

# Integration of AWS Redshift with CADP

## AWS Redshift [Overview]

This document describes how to configure and integrate CipherTrust Manager with AWS Redshift.

Amazon Redshift is a data warehouse product which forms part of the larger cloud-computing platform Amazon Web Services. It is built on top of technology from the massive parallel processing (MPP) data warehouse company ParAccel to handle large scale data sets and database migrations. Redshift differs from Amazon's other hosted database offering, Amazon RDS, in its ability to handle analytic workloads on big data data sets stored by a column-oriented DBMS principle. Redshift allows up to 16 petabytes of data on a cluster compared to Amazon RDS Aurora's maximum size of 128 terabytes. Thales provides a couple of different methods to protect sensitive data in AWS Redshift.

### Bring Your Own Encryption (BYOE)

- **Data Ingest** – with Thales Batch Data Transformation (BDT)
- **Data Access** –user defined functions for column level encrypt and decryption using Thales CADP and tokenization using Thales CT-VL.

### Bring/Hold Your Own Key (BYOK) (HYOK)

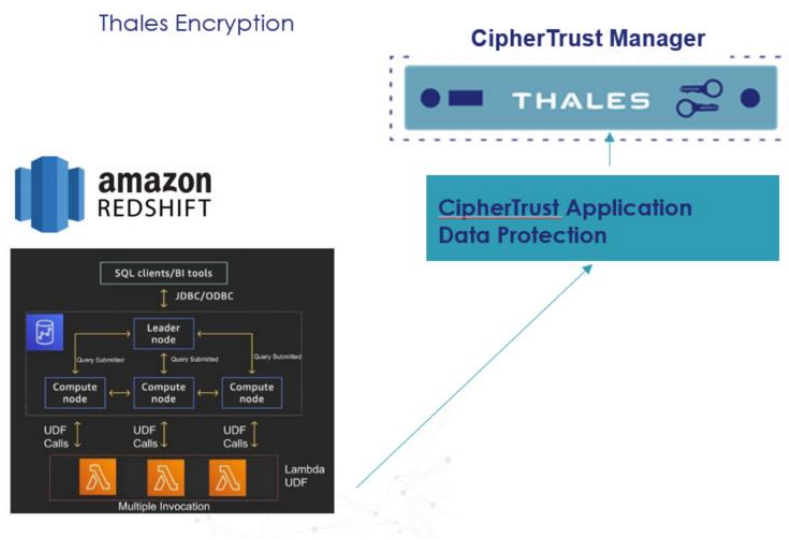
- **AWS Redshift Customer Managed Keys-** with Thales CM CCKM BYOK and HYOK.

**The above methods are NOT mutually exclusive. All methods can be used to build a strong defense in depth strategy to protect sensitive data in the cloud. The focus of this integration will be on Data Access protecting sensitive data in AWS Redshift columns by using CADP to create User Defined Functions (UDF) for encryption and decryption of sensitive data.**

### Architecture

The examples provided in this document use a capability AWS Redshift called “External Function”. A Redshift external function lets you incorporate functionality with software outside of Redshift by providing a direct integration with AWS Lambda Functions. A Redshift external function allows you to implement your function in other languages other than SQL. Listed

below is a diagram of how this integration works.



## Supported Product Versions

- **CipherTrust Manager** CipherTrust Manager 2.14 and higher.
- **CADP** CADP Java 8.13 and higher. Note: CADP 8.15.0.001 was tested. If a higher version of CADP is used you may encounter a class not found error with AWS Lambda Functions.
- **AWS Redshift**

This integration is validated using 2<sup>nd</sup> generation Google Cloud Functions and Java 11 along with CM 2.14 and CADP 8.15.0.001

## Prerequisites

Steps performed for this integration were provided by this AWS link:

<https://docs.aws.amazon.com/redshift/latest/dg/udf-creating-a-lambda-sql-udf.html>

[https://docs.aws.amazon.com/redshift/latest/dg/r\\_CREATE\\_EXTERNAL\\_FUNCTION.html](https://docs.aws.amazon.com/redshift/latest/dg/r_CREATE_EXTERNAL_FUNCTION.html)

- Ensure that CADP for Java is installed and configured. Refer to <https://thalesdocs.com/ctp/con/cadp/cadp-java/latest/admin/cadp-for-java-quick-start/cadp-for-java-installer/index.html>
- Ensure that the CipherTrust Manager is installed and configured. Refer to the [CipherTrust Manager documentation](#) for details.

- AWS Lambda functions communicates with the CipherTrust Manager using the Network Attached Encryption (NAE) Interface. Ensure that the NAE interface is configured. Refer to the [CipherTrust Manager documentation](#) for details.

## Steps for Integration

- [Installing and Configuring Thales CADP for Java]
- [Download code from Thales github and compile]
- [Publish jar/zip file to AWS Lambda Function (endpoint)]
- [Create AWS Redshift External Function]
- [Integration with Thales CipherTrust Manager]
- [Environment Variables]

## Installing and Configuring CADP for Java

To install and configure **CADP for Java**, refer to [Quick Start](#).

As noted in the link above Thales provides CADP in the Maven repository. This can be found at: [https://central.sonatype.com/artifact/io.github.thalescpl-io.cadp/CADP\\_for\\_JAVA](https://central.sonatype.com/artifact/io.github.thalescpl-io.cadp/CADP_for_JAVA)

Note: As indicated earlier please use the 8.15.0.001 version with Redshift. Eclipse development tool was used for these examples. Here is the version used for testing along with the Maven plugin for Eclipse.

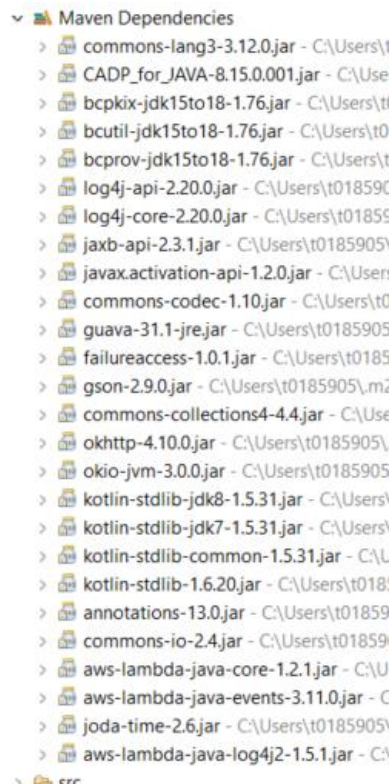
```
eclipse.buildId=4.15.0.I20200305-0155
```

```
m2e - Maven Integration for Eclipse
```

Example:

```
<dependency>
  <groupId>io.github.thalescpl-io.cadp</groupId>
  <artifactId>CADP_for_JAVA</artifactId>
  <version>8.15.0.001</version>
</dependency>
```

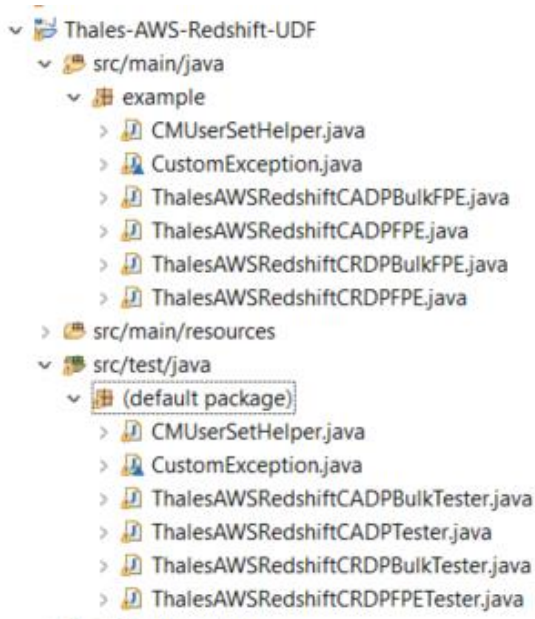
Here is a screenshot in eclipse of the jars from the pom file used for these examples:



## Download code from github and compile.

git clone [https://github.com/ThalesGroup/CipherTrust\\_Application\\_Protection.git](https://github.com/ThalesGroup/CipherTrust_Application_Protection.git)

The database directory has all the code for AWS Redshift. You should see the classes in your project.



CADP supports a bulk API which allows for CADP to batch requests before calling encrypt or decrypt. . A single class file that accepts environment variables will allow the UDF to be used by more than one function. There is a limit of 10,000 items to encrypt but it also depends on the size of the items as well so review the CADP documentation to make sure those thresholds are not exceeded. Typically, the point of diminishing returns is around 200.

Assuming you have your CM already configured the `ThalesAWSRedshiftCADPTester` can be used to test basic connection to CM to make sure your CM environment is configured correctly. You will need to modify the keyname to make sure it exists in CM and there are environment variables that must be set as well. Please see the appendix for all the environment variables and descriptions.

If you plan on using UserSets the user must also be in the Application Data Protection Client Group. See section on Environment Variables for more details. If you have already installed CM then you will need to update the `CADP_for_JAVA.properties` file with all the necessary settings such as IP/NAE Port, etc. [The file is located under the resource's directory in the java project for eclipse.](#)

### **Generate the jar file to upload to the CSP.**

To compile and generate the target jar file to be uploaded to AWS Lambda select the project and choose "Run As" "maven install" to generate the target.

```
[INFO] Replacing C:\Users\t0185905\workspace\CADP-AWS-Redshift-UDF\target\CADP-AWS-Redshift-UDF-0.0.2-SNAPSHOT.jar with C:\Users\t0185905\workspace\CADP-AWS-Redshift-UDF\target\CADP-AWS-Redshift-UDF-0.0.2-SNAPSHOT-shaded.jar
[INFO]
[INFO] --- maven-install-plugin:3.0.1:install (default-install) @ CADP-AWS-Redshift-UDF ---
[INFO] Installing C:\Users\t0185905\workspace\CADP-AWS-Redshift-UDF\pom.xml to C:\Users\t0185905\.m2\repository\Thales\CADP-AWS-Redshift-UDF\0.0.2-SNAPSHOT\CADP-AWS-Redshift-UDF-0.0.2-SNAPSHOT.pom
[INFO] Installing C:\Users\t0185905\workspace\CADP-AWS-Redshift-UDF\target\CADP-AWS-Redshift-UDF-0.0.2-SNAPSHOT.jar to C:\Users\t0185905\.m2\repository\Thales\CADP-AWS-Redshift-UDF\0.0.2-SNAPSHOT\CADP-AWS-Redshift-UDF-0.0.2-SNAPSHOT.jar
[INFO] -----
[INFO] BUILD SUCCESS
[INFO] -----
[INFO] Total time: 4.191 s
[INFO] Finished at: 2024-05-15T15:59:06-04:00
[INFO] -----The
```

code provided has uses Googles userDefinedContext capability. This code accepts a mode and a datatype as a keyvalue pair.

## **Publish jar/zip file to AWS Lambda Function. (endpoint)**

Once you have generated the jar file to upload you can then create the CSP function.

## AWS Lambda Function

The screenshot shows the 'Create function' page in the AWS Lambda console. At the top, there are three radio buttons for creating a function: 'Author from scratch' (selected), 'Use a blueprint', and 'Container image'. Below this is the 'Basic information' section, which includes a 'Function name' field with the placeholder 'YourFunctionName', a 'Runtime' dropdown menu set to 'Java 11', and an 'Architecture' section with 'x86\_64' selected. There are also links for 'Permissions' and 'Advanced settings'. At the bottom right, there are 'Cancel' and 'Create function' buttons.

Click Create function.

Set environment variables appropriate values (See section on Environment Variables for details) and then select Next. Upload the jar or zip file on the next screen. Be sure to change the entry point to reflect your class name. ***The example below is ThalesAWSRedshiftCADPBulkFPE This should match your code path.*** Be sure to have a role that allows for AWSLambdaBasicExecutionRole. 512MB of RAM is plenty. You can monitor your tests and possibly lower it.

Once you have created the functions above and if you have already configured and setup CM with the key and all the environment variables you can test the function with the test tab. You will need to provide the appropriate json to test. Here is an example:

```
{
  "request_id": "23FF1F97-F28A-44AA-AB67-266ED976BF40",
  "cluster": "arn:aws:redshift:xxxx",
  "user": "adminuser",
  "database": "db1",
  "external_function": "public.foo",
  "query_id": 5678234,
  "num_records": 3,
  "arguments": [
    "The. king. and I me"
```

```
    ],  
    [  
        "The mroe data "  
    ],  
    [  
        "test dada"  
    ]  
]  
}
```

## Create AWS Redshift External Function.

Here are some links that provide details.

<https://docs.aws.amazon.com/redshift/latest/dg/udf-creating-a-lambda-sql-udf.html>

[https://docs.aws.amazon.com/redshift/latest/dg/r\\_CREATE\\_EXTERNAL\\_FUNCTION.html](https://docs.aws.amazon.com/redshift/latest/dg/r_CREATE_EXTERNAL_FUNCTION.html)

<https://aws-dojo.com/excercises/excercise31/>

As noted above the steps are:

1. Create the AWS Lambda Function. (should already be done from above)
2. Create external function object in AWS Redshift

### Examples:

Redshift Function Definitions.

The functions are created using the following format:

```
CREATE OR REPLACE EXTERNAL FUNCTION thales_token_cts_char(sensitivedata  
varchar)  
RETURNS varchar  
STABLE  
LAMBDA 'thales_aws_lambda_redshift_cts_token_char'  
IAM_ROLE 'arn:aws:youriam/hiftpoweruseraccess';
```

## Integration with CipherTrust Manager.

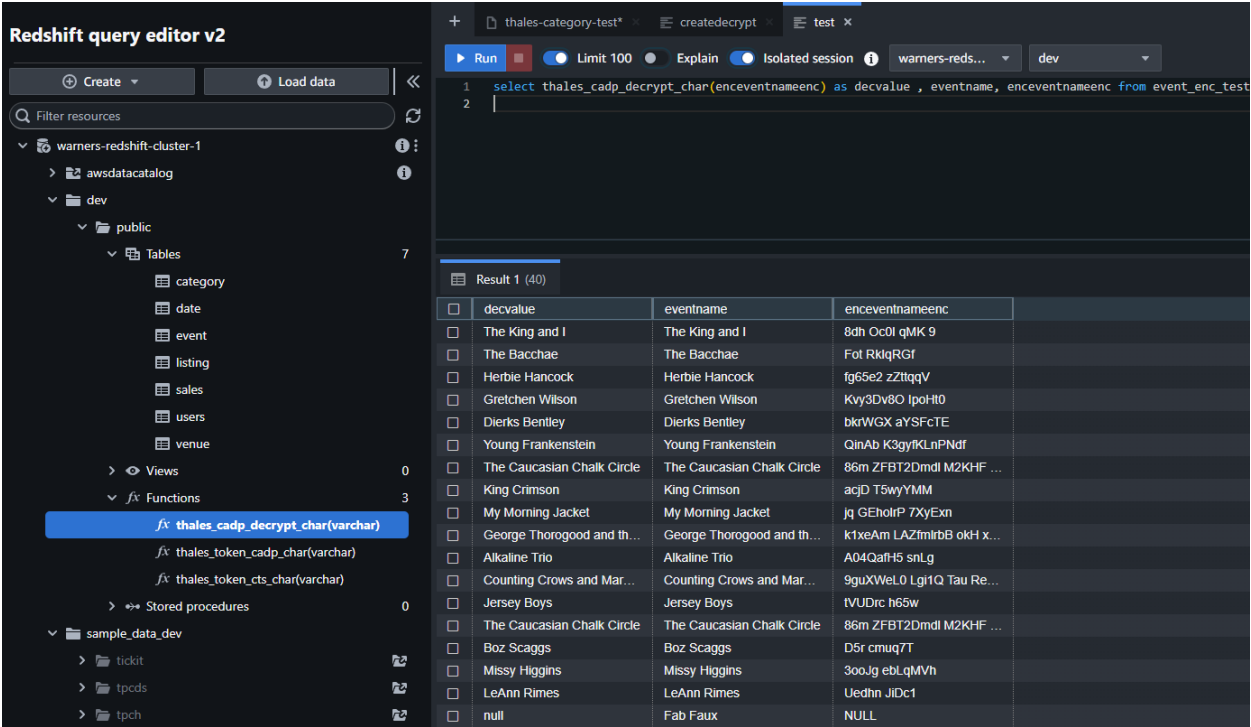
Assuming the NAE interface has already been setup based on the prerequisites the only other setup required is to have a user and key created in CM. This user should have the ability to export the key. This user will also be needed to be provided as an environment variable for the function. The examples provided have the key as a hardcoded value, but this can be easily altered to be provided as an environment variable, obtained from a secrets manager.

As noted above there is a test class that can be used to test connectivity with CM without having to publish the Function. If you have not updated the CADP\_for\_JAVA.properties file with the CM settings such as IP/NAE Port, etc. then do so now. The file is located under the resource’s directory in the eclipse project. These properties can also be overwritten with CADP code as well if your desire is to pass them in as environment variables or headers of the JSON request. Here is an example:

```
System.setProperty("com.ingrian.security.nae.NAE_IP.1", "10.20.1.9");
```

When all the above steps are performed you should see your UDF’s in AWS Redshift under Routines in the UI. Here is a sample query using one of the UDF’s.

*Sample Results:*



## Environment Variables

Listed below are the environment variables currently required for the Cloud Function.

Key	Value	Desc
BATCHSIZE	200	Nbr of rows to chunk when using batch mode
CMPWD	Yourpwd!	CM PWD if using CADP
CMUSER	apiuser	CM USERID if using CADP
CRDPID	20.221.216.666	CRDP Container IP if using CRDP



datatype	charint	datatype of column (char or charint)
keymetadata	1001000	policy and key version if using CRDP
keymetadatalocation	external	location of metadata if using CRDP (internal,external)
mode	revealbulk	mode of operation(protect,reveal,protectbulk,revealbulk) CRDP
protection_profile	plain-nbr-ext	protection profile in CM for CRDP
returnciphertextforuserwithnokeyaccess	yes	if user in CM not exist should UDF error out or return ciphertext
useridincm	716f01a6-5cab-4799-925a-6dc2d8712fc1	userid in cm if user lookup is done
useridlookup	no	should useridlookup be done (yes,no)
useridlookupip	20.241.70.666	userid lookup
showrevealinternalkey	yes	show keymetadata when issuing a protect call (CRDP)

## Advanced Topics.

### Application Data Protection UserSets

Application Data Protection UserSets are currently used for DPG to control how the data will be revealed to users. These UserSets can also be independent of any Access Policy. Most cloud databases have some way to capture who is running the query and this information can be passed to CM to be verified in a UserSet to ensure the person running the query has been granted proper access. In github there is a sample class file called CMUserSetHelper that can be used to load a userset with values from an external identity provider such as LDAP. The name of this method is `addAUserToUserSetFromFile`. Once users have been loaded into this userset the userid must be captured and used as an environment variable to the Function. The function has a number of environment variables that must be provided for the function to work. Please review the section on Environment Variables for more details.

### Options for keyname

Currently the sample code uses a hard coded key name. Other options to be considered are:

- 1.)keyname as an environment variable.
- 2.)keyname as a AWS secret.
- 3.)keyname passed in an attribute to the function. Example with a database view.  
create view employee as select first\_name, last\_name,  
thales\_cadp\_aws\_decrypt\_char('testfaas', email) as email from emp\_basic

### Options for handling null values.

Since it is not possible to encrypt a column that contains null values or any column that has 1 byte it is necessary to skip those to avoid getting an error when running the query. There are a couple of ways to handle this use case.

## Option 1. Modify the queries.

Many times, simply adding a where clause to exclude values that have nulls or less than 2 bytes can avoid query errors. For example: select \* from FROM

mw\_demo\_dataset\_US.plaintext50cadpemailprotected\_nulltest where email is not null and length(email) > 1; For those scenarios where that is not suffice some other examples are listed below. Here is an example of a select statement that can be modified to handle null values: Here is an example of a select statement that can be modified to handle null values:

```
SELECT
  name,
CASE
  WHEN email IS NULL THEN 'null'
  WHEN length(email) < 2 then email
  WHEN email = 'null' then email
  ELSE `cpl-sales-l-app-us-01.mw_demo_dataset_US.thales_cadp_encrypt_char`(email)
END AS email
FROM
  mw_demo_dataset_US.plaintext50cadpemailprotected_nulltest;
```

The above use case was for situations where the column contained both null and the word 'null'.

Here is an example of a select statement that can be modified to handle null values in the where clause.

```
SELECT name, email, email_enc
FROM
  mw_demo_dataset_US.plaintext50cadpemailprotected_nulltest
where CASE
  WHEN email_enc IS NULL THEN 'null'
  WHEN length(email_enc) < 2 then email_enc
  WHEN email_enc = 'null' then email_enc
  ELSE `cpl-sales-l-app-us-01.mw_demo_dataset_US.thales_cadp_decrypt_char`(email_enc)
END like "%gmail%"
```

Row	name	email	email_enc
1	Dr. Lemmie Zboncak	ikris@gmail.com	1IOEk@RPDqC.GHd
2	Zillah Leuschke	scronin@gmail.com	53mWY7Y@4bzb6.2D4
3	Troy Gaylord	devon49@gmail.com	EsTMziF@ISZg2.yN6
4	Dr. Spurgeon Wintheiser	hilah17@gmail.com	PIM4vAd@qAIV9.oJC
5	Tatyana Bernhard	denny56@gmail.com	yVQ7ITA@3idYW.GqV
6	Renata Hilpert	tdickens@gmail.com	plzg3zLb@oetcj.Opi
7	Deliah Douglas	sconnelly@gmail.com	jpXAUZWpX@koaH2.yS8
8	Amare Feeney	batz.geary@gmail.com	G04W.CRNlb@riDnQ.DHT
9	Dr. Cristy Schinner	schulist.garfield@gmail.com	Hll90wCU.LPtU3p6o@F9Nb7.G...
10	Vena Douglas	huston.christiansen@gmail.com	7OHpfz.PXI0PoJJJaWlQ@AvTP...

## Option 2. Modify the Cloud Function to skip encrypting these values.

Please review the method checkvalid for the logic of handling these use cases.

*Note: When using this logic it is important to know that for any column that is null will return 'null'. This should be fine for use cases where the query is not updating the column. For use cases where the source system is expecting null vs the word 'null' additional testing should be conducted on the systems that rely on this data type as being null vs the word 'null'.*