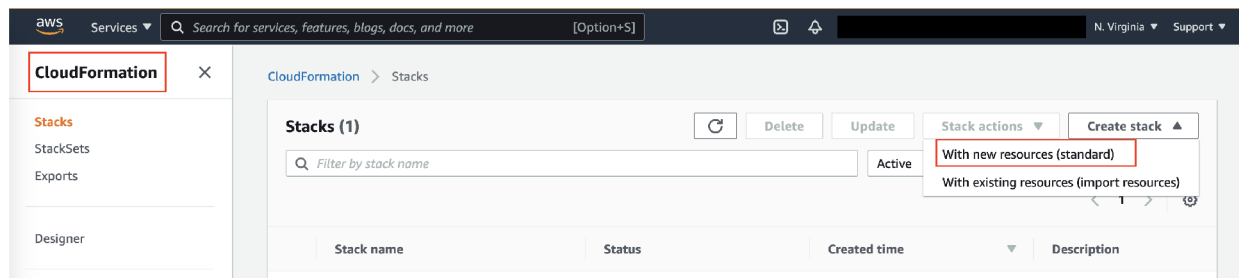
FeatureStore CompactUtil Deployment Guide

Deployment steps:

The CloudFormation (CFN) template creates 2 Lambda functions, a Step function, an Event bridge rule and IAM roles necessary for each service to invoke the next service.

Follow these steps to deploy the stack:

- There are 4 artifacts provided
 - sm-featurestore-offline-compact.zip - Lambda deployment package for compacting function.
 - sm-featurestore-offline-job-monitor.zip - Lambda deployment package for processing job status check.
 - sm-featurestore_offline_compact_spark.py - PySpark script that is fed to sagemaker processing to run compaction.
 - sm-featurestore-compact-util-cf.yml - Cloudformation template that creates all necessary resources.
- Upload the two zip files sm-featurestore-offline-compact.zip and sm-featurestore-offline-job-monitor.zip and the pyspark script sm-featurestore_offline_compact_spark.py to an S3 bucket. The S3 bucket has to be created in the same AWS Region as where Cloudformation stack will be deployed.
- Upload the sm-featurestore-compact-util-cf.yml via AWS CloudFormation console to deploy the stack. Input the S3 bucket name and prefix path as parameters to Cloudformation.
 - Go to CloudFormation console and create stack with new resources



- Upload the yml file provided and click Next

Create stack

Step 1
Specify template

Step 2
Specify stack details

Step 3
Configure stack options

Step 4
Review

Prerequisite - Prepare template

Prepare template
Every stack is based on a template. A template is a JSON or YAML file that contains configuration information about the AWS resources you want to include in the stack.

☒ Template is ready ☐ Use a sample template ☐ Create template in Designer

Specify template
A template is a JSON or YAML file that describes your stack's resources and properties.

Template source
Selecting a template generates an Amazon S3 URL where it will be stored.

☐ Amazon S3 URL ☒ Upload a template file

Upload a template file

Choose file **sm-featurestore-compact-util-cf.yml**

JSON or YAML-formatted file

S3 URL: <https://s3-external-1.amazonaws.com/cf-templates-1ekcp88p86gol-us-east-1/2021314rTa-sm-featurestore-compact-util-cf.yml> [View in Designer](#)

Cancel **Next**

- Enter a name for the stack, the S3 bucket name where code is uploaded to, and the prefix location if any where code is located. Leave prefix empty if the files are directly uploaded to a bucket and not within folders. Note that prefix should not have leading or trailing slash. Click Next.

CloudFormation > Stacks > Create stack

Step 1
[Specify template](#)

Step 2
Specify stack details

Step 3
Configure stack options

Step 4
Review

Specify stack details

Stack name

Stack name

Stack name can include letters (A-Z and a-z), numbers (0-9), and dashes (-).

Parameters
Parameters are defined in your template and allow you to input custom values when you create or update a stack.

S3 Location to lambda code and pyspark script

Bucket name where code is located

Prefix path where code is located - (Folder path within the bucket without leading or trailing slash)

Cancel Previous **Next**

- Leave the default settings on the next page and Click Next.
- On the last page, acknowledge the access permissions and click Create Stack.

Stack creation options

Timeout
-

Termination protection
Disabled

Capabilities and transforms

Transforms might require access capabilities

A transform might add Identity and Access Management (IAM) resources that could provide entities access to make changes to your AWS account. If a transform adds IAM resources, you must acknowledge their capabilities to create or update them. Ensure that you want to create or update the IAM resources, and that they have the minimum required permissions. In addition, if they have custom names, check that the names are unique within your AWS account. [Learn more](#)

☒ I acknowledge that AWS CloudFormation might create IAM resources.

☐ I acknowledge that AWS CloudFormation might create IAM resources with custom names.

☒ I acknowledge that AWS CloudFormation might require the following capability:
CAPABILITY_AUTO_EXPAND

Cancel Previous Create change set **Create stack**

- SageMaker Processing job is created and run by sm-featurestore-offline-compact Lambda function. The lambda function is created with these environment configuration that can be modified anytime via Lambda console.
 - PYSPARK_SCRIPT_PATH - Pyspark script location in S3
 - SAGEMAKER_INSTANCE_COUNT - Instance count for the job, default is set to 1.
 - SAGEMAKER_INSTANCE_TYPE - Instance type, by default is configured with ml.m5.4xlarge
 - SAGEMAKER_INSTANCE_VOLUME_SIZE - Size of ML storage volume in GB, default is 30.
 - SAGEMAKER_ROLE - IAM role to execute the job, default is what gets created while deploying the stack.
 - SPARK_CONTAINER_IMAGE - Container image that the SageMaker Processing job should use. The location is configured in the Cloudformation template based on the region that the stack gets created in.

Using the util:

Refer to FeatureStore-CompactUtil-Solution.pdf for detailed architecture and input JSON structure.

The Cloudformation creates an event bridge rule disabled by default and with a default schedule to run every 24 hours. The target for the rule is the step function.

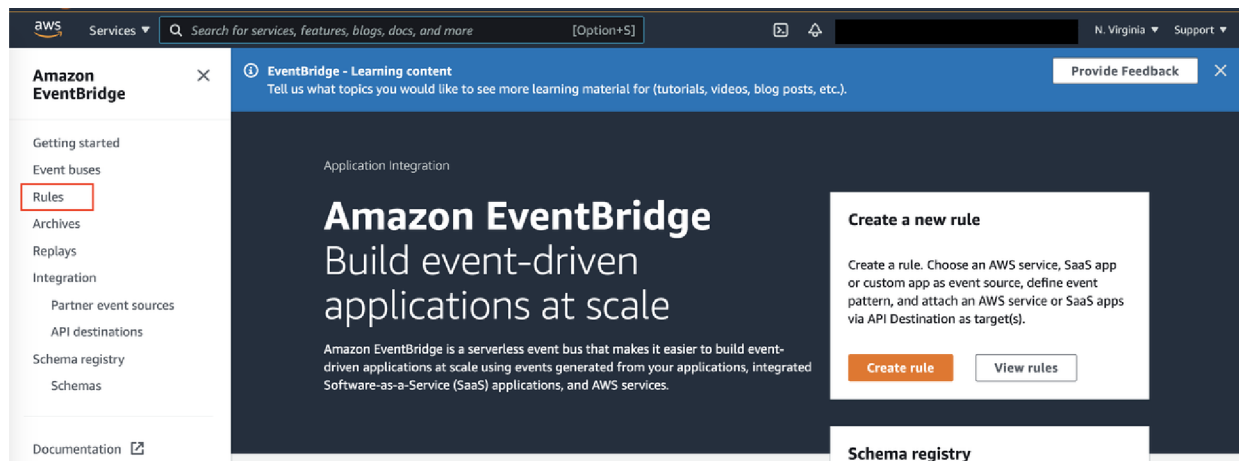
Go to EventBridge console and edit the event rule to

1. Run the util in full mode the first time
2. Run the util in incremental mode after full mode

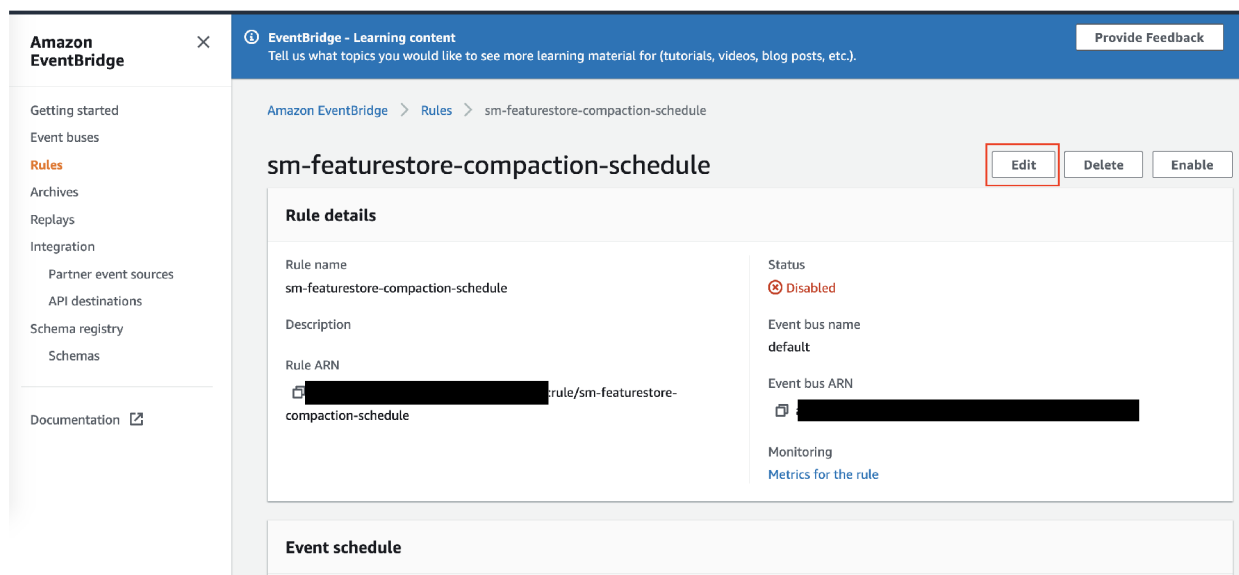
For #1 Run the util in full mode the first time

:

- **Navigate to EventBridge console and click Rules.**



- Click **sm-featurestore-compaction-schedule** rule and click **Edit**.



- Modify the input JSON with the feature group name, compact_mode as “full”, partiton_mode as “hour” or “day” based on needs, “compact_uri” if compated files need to go into user specified location and click Update.

Use the event JSON structure shown here. All keys are mandatory but optional values can be blank.

```
{
  "feature_group_name": "REPLACE_WITH_FEATURE_GROUP_NAME",
  "compact_mode": "full",
  "partition_mode": "hour",
  "compact_uri": "",
  "year": "",
}
```

```

"month": "",
"day": ""
}

```

Target
Select target(s) to invoke when an event matches your event pattern or when schedule is triggered (limit of 5 targets per rule).

Step Functions state machine

State machine
sm-featurestore-offline-compact

▼ Configure input

☐ Matched events [Info](#)

☐ Part of the matched event [Info](#)

☒ Constant (JSON text) [Info](#)

`{\"FeatureGroupName\":\"featuregroupname\", \"CompactMode\":\"full\", \"PartitionMode\":\"\", \"CompactURI\":\"\", \"Year\":\"2021\"}`

☐ Input transformer [Info](#)

☐ Create a new role for this specific resource

☒ Use existing role

Amazon_EventBridge_Invoke_Step_Functions_2093064009

[Learn more](#) about EventBridge identity-based policies.

► Retry policy and dead-letter queue

- **Enable the event**

Amazon EventBridge

EventBridge - Learning content
Tell us what topics you would like to see more learning material for (tutorials, videos, blog posts, etc.). [Provide Feedback](#)

Amazon EventBridge > Rules > sm-featurestore-compaction-schedule

sm-featurestore-compaction-schedule [Edit](#) [Delete](#) [Enable](#)

Rule details

Rule name sm-featurestore-compaction-schedule	Status Disabled
Description	Event bus name default
Rule ARN [redacted]/sm-featurestore-compaction-schedule	Event bus ARN [redacted]
	Monitoring Metrics for the rule

- **The event will fire in 30-60 seconds. Navigate to Step functions and check state machine sm-featurestore-offline-compact if execution has begun (Refer to monitoring section below).**

For #2 Run the util in incremental mode after full mode

- Once the util runs in “full” mode, it is time to setup the util to run in incremental manner. Navigate to EventBridge and edit the schedule to set schedule to run at midnight UTC everyday. The util will run every night at 12 AM UTC and compacts previous day’s files.

The screenshot shows the AWS EventBridge console interface. The 'Define pattern' section is active, with the 'Schedule' radio button selected. Under the 'Schedule' section, the 'Cron expression' radio button is selected. The text input field for the cron expression contains '0 0 * * ? *'. A red rectangular box highlights the 'Cron expression' section, including the text input field and the 'Next 10 trigger date(s)' dropdown.

Change the Schedule from Fixed rate every to Cron expression and set an expression to fire event at 12 AM UTC.

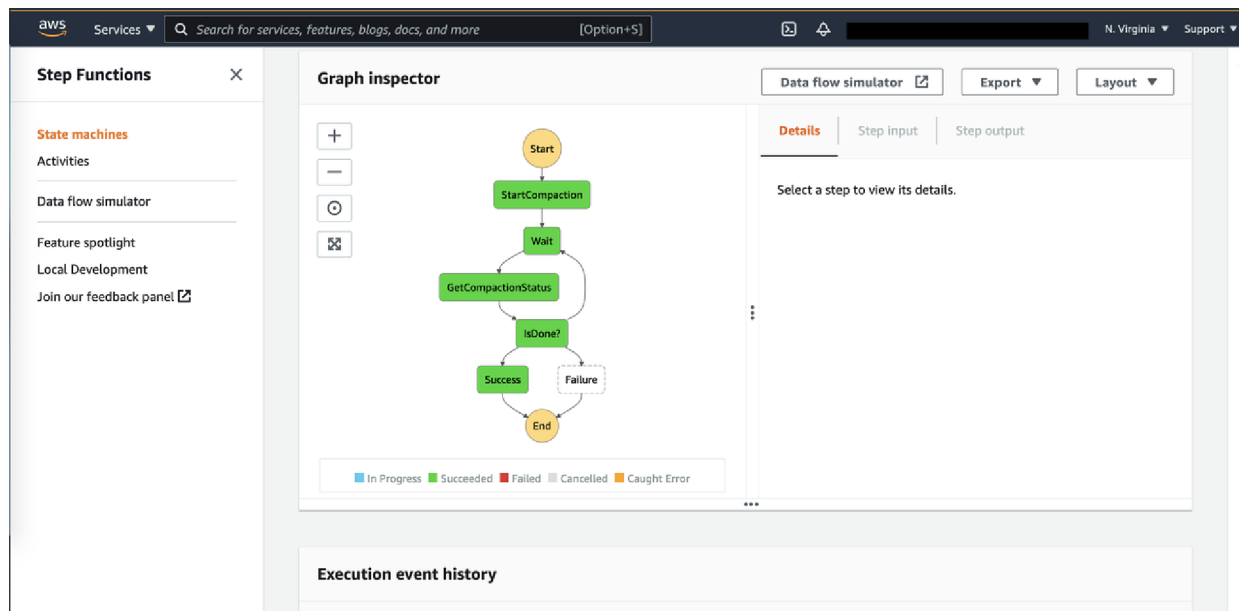
0 0 * * ? *

Set the input (Constant JSON text) as:

```
{
  "feature_group_name": "REPLACE_WITH_FEATURE_GROUP_NAME",
  "compact_mode": "incremental",
  "partition_mode": "",
  "compact_uri": "",
  "year": "",
  "month": "",
  "day": ""
}
```

Monitoring:

- Lambda and SageMaker processing jobs will create logs in Cloudwatch.
- Step functions can be monitored to check if all steps executed without errors. Navigate to Step Functions via AWS Console → StateMachines → Click on sm-featurestore-offline-compact → Click on the latest execution. A graph should be displayed to show the status of the execution. Each step in the graph shows details of input, output, execution success/failure and errors if any.



Validation:

- Check the S3 compact URI location for the compacted files. They will be partitioned by year, month day (and hour if partition_mode is hour). If compact_uri was not specified as an input, the files will be in the same bucket as the offline store. The S3 location will be in the below format
 - `s3://<bucket-name>/<customer-prefix>/<account-id>/sagemaker/<aws-region>/compact-offline-store/<feature-group-name>-<feature-group-creation-time>/data/year=<event-time-year>/ month=<event-time-month>/day=<event-time-day>/hour=<event-time-hour>/`
- Create glue crawler and tables for this location and run queries to validate the data.

Helpful resources:

EventBridge rules - <https://docs.aws.amazon.com/eventbridge/latest/userguide/eb-create-rule-schedule.html>

EventBridge schedule expressions

- <https://docs.aws.amazon.com/AmazonCloudWatch/latest/events/ScheduledEvents.html>

EventBridge targets - <https://docs.aws.amazon.com/eventbridge/latest/userguide/eb-targets.html>

SageMaker Processing job - <https://docs.aws.amazon.com/sagemaker/latest/dg/processing-job.html>

Step Functions invoking Lambda - <https://docs.aws.amazon.com/step-functions/latest/dg/connect-lambda.html>

Lambda - <https://docs.aws.amazon.com/lambda/latest/dg/welcome.html>