

Secure Data Commons

Data Analyst User Guide

www.its.dot.gov/index.htm
Draft Report — October 29, 2019
FHWA-JPO-18-xxx

Notice

This document is disseminated under the sponsorship of the Department of Transportation in the interest of information exchange. The United States Government assumes no liability for its contents or use thereof. The U.S. Government is not endorsing any manufacturers, products, or services cited herein and any trade name that may appear in the work has been included only because it is essential to the contents of the work.

Revision History

| # | Name | Version | Revision Date | Revision Description |
|---|------------|---------|---------------|--|
| 1 | REAN Cloud | 1.0 | 08/02/2018 | Initial Draft |
| 2 | REAN Cloud | 2.0 | 09/13/2018 | Add Export Functionality |
| 3 | REAN Cloud | 3.0 | 10/29/2019 | Add Manage Workstations section and update per feedback comments |

Table of Contents

| | |
|--|----|
| Table of Contents..... | iv |
| Chapter 1. Introduction and Document Overview..... | 1 |
| Prerequisites..... | 2 |
| Chapter 2. Initial Setup and Validation..... | 3 |
| Accessing Secure Data Commons Portal..... | 3 |
| Request Access to Datasets..... | 5 |
| Upload User Data to S3 Bucket Through Portal..... | 8 |
| Download User Data from S3 Bucket Through Portal..... | 9 |
| Chapter 3. Accessing and Launching Workstations..... | 10 |
| Launch Workstations..... | 10 |
| Software Validation..... | 12 |
| Connecting to the Data Warehouse Using SQL Workbench..... | 13 |
| Connecting to Waze Data in Redshift..... | 13 |
| Connecting to CVP Data in Hive Metastore..... | 14 |
| Update Data Formatting Settings in SQL Workbench..... | 14 |
| Connecting to Redshift from Linux Environments..... | 15 |
| Accessing Jupyter Notebook and RStudio Server..... | 16 |
| Manage Workstations..... | 16 |
| Resize Workstation..... | 17 |
| Schedule/Extend Uptime..... | 19 |
| Stop Workstations..... | 21 |
| Chapter 4. Sample Queries for SDC Datasets..... | 22 |
| Connected Vehicle Data..... | 22 |
| Overview..... | 22 |
| WYDOT..... | 22 |
| WYDOT BSM and TIM Metadata RECORDGENERATEDAT vs ODERECEIVEDAT...23 | |
| THEA..... | 24 |

| | |
|---|----|
| Geospatial Queries | 25 |
| Waze Data | 25 |
| Overview | 25 |
| Chapter 6. Exporting Datasets from the SDC | 27 |
| Chapter 7. Technical Documentation and Contact Information | 31 |
| Architecture Diagram..... | 31 |
| Workstation Details | 31 |
| Tools and Versions..... | 32 |
| Contact Information | 32 |
| Useful Links..... | 32 |
| AWS S3 CLI Commands..... | 32 |
| Chapter 8. Frequently Asked Questions..... | 34 |
| How can I get access to the SDC Datasets? | 34 |
| How will I understand what a particular dataset consists of?..... | 35 |
| How can I launch a workstation? | 35 |
| Where do I store my data? | 36 |
| How can I bring my own datasets/algorithm to my workstation?..... | 36 |
| How can I publish my dataset/algorithm? | 36 |

List of Tables

| | |
|--|----|
| Table 1: Default Workstation Details | 31 |
| Table 2: List of Tools Used and Their Versions | 32 |

List of Figures

| | |
|--|----|
| Figure 1: SDC Starting Homepage | 3 |
| Figure 2: Register/Login Page | 4 |
| Figure 3: Login Page | 4 |
| Figure 4: Landing Page After Login..... | 5 |
| Figure 5: Datasets Option..... | 6 |
| Figure 6: Request Dataset Access..... | 6 |
| Figure 7: Send Data Access Request..... | 7 |
| Figure 8: Uploading Files..... | 8 |
| Figure 9: Upload Success..... | 8 |
| Figure 10: Selecting Files for Download | 9 |
| Figure 11: Starting Workstations | 10 |
| Figure 12: Run Instance Success | 10 |
| Figure 13: Launch Workstations..... | 11 |
| Figure 14: Initialization and Login Screens | 11 |
| Figure 15: SQL Workbench Icon | 13 |
| Figure 16: Create Redshift Connection Profile..... | 13 |
| Figure 17: Hive Connection Settings | 14 |
| Figure 18: Tools → Options → Data formatting..... | 15 |
| Figure 19: Manage Workstation | 16 |
| Figure 20: Manage Workstation Options | 17 |
| Figure 21: Resize Workstation Option | 17 |
| Figure 22: Workstation Stopped for Resize Changes..... | 17 |
| Figure 23: Resizing Options | 18 |
| Figure 24: Schedule Resize | 19 |
| Figure 25: Schedule Workstation Uptime Option | 19 |
| Figure 26: Schedule Uptime | 20 |
| Figure 27: Tooltip with Existing Scheduled Uptime..... | 20 |
| Figure 28: New Tooltip with Extended Uptime Schedule | 20 |
| Figure 29: Stop Workstation | 21 |
| Figure 30: WYDOT BSM Tables | 22 |
| Figure 31: Request Export | 27 |
| Figure 32: Request Export Form | 28 |
| Figure 33: Approval Form Fields | 29 |
| Figure 34: Acceptable Use Policy | 30 |
| Figure 35: SDC Architecture Overview..... | 31 |
| Figure 36: Datasets Tab..... | 34 |
| Figure 37: Dataset Access Request | 35 |
| Figure 38: Workstations Tab..... | 35 |
| Figure 39: Login for Programming Environment..... | 36 |
| Figure 40: Upload Button..... | 37 |

| | |
|---|----|
| Figure 41: Publish Dataset Form | 37 |
| Figure 42: Algorithm Request Form | 38 |

Chapter 1. Introduction and Document Overview

The Secure Data Commons (SDC) is a United States Department of Transportation (U.S DOT) sponsored cloud-based analytical sandbox designed to create wider access to sensitive transportation datasets, with the goal of advancing the state of the art of transportation research and state/local traffic management.

The SDC stores sensitive transportation data made available by participating Data Providers, and grants access to approved researchers to these datasets. The SDC also provides access to open source tools and allows researchers to collaborate and share code with other system users.

The SDC platform is a research environment that allows users to conduct analyses and do development and testing of new tools and software products. It is not intended to be an alternative to any local jurisdiction's traffic management center or local data repository. The current SDC platform provides users with the following data, tools, and features:

- **Data:** The SDC is ingesting several datasets currently. Additional datasets will be added to the environment over time. Users can bring their own data into the environment to use along with the Waze data.
- **Tools:** The environment provides access to open source tools including Python, RStudio, Microsoft R, SQL Workbench, Power BI, and Jupyter Notebook. These tools are available on a virtual machine in the system enabling data analytics in the cloud.
- **Functionality:** Users can access and analyze data within the environment, save their work to a virtual machine, and publish processes and results to share with other SDC users.

The SDC platform supports two major roles:

- **Data Providers** – These are entities that provide data hosted on the SDC platform. The Data Provider establishes the data protection needs and acceptable use terms for the data analysts.
- **Data Analysts** – These are entities that conduct analysis using the datasets hosted within the SDC system. Note that analysts can bring their own data and tools into the SDC system.

This document provides guidance for the **Data Analyst** role. A similar guide will be prepared for the Data Providers. The document is organized as follows:

- Initial Setup and Validation
- Workstation Access
- Sample Queries
- Exporting Data
- Importing and Exporting Code
- Accessing External Data Sources within SDC

- Technical Support and Contact Information
- Frequently Asked Questions

Prerequisites

Workstation access will not be granted for a Data Analyst user until the user has:

1. submitted a completed Access Request Form,
2. received approval for the request,
3. received an email message with onboarding instructions from the support team, and
4. received a walkthrough of the system from the support team.

Refer to the [Useful Links](#) section later in the document for further information on technologies relevant to SDC.

Chapter 2. Initial Setup and Validation

This chapter provides guidance on the initial setup and validation of the user into the SDC system.

Accessing Secure Data Commons Portal

Users can access the SDC web portal by navigating to <https://portal.securedatacommons.com>.

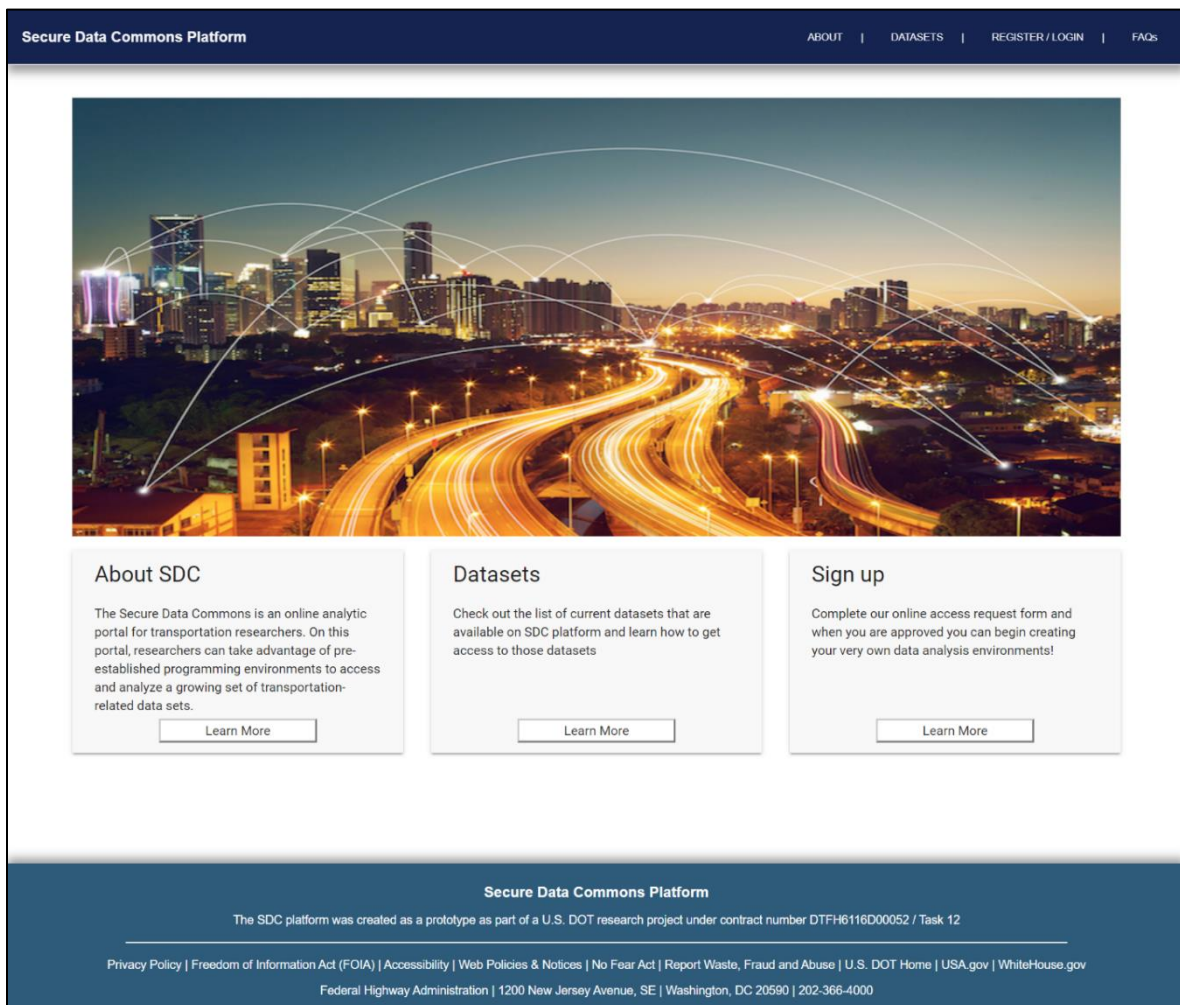


Figure 1: SDC Starting Homepage

Chapter 2. Initial Setup and Validation

Users can click on Register/Login from the top menu to access the Access Request Form and Privacy Policy links, as well as the Login button.

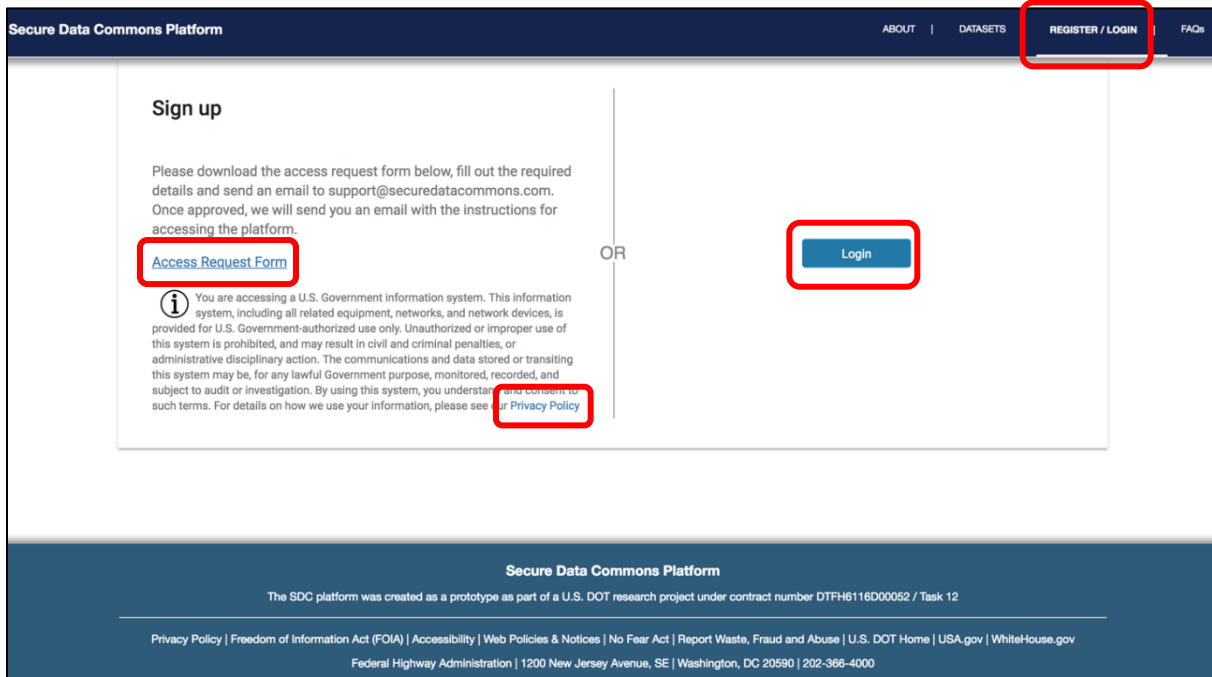


Figure 2: Register/Login Page

Clicking on Login redirects the user to the Secure Data Commons platform login page:

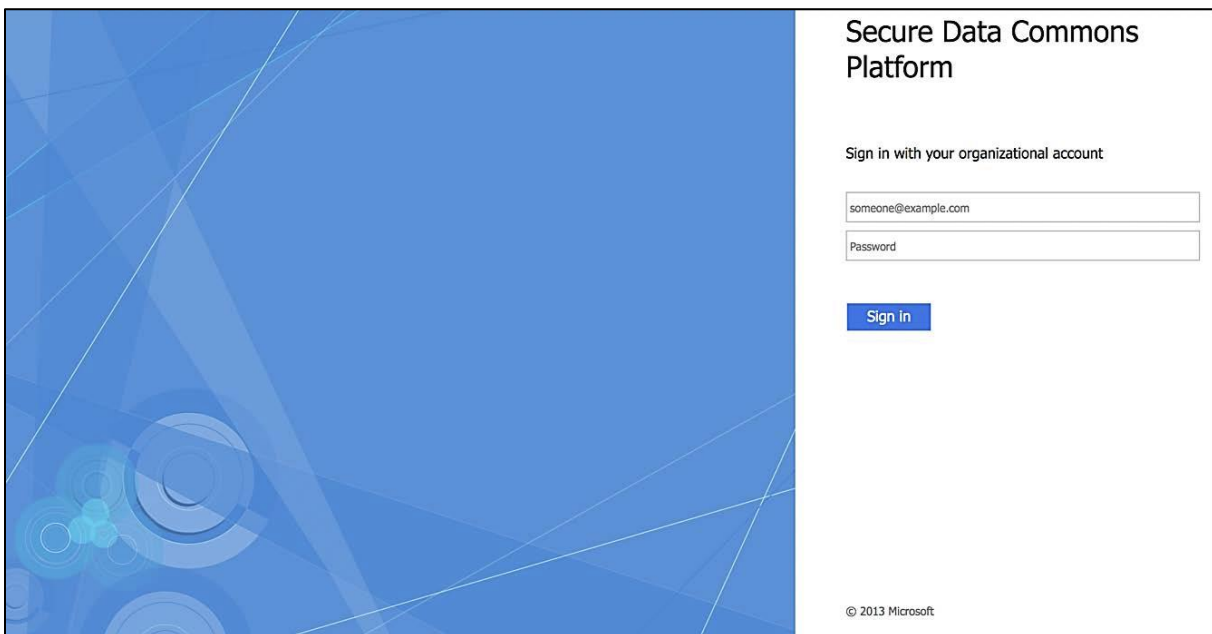


Figure 3: Login Page

If the user is accessing the portal for the first time, they will be prompted to change their password after entering the credentials provided in the welcome email.

U.S. Department of Transportation
Intelligent Transportation Systems Joint Program Office

Upon successfully logging in, the user is redirected to the landing page, which provides an overview of Secure Data Commons and the different actions the user can perform from the web portal:

1. Curated and published datasets
2. Access to workstations with programming tools
3. Bring your own datasets / algorithms
4. Publish your datasets / algorithms

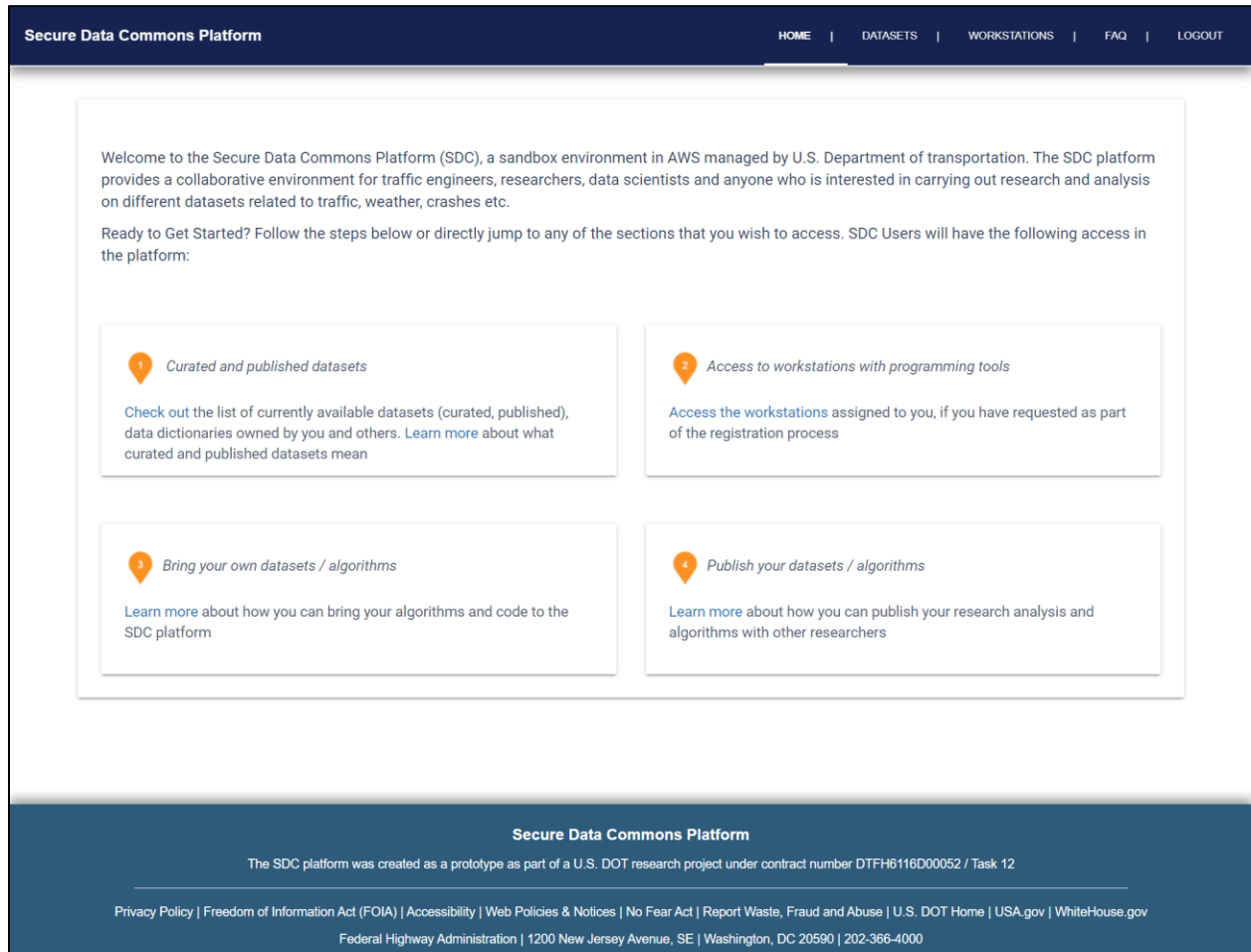


Figure 4: Landing Page After Login

Request Access to Datasets

Users can request access to the datasets that are available within the SDC platform as published / enabled by the SDC team or published by other users.

Once you are logged in, go to Datasets in the top menu.

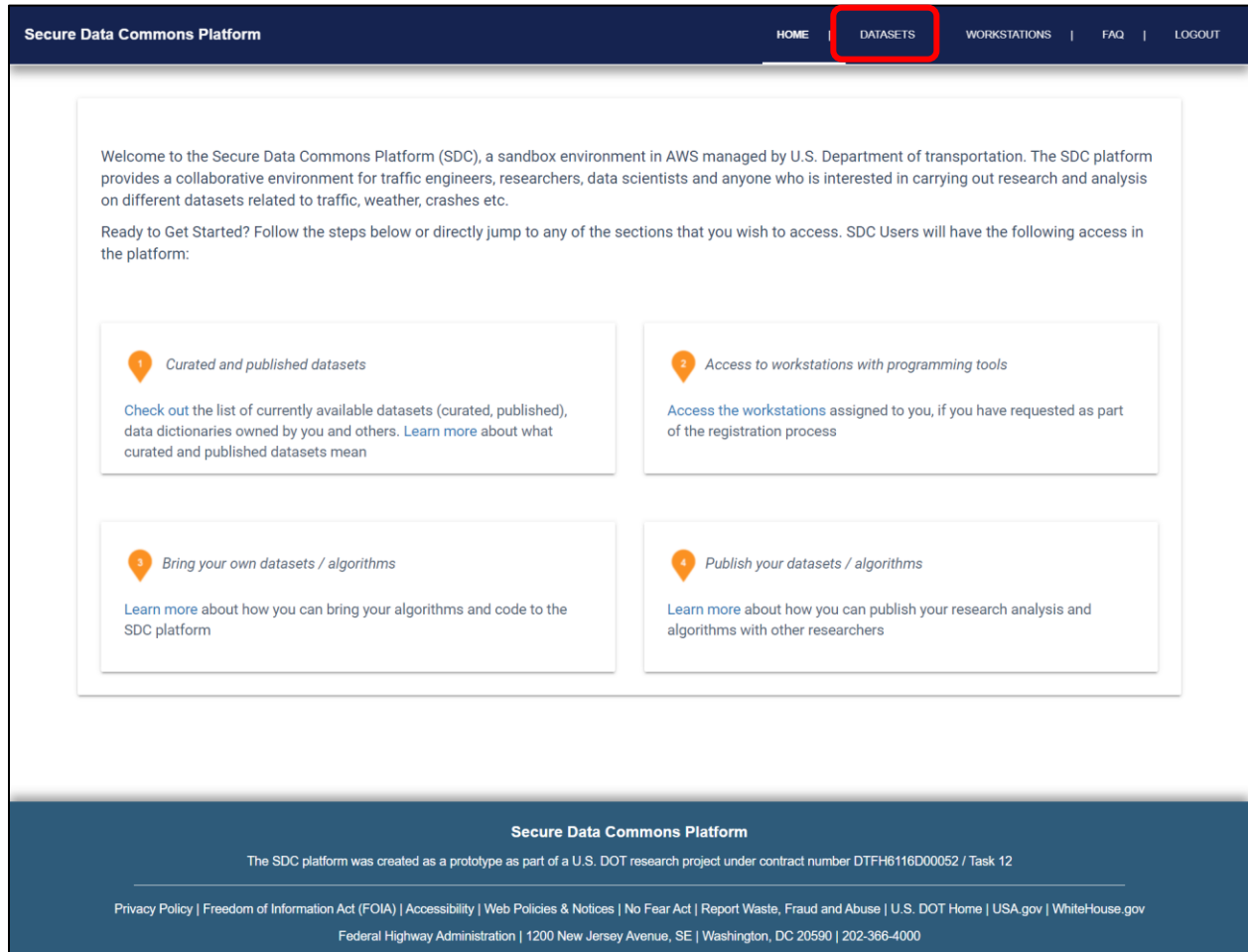


Figure 5: Datasets Option

Expand the SDC Datasets. You will be able to see all available datasets in the SDC platform. To access a dataset, click on Request.

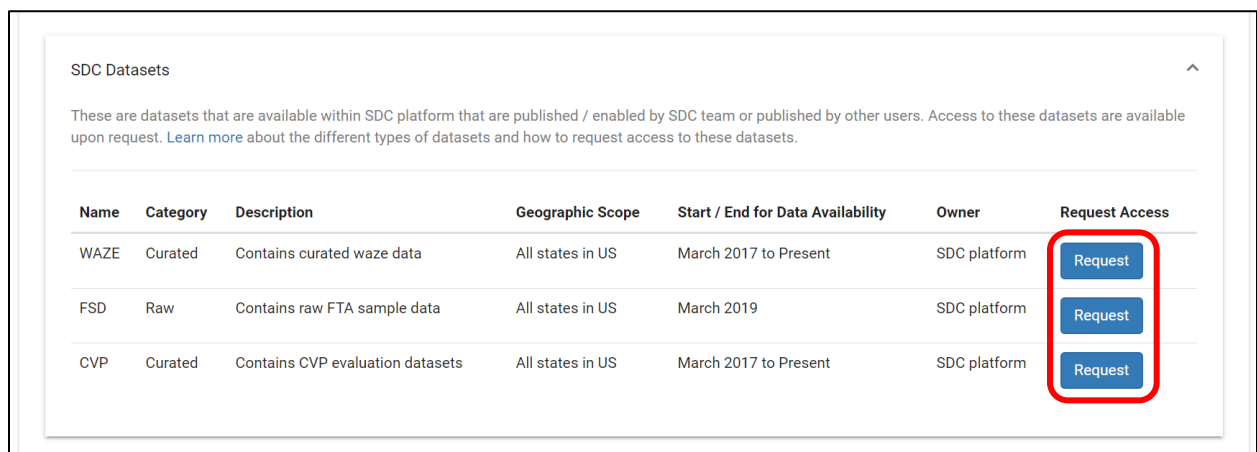


Figure 6: Request Dataset Access

Complete the SDC Data Access Request form that appears. Once completed, click on Send Request.

The screenshot shows the 'Secure Data Commons Platform' interface. A modal window titled 'SDC Data Access Request Form' is open. The form contains the following sections:

- Basis for Access**: A question 'Are you a USDOT Employee or researcher under contract to USDOT?' with radio button options for 'Yes' and 'No'.
- Geographic Extent of Access**: A text input field labeled 'Enter the state(s) for which you would like to request access. *' with a note 'use commas when requesting multiple states (e.g. KN, VI)'.
- Agreement**: A checkbox followed by the text: 'By clicking accept, the user agrees to the conditions set in the data provider agreement. If you are a USDOT employee or engaged in USDOT-sponsored research activity, you are agreeing to the overall data agreements that USDOT has with the data providers.'
- Buttons**: At the bottom of the form are two buttons: 'SEND REQUEST' (highlighted with a red box) and 'CANCEL'.

The background shows a sidebar with 'SDC Datasets' and 'SDC Algorithms' sections, and a main content area with a table of datasets.

Figure 7: Send Data Access Request

The request will be sent to the support team and access to the requested dataset will be given upon validation and approval of the information in the form.

Upload User Data to S3 Bucket Through Portal

Users can upload their own data to their assigned team/individual buckets through the portal.

1. Click on Datasets from the home page.
2. Click on Upload Files under “My Datasets / Algorithm.”
3. A pop-up window appears prompting you to choose one or more files for upload to the assigned bucket. (The assigned bucket name will be displayed on the upload pop-up window.)

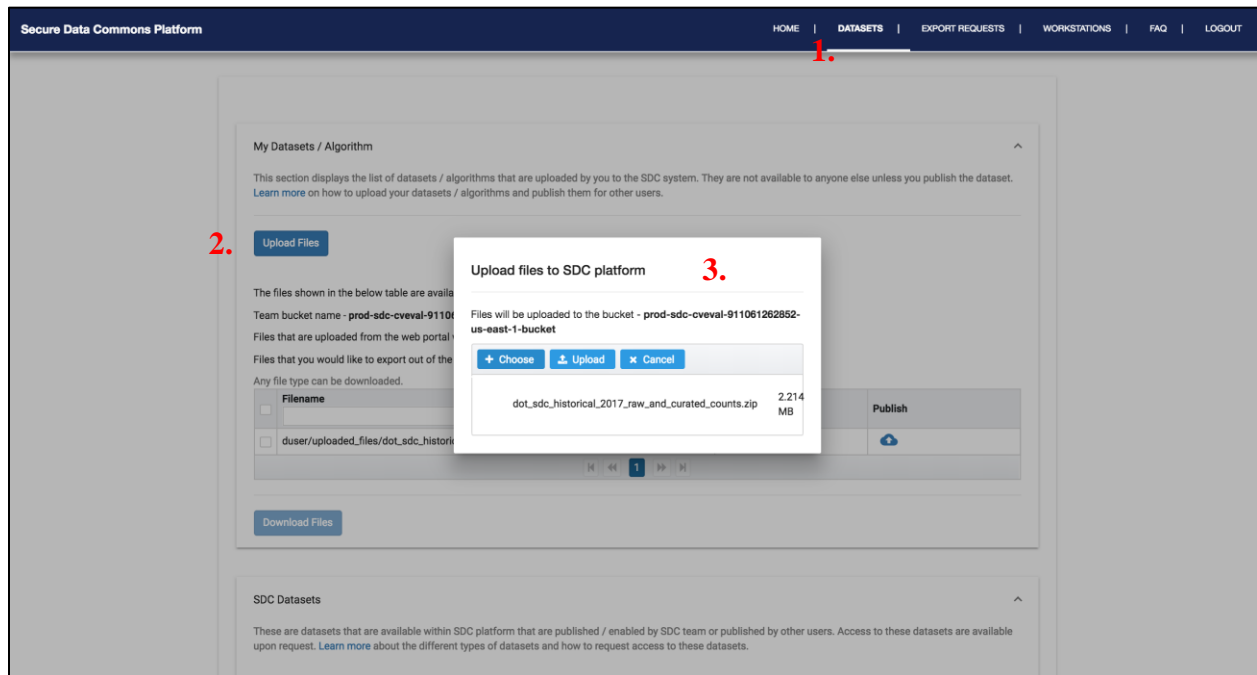


Figure 8: Uploading Files

4. A success message will be displayed upon a successful upload.

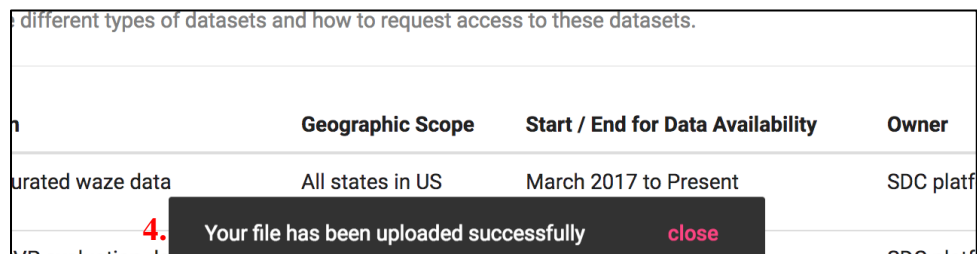


Figure 9: Upload Success

5. Files that are uploaded from the web portal will be saved in the folder – **username/uploaded_files**
6. Users would be able to access only the files that are under **uploaded_files**, **export_requests** folder.

Download User Data from S3 Bucket Through Portal

Users can download their data from their assigned team/individual buckets through the portal.

1. Click on Datasets from the home page.
2. All the available files under **username/uploaded_files** in the assigned bucket will be displayed along with the assigned bucket name under My Datasets / Algorithm.
3. Select the files that you want to download and then click on Download Files.
4. Users should go through the export request workflow to download files that are uploaded under **export_requests** folder. Export requests workflow can be found by clicking [here](#).

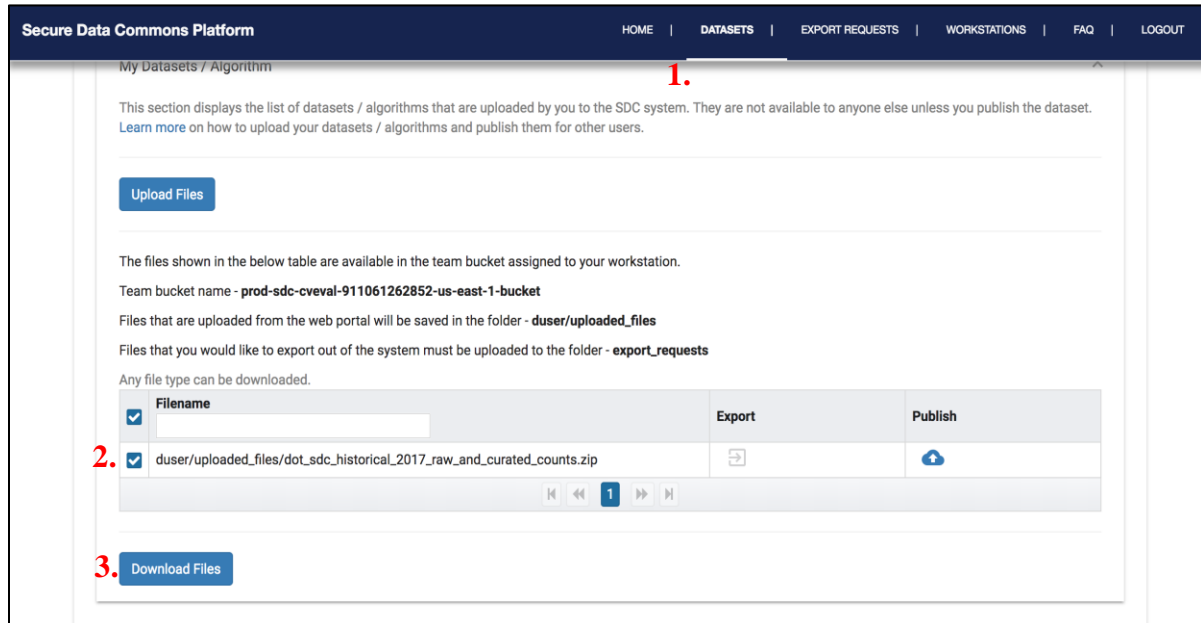


Figure 10: Selecting Files for Download

Notes:

- Not all the files are downloaded directly. Files with extensions such as .txt, .png, or .pdf will be opened in a separate tab from where they can be downloaded. All other files with extensions like .csv, .zip, etc. can be downloaded directly.
- Files are downloaded individually.
- The Filename box allows searches for partial filenames. This can be used to download all the contents of a sub-folder in an S3 bucket by searching for the sub-folder name and then clicking the box next to Filename to select all objects.

Chapter 3. Accessing and Launching Workstations

Users are assigned cloud-based workstations to perform analysis on the datasets. This section provides a description of how to launch and use these workstations.

Launch Workstations

1. Users can see the assigned workstations by clicking on WORKSTATIONS from the top menu. By default, all the workstations are in an inactive state.
2. Click on Start to start the workstation.

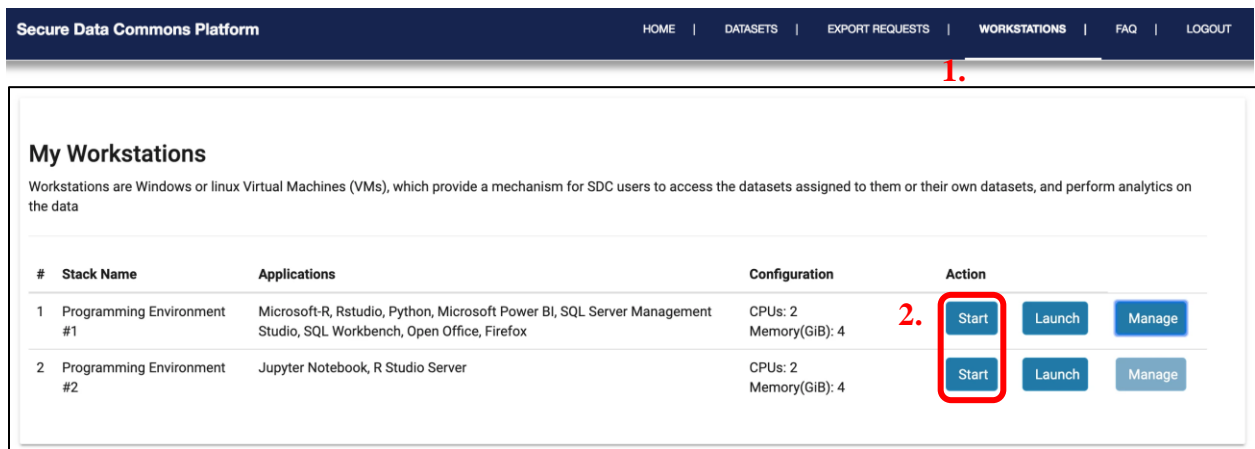


Figure 11: Starting Workstations

3. The workstation should become available within five minutes; you may not see any change immediately. A message will appear when the workstation has been successfully started.

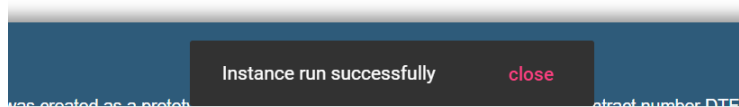


Figure 12: Run Instance Success

- Now click Launch for the workstation.

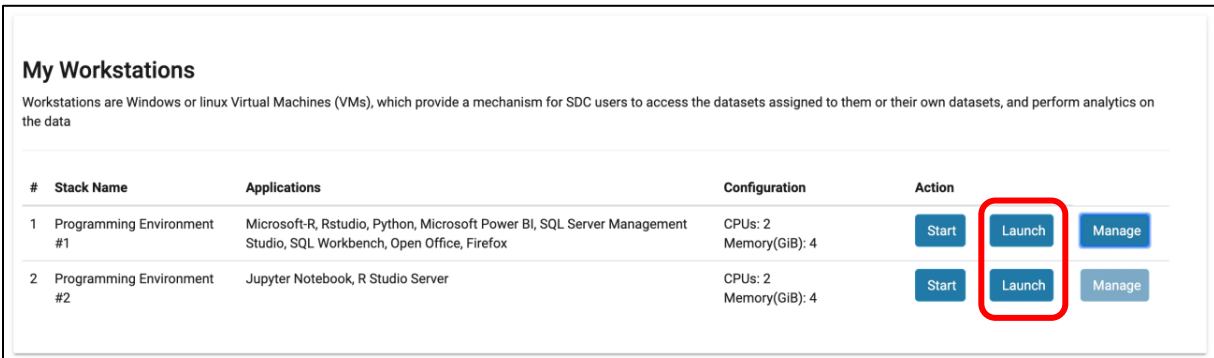


Figure 13: Launch Workstations

- This will provide a user access to their workstation within the browser. The workstation may take a few minutes to initialize. When complete, a login screen will appear. User is prompted to re-enter a valid username and password.

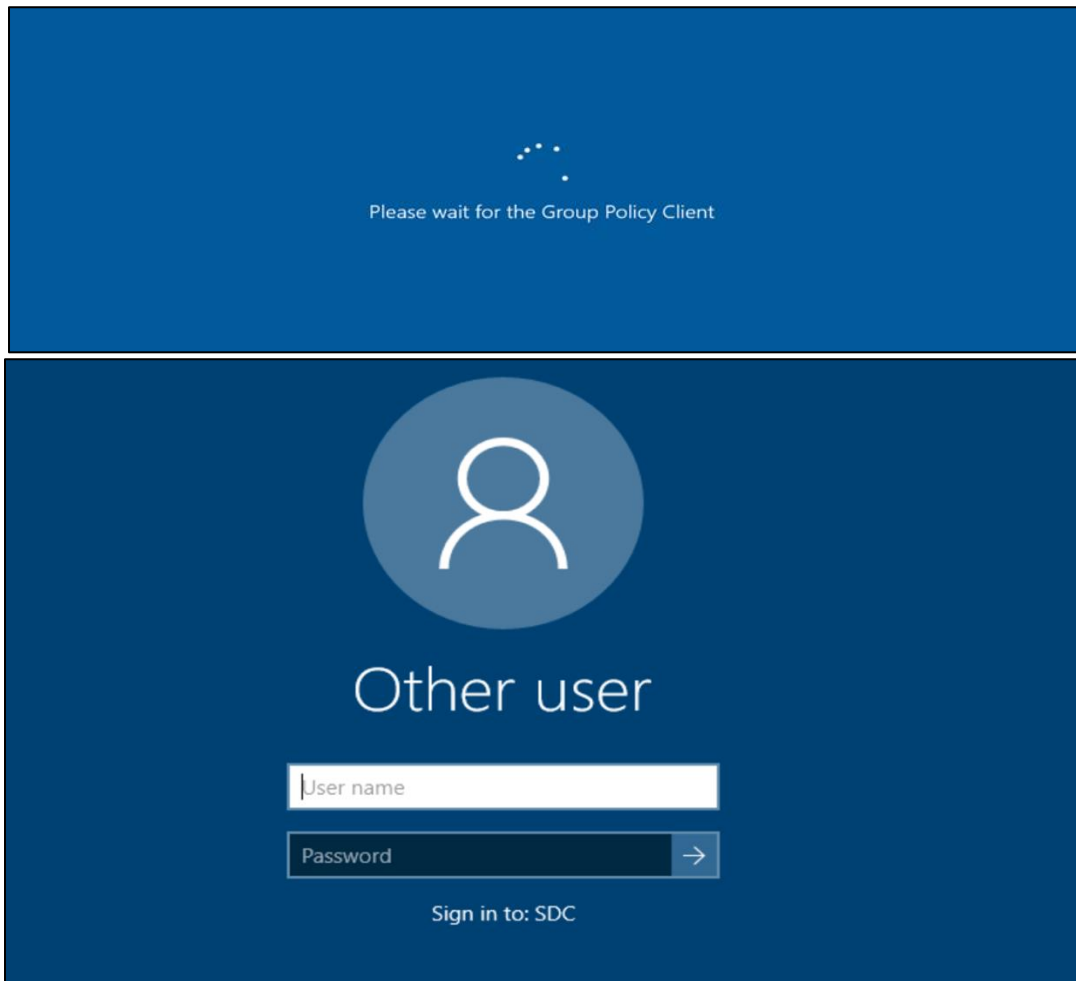


Figure 14: Initialization and Login Screens

Software Validation

By default, users will have the following installed on their workstations:

- Java
- Python
- R, RStudio
- SQL Workbench
- Power BI
- AWS CLI
- Adobe
- Libre Office
- Visual Studio
- PuTTY
- Firefox

To test Python connectivity to the data warehouse, open the IDLE python editor and execute:

```
from impala.dbapi import connect
conn = connect(
    host='172.18.1.20',
    port=10000,
    auth_mechanism='PLAIN',
    user='<your_username>', password='<your_password>')
cursor = conn.cursor()
cursor.execute('SHOW TABLES')
print cursor.fetchall()
```

This should result in an array of tables displayed to the user.

Connecting to the Data Warehouse Using SQL Workbench

The following sections illustrate how the user can connect to the data stores available to the SDC.

Connecting to Waze Data in Redshift

Launch SQL Workbench by double-clicking the SQL Workbench shortcut on the desktop.

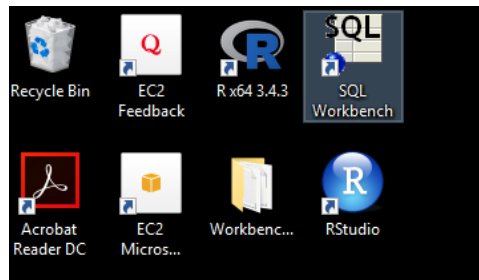


Figure 15: SQL Workbench Icon

Create a Redshift connection profile to connect to Waze data:

1. Create a new connection profile by selecting the top left corner icon on the “Select Connection Profile” window.
2. Select “Amazon Redshift Driver” from the Driver drop-down.
3. Update the URL section with the Redshift URL provided in the email from the support desk detailing Redshift login credentials.
4. Provide your username and password received in the welcome email.
5. Click on the Test button at the bottom to test the connection. A pop-up dialog will appear confirming a successful or failed connection.

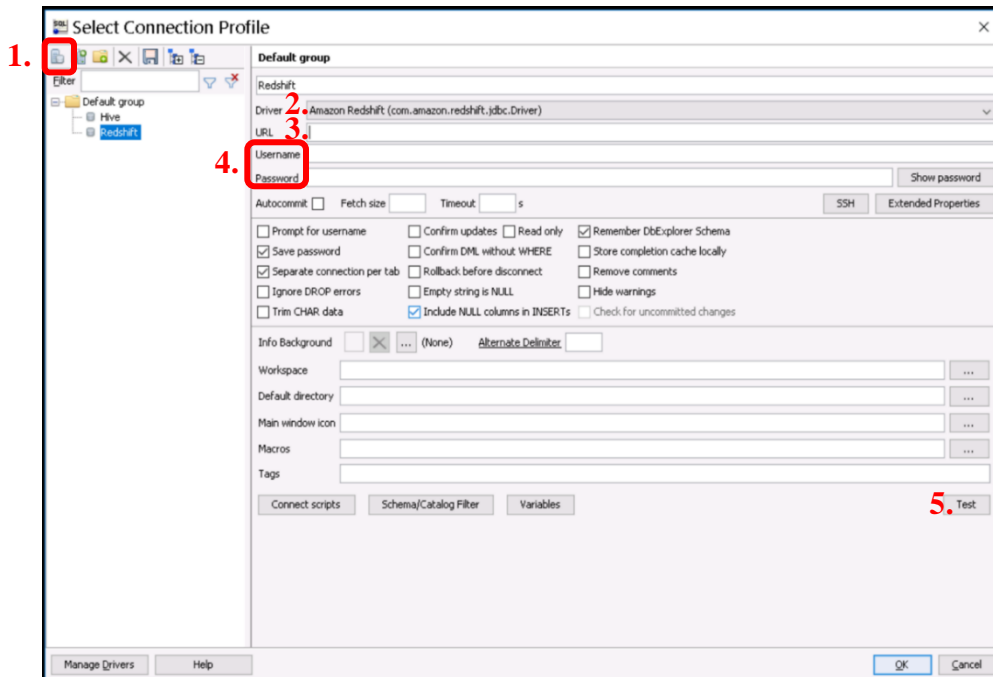


Figure 16: Create Redshift Connection Profile

U.S. Department of Transportation
Intelligent Transportation Systems Joint Program Office

Connecting to CVP Data in Hive Metastore

Launch SQL Workbench by double-clicking on the SQL Workbench shortcut on the desktop:

1. Create a new connection profile by selecting the top left corner icon on the “Select Connection Profile” window.
2. Select “Hive JDBC” from the Driver drop-down.
3. Update URL section with the Hive URL provided in the [cheat sheet](#).
4. Provide your username and password received in the welcome email.
5. Click on the Test button at the bottom to validate your connection. A pop-up dialog will appear confirming a successful or failed connection. If you continue running into a failed connection, contact the [SDC support desk](#) for assistance.

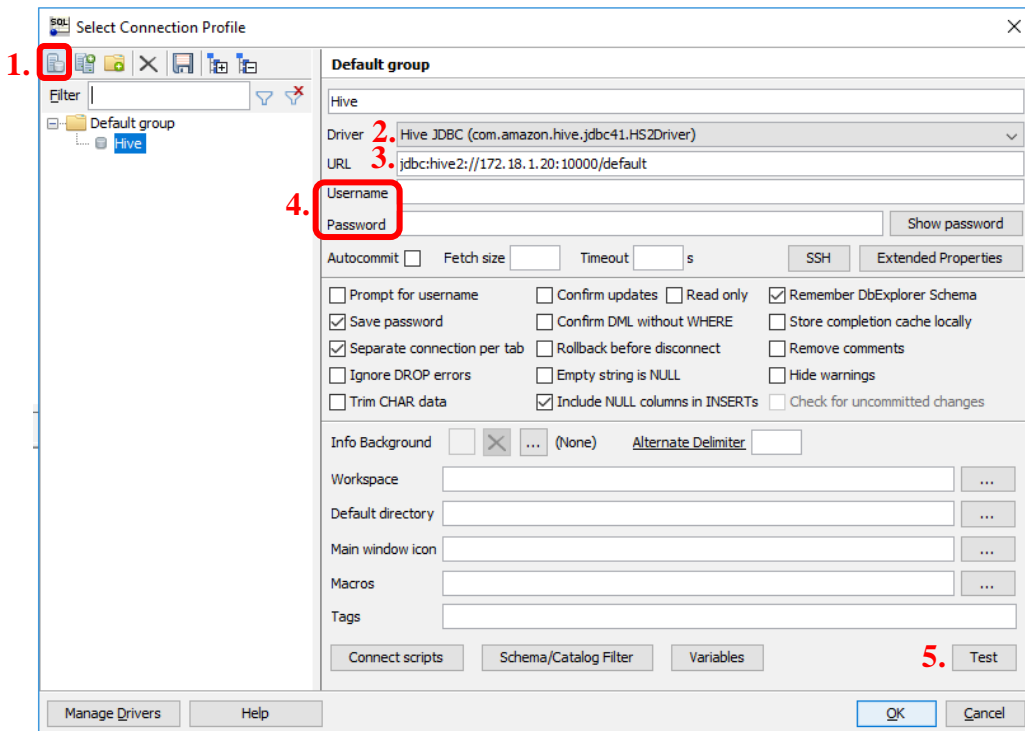


Figure 17: Hive Connection Settings

Update Data Formatting Settings in SQL Workbench

Once the connection has been established, navigate to Tools | Options | Data formatting and update the Decimal digits value to 0.



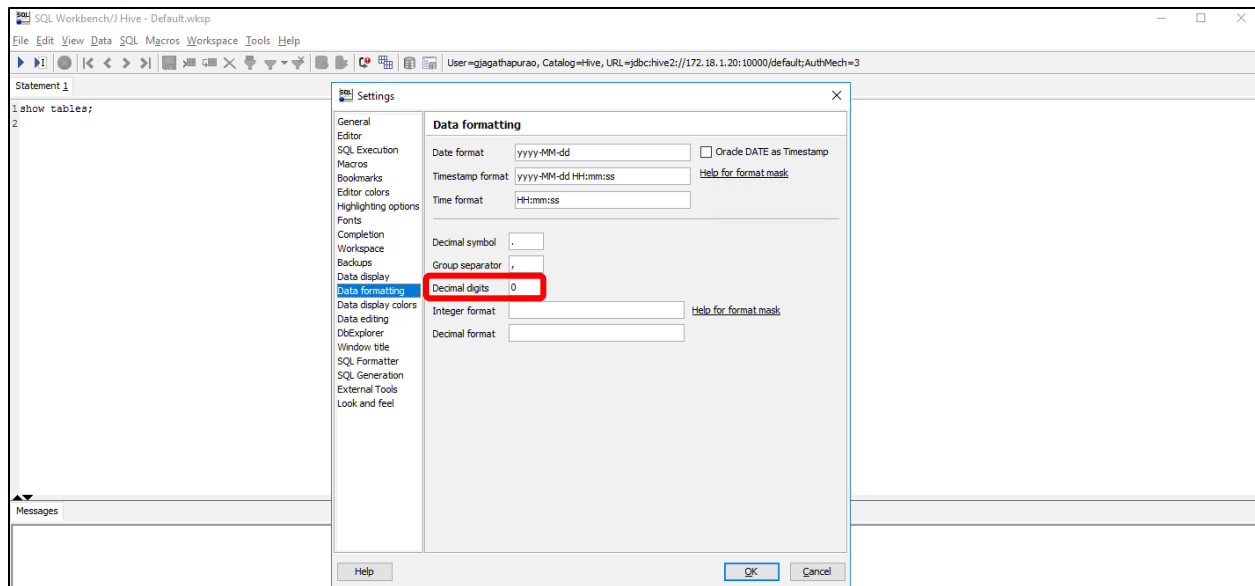


Figure 18: Tools → Options → Data formatting

Connecting to Redshift from Linux Environments

Credentials to access the Waze Redshift database are communicated from the SDC Administrator (support@securedatacommons.com) by a secure email service. Users being granted access will receive an email in the normal email client (such as Outlook) with an “Unlock Message” link. Clicking on this will link to a secure reader that displays the credentials as an email in the web browser.

- In R, it is possible to connect to Redshift using multiple packages. The RPostgreSQL package provides a simple method. This package requires the PostgreSQL library to be installed at the system level; if it is not installed, it would be necessary to install as root in the terminal:
\$ sudo yum install postgresql-devel
- In R, you may need to install.packages("RPostgreSQL", dep=T) if you do not already have the package installed.
- Connect to Redshift using the following code as a guide:

```
library(RPostgres)
# Specify username and password manually, once:
if(Sys.getenv("sdc_waze_username")== "") {
  cat("Please enter SDC Waze username and password
manually, in the console, the first time accessing the
Redshift database, using: \n Sys.setenv('sdc_waze_username'
= <see email from SDC Administrator>) \n
Sys.setenv('sdc_waze_password' = <see email from SDC
Administrator>) ")
}
```

```
redshift_host <- "(details provided by SDC Support to
registered SDC Redshift Users)"
```

U.S. Department of Transportation
Intelligent Transportation Systems Joint Program Office

```
redshift_port <- "5439"
redshift_user <- Sys.getenv("sdc_waze_username")
redshift_password <- Sys.getenv("sdc_waze_password")
redshift_db <- "dot_sdc_redshift_db"

#drv <- dbDriver("PostgreSQL")
conn <- dbConnect(
  RPostgres::Postgres(),
  host=redshift_host,
  port=redshift_port,
  user=redshift_user,
  password=redshift_password,
  dbname=redshift_db)
```

- A database can then be queried using the `dbGetQuery()` function.

Accessing Jupyter Notebook and RStudio Server

Linux users can access their Jupyter Notebook and RStudio Server using the Firefox web browser through windows workstation using below URLs.

- RStudio – <http://<username>-workspace.securedatacommons.internal:8787>
- Jupyter Notebook – <http://<username>-workspace.securedatacommons.internal:8888>



Windows users can click on the “RStudio” shortcut icon present on the desktop to open RStudio console.

Manage Workstations

After launching their workstations, users can manage resizing CPU/RAM and scheduling uptime for a workstation by clicking on its Manage button as shown below.

| My Workstations | | | | |
|--|----------------------------|---|---------------------------|----------------------------|
| Workstations are Windows or linux Virtual Machines (VMs), which provide a mechanism for SDC users to access the datasets assigned to them or their own datasets, and perform analytics on the data | | | | |
| # | Stack Name | Applications | Configuration | Action |
| 1 | Programming Environment #1 | Microsoft-R, Rstudio, Python, Microsoft Power BI, SQL Server Management Studio, SQL Workbench, Open Office, Firefox | CPUs: 2 Memory(GiB): 4 | Start Launch Manage |
| 2 | Programming Environment #2 | Jupyter Notebook, R Studio Server | CPUs: 2 Memory(GiB): 4 | Start Launch Manage |

Figure 19: Manage Workstation

A dialogue window appears with two checkbox options:

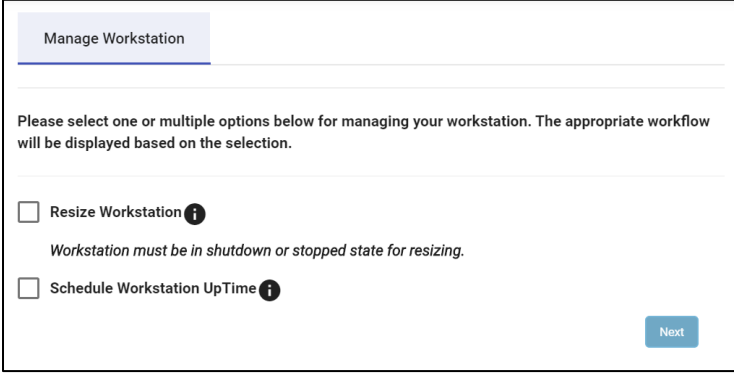
A dialog box titled "Manage Workstation" with a subtitle "Please select one or multiple options below for managing your workstation. The appropriate workflow will be displayed based on the selection." It contains two unchecked checkboxes: "Resize Workstation" and "Schedule Workstation UpTime". Each checkbox has an information icon (i) to its right. Below the "Resize Workstation" checkbox, there is a note: "Workstation must be in shutdown or stopped state for resizing." A "Next" button is located at the bottom right of the dialog box.

Figure 20: Manage Workstation Options

Selecting each option renders the appropriate tabs in the dialogue window. The **i** icon shown next to each option provides an informational tooltip on their functions.

Resize Workstation

1. To resize the workstation, select the checkbox for Resize Workstation and then Next to continue.

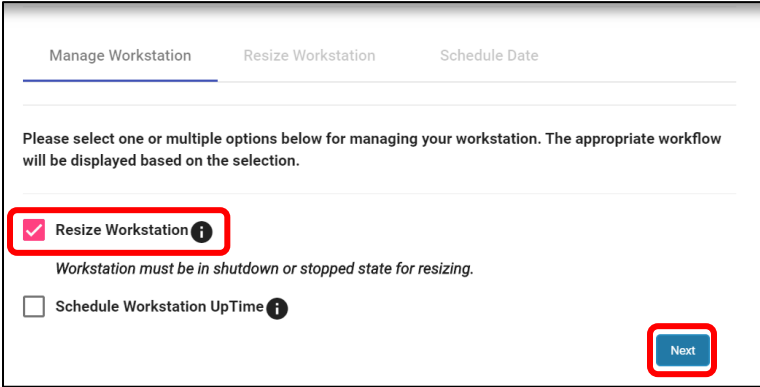
The dialog box now shows three tabs: "Manage Workstation", "Resize Workstation", and "Schedule Date". The "Resize Workstation" tab is selected. The "Resize Workstation" checkbox is now checked and highlighted with a red box. The "Schedule Workstation UpTime" checkbox remains unchecked. The "Next" button at the bottom right is also highlighted with a red box. The instructional text and note about the workstation state remain the same.

Figure 21: Resize Workstation Option

2. A message is shown at the bottom of the screen indicating that the workstation will be stopped before applying the resize.

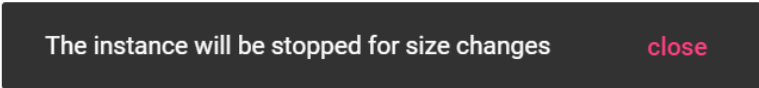
A dark gray message box with white text that reads "The instance will be stopped for size changes". To the right of the text is a red "close" button.

Figure 22: Workstation Stopped for Resize Changes

3. The Resize Workstation tab allows users to select desired CPU/RAM for their workstation. Current configurations will be grayed out and unavailable. Users can also explore pricing details using the link provided under “click here.”
4. Select the “Please start my workstation after resizing to the new configuration” checkbox to automatically start the workstation with the new configuration after saving changes.
5. Select Submit after all details are entered.
6. A Recommended List of instances will appear. Select the desired instance and then the Next button.

The screenshot shows the 'Resize Workstation' interface. At the top, there are three tabs: 'Manage Workstation', 'Resize Workstation' (which is active), and 'Schedule Date'. Below the tabs, a message states: 'Your current workstation configuration is 2 CPU and 4 GB RAM'. A note follows: 'Please choose the appropriate instance type by selecting the desired CPU and RAM. For pricing related details click here.' (Annotation 3 points to 'click here'). Another note says: 'Note: Your workstation will be stopped, if it's currently running in order to resize the workstation. Please save your work before requesting the resize for workstation.' Below this, there's a section 'Select desired CPU / Memory'. It has two dropdown menus: 'CPU' with '4' selected (Annotation 4 points to the dropdown) and 'RAM' with '8 GB' selected (Annotation 5 points to the dropdown). To the right of these is a 'Submit' button. Below the selection area is a checkbox labeled 'Please start my workstation after resizing to the new configuration' (Annotation 4 points to the checkbox). Underneath is a 'Recommended list' section with a table. The table has four columns: 'Instance Name', 'CPU', 'Memory', and 'COST'. One instance is listed: 'c5.xlarge' with CPU '4', Memory '8 GB', and COST '\$0.35 per hour' (Annotation 6 points to the instance name). Below the table is a 'Full List' link. At the bottom right are 'Cancel' and 'Next' buttons.

Resize Workstation

Manage Workstation Resize Workstation Schedule Date

Your current workstation configuration is 2 CPU and 4 GB RAM

Please choose the appropriate instance type by selecting the desired CPU and RAM. For pricing related details [click here](#).

Note: Your workstation will be stopped, if it's currently running in order to resize the workstation. Please save your work before requesting the resize for workstation.

Select desired CPU / Memory

CPU: 4 RAM: 8 GB **Submit**

☒ Please start my workstation after resizing to the new configuration

Recommended list

| Instance Name | CPU | Memory | COST |
|---------------|-----|--------|-----------------|
| c5.xlarge | 4 | 8 GB | \$0.35 per hour |

[Full List](#)

Cancel **Next**

Figure 23: Resizing Options

7. On the Schedule Date tab, users are prompted to enter a date range for how long the resize should last for the workstation instance. Enter the From and To dates and then select Submit.

The screenshot shows a web interface titled "Select schedule". It has three tabs: "Manage Workstation", "Resize Workstation", and "Schedule Date". The "Schedule Date" tab is active. Below the tabs, there is a message: "Please select the schedule between what dates you would want the workstation to be in resized state. An email notification will be sent after the resize is completed." Underneath, there is a section titled "Workspace resize schedule". It contains two date pickers: "From date" with the value "10/5/2019" and "To date" with the value "10/11/2019". At the bottom right, there are two buttons: "Cancel" and "Submit".

Figure 24: Schedule Resize

8. Users will be returned to the Workstations tab with updated CPU and memory information. They will also receive a success email message from the system confirming the resize expiration date.

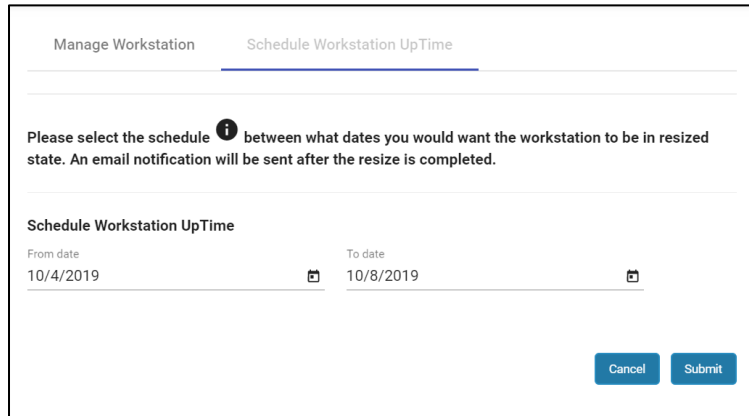
Schedule/Extend Uptime

1. By default, all workstations are shut down at 11 pm EST. If you want to schedule your workstations to be up for a longer period to accommodate analysis runs, select the checkbox for Schedule Workstation Uptime and then Next to continue.

The screenshot shows a web interface titled "Schedule Workstation UpTime". It has two tabs: "Manage Workstation" and "Schedule Workstation UpTime". The "Schedule Workstation UpTime" tab is active. Below the tabs, there is a message: "Please select one or multiple options below for managing your workstation. The appropriate workflow will be displayed based on the selection." Underneath, there are two checkboxes: "Resize Workstation" (unchecked) and "Schedule Workstation UpTime" (checked). Below the "Schedule Workstation UpTime" checkbox, there is a note: "Workstation must be in shutdown or stopped state for resizing." At the bottom right, there is a "Next" button.

Figure 25: Schedule Workstation Uptime Option

2. The Schedule Workstation Uptime tab allows users to enter a date range for how long the workstation uptime should last to skip shutdown. Enter the From and To dates and then select Submit.



The screenshot shows a web interface with two tabs: 'Manage Workstation' and 'Schedule Workstation UpTime'. The 'Schedule Workstation UpTime' tab is active. Below the tabs, there is a message: 'Please select the schedule between what dates you would want the workstation to be in resized state. An email notification will be sent after the resize is completed.' Below this message is a section titled 'Schedule Workstation UpTime' with two date pickers: 'From date' and 'To date'. The 'From date' is set to '10/4/2019' and the 'To date' is set to '10/8/2019'. At the bottom right of the form are two buttons: 'Cancel' and 'Submit'.

Figure 26: Schedule Uptime

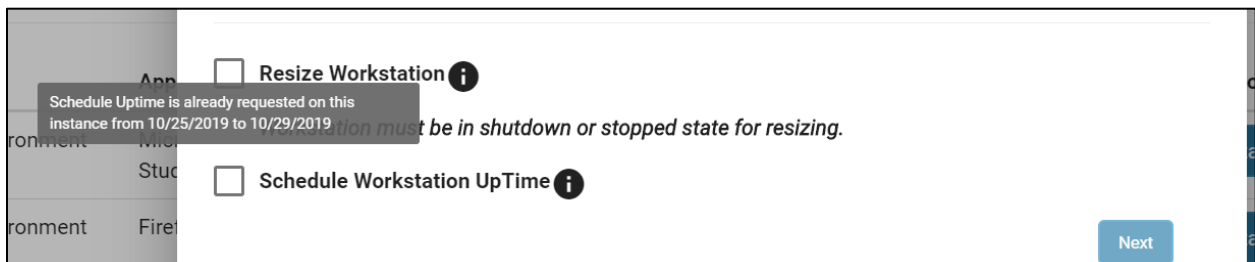
3. To extend any currently scheduled uptime for the workstation, select the Workstations tab and then select Manage again for the workstation. A new tooltip is now shown for the Schedule Workstation Uptime checkbox on mouse hover that indicates previously scheduled uptime.



The screenshot shows a tooltip for the 'Schedule Workstation UpTime' checkbox. The tooltip text reads: 'Schedule Uptime is already requested on this instance from 10/25/2019 to 10/26/2019. Instance must be in shutdown or stopped state for resizing.' The tooltip is displayed over the 'Schedule Workstation UpTime' checkbox, which is currently unchecked. A 'Next' button is visible in the bottom right corner of the interface.

Figure 27: Tooltip with Existing Scheduled Uptime

4. Repeat steps 1-2. For step 2, the From date will already include the date from the previously scheduled uptime. Add a new To date later in the calendar and then submit the update. The previously scheduled uptime goes inactive while the new one becomes active.
5. After selecting Submit, return to the Workstations tab and then select Manage for the workstation. The tooltip shown on hover for the Schedule Workstation Uptime checkbox now displays the extended uptime.



The screenshot shows a tooltip for the 'Schedule Workstation UpTime' checkbox, similar to Figure 27, but with an extended schedule. The tooltip text reads: 'Schedule Uptime is already requested on this instance from 10/25/2019 to 10/29/2019. Instance must be in shutdown or stopped state for resizing.' The tooltip is displayed over the 'Schedule Workstation UpTime' checkbox, which is currently unchecked. A 'Next' button is visible in the bottom right corner of the interface.

Figure 28: New Tooltip with Extended Uptime Schedule

Stop Workstations

Users can see the assigned workstations by clicking on the workstations tab on the top right corner of the page. By default, all the workstations are scheduled to stop every day at 11 PM EST. Users can stop the workstations manually by clicking on the Stop button as shown below. A message will appear when the instance is successfully stopped.

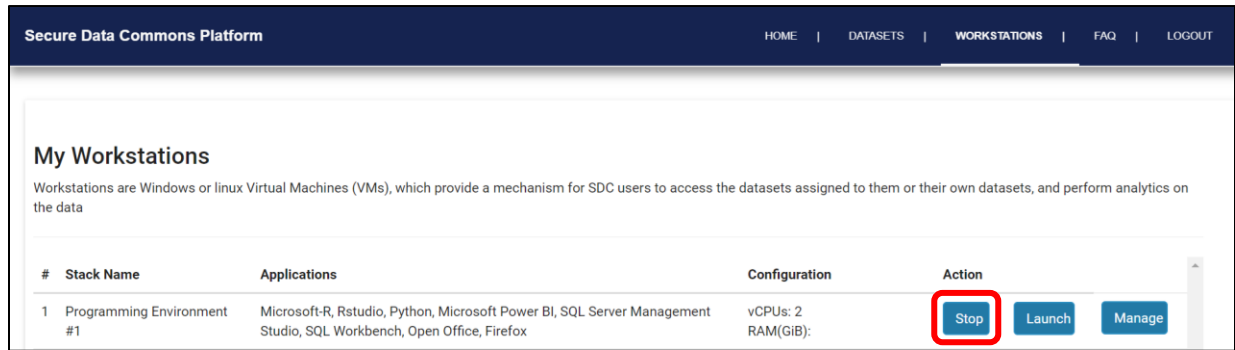


Figure 29: Stop Workstation

Chapter 4. Sample Queries for SDC Datasets

The following section provides some sample queries for anchor datasets hosted by the SDC. Two datasets which are currently available in the SDC system:

- Data from USDOT-sponsored Connected Vehicle Pilot (CVP) program, provided by the following pilot sites:
 - Tampa Hillsborough Expressway Authority (THEA)
 - Wyoming Department of Transportation (WYDOT)
- Data provided by Waze on traffic jams and alerts

Sample queries are provided for each of the two datasets.

Connected Vehicle Data

Overview

The CVP Data Warehouse is built on top of a Hadoop's HDFS cluster and utilizes Hive with a HiveQL querying language as a front-end querying package. Information about the HiveQL language and the language manual can be found online:

- Wikipedia info on Apache Hive: https://en.wikipedia.org/wiki/Apache_Hive
- Language manual: <https://cwiki.apache.org/confluence/display/Hive/LanguageManual>

Key components of the CVP dataset are the basic safety messages (BSM) and Traveler Information Messages (TIM). Both THEA and WYDOT data include BSM and TIM.

WYDOT

Basic Safety Message (BSM) Relational Tables

WYDOT BSM data are stored in the following set of relational tables.

```
wydot_bsm_core
wydot_bsm_partii
wydot_bsm_partii_crumbdata
```

The following figure illustrates the foreign key relationship of WYDOT BSM tables.

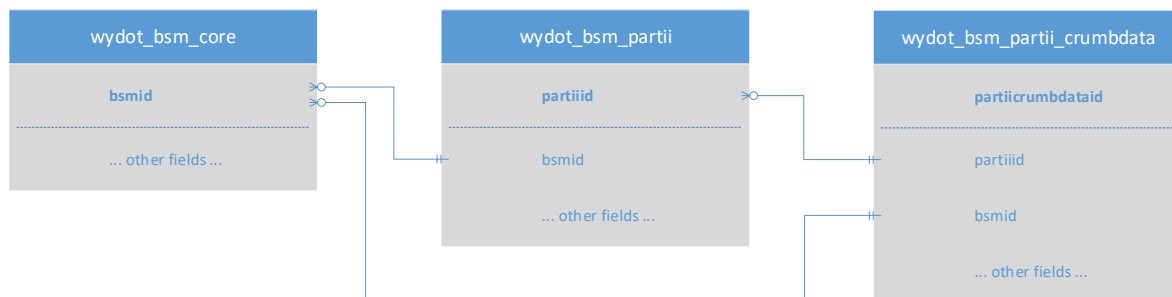


Figure 30: WYDOT BSM Tables

U.S. Department of Transportation
Intelligent Transportation Systems Joint Program Office

Sample queries using WYDOT BSM tables.

These queries can be run using the SQL Workbench which is configured to connect to the data warehouse. Note: running the last three queries will take longer, because they are joining multiple tables.

```
select * from wydot_bsm_core limit 5;
select * from wydot_bsm_partii limit 5;
select * from wydot_bsm_partii_crumbdata limit 5;

select coredatalatitude, coredatalongitude from
wydot_bsm_core limit 5;

select * from wydot_bsm_core join wydot_bsm_partii where
wydot_bsm_core.bsmid = wydot_bsm_partii.bsmid limit 5;

select * from wydot_bsm_partii join
wydot_bsm_partii_crumbdata where wydot_bsm_partii.partiiid
= wydot_bsm_partii_crumbdata.partiiid limit 5;

select * from wydot_bsm_core join wydot_bsm_partii join
wydot_bsm_partii_crumbdata where wydot_bsm_core.bsmid =
wydot_bsm_partii.bsmid and wydot_bsm_partii.partiiid =
wydot_bsm_partii_crumbdata.partiiid limit 10;
```

WYDOT BSM and TIM Metadata RECORDGENERATEDAT vs ODERECEIVEDAT

WYDOT BSM and TIM messages include a metadata section, appended by the Operational Data Environment (ODE) component. Among others, there are two timestamps of interest are presented in this section:

recordGeneratedAt: Closest time to which the record was created by a Vehicle.

odeReceivedAt: Time ODE received the data in UTC format. This time is the closest to which the CV PEP system received the record.

- Based on real-life conditions, odeReceivedAt may be days or even weeks after the recordGeneratedAt timestamp.
- The following queries result in distribution of recordGeneratedAt by odeReceivedAt times:
- WYDOT BSM messages:

```
select SUBSTR(metadataarecordgeneratedat, 0, 10) as
RECORDGENERATEDAT
, SUBSTR(metadataodereceivedat, 0, 10) as ODERECEIVEDAT
, count(SUBSTR(metadataodereceivedat, 0, 10)) as CNT
from wydot_bsm_core
group by SUBSTR(metadataodereceivedat, 0, 10)
, SUBSTR(metadataarecordgeneratedat, 0, 10)
order by RECORDGENERATEDAT limit 10000;
```

U.S. Department of Transportation
Intelligent Transportation Systems Joint Program Office

- **WYDOT TIM messages:**

```
select SUBSTR(metadataarecordgeneratedat, 0, 10) as  
RECORDGENERATEDAT  
  , SUBSTR(metadataodereceivedat, 0, 10) as ODERECEIVEDAT  
  , count(SUBSTR(metadataodereceivedat, 0, 10)) as CNT  
from wydot_tim  
group by SUBSTR(metadataodereceivedat, 0, 10)  
  , SUBSTR(metadataarecordgeneratedat, 0, 10)  
order by RECORDGENERATEDAT limit 10000;
```

WYDOT Speed Data

There are two WYDOT speed datasets in the CV PEP Data Warehouse: `wydot_speed_unprocessed` and `wydot_speed_processed`. The following sample query displays average vehicle speed distribution by lane and it will work against either of the tables:

```
select lane, avg(speedmph) as speed_average  
from wydot_speed_unprocessed  
group by lane  
order by lane;
```

THEA

THEA Basic Safety Message (BSM) Relational Tables

THEA BSM data are stored in the following set of relational tables.

| THEA BSM Relational Tables |
|--|
| <code>thea_bsm_core</code> <code>thea_bsm_partii</code> <code>thea_bsm_partii_crumbdata</code> |

Sample queries using THEA BSM tables.

These queries can be run using the SQL Workbench which is configured to connect to the data warehouse. Note: running the last three queries will take longer, because they are joining multiple tables.

```
select * from thea_bsm_core limit 5;  
  
select * from thea_bsm_partii limit 5;  
  
select * from thea_bsm_partii_crumbdata limit 5;  
  
select coredatalat, coredatalong from thea_bsm_core limit  
5;
```



```
select * from thea_bsm_core join thea_bsm_partii where
thea_bsm_core.bsmid = thea_bsm_partii.bsmid limit 5;

select * from thea_bsm_partii join
thea_bsm_partii_crumbdata where thea_bsm_partii.partiid =
thea_bsm_partii_crumbdata.partiid limit 5;

select * from thea_bsm_core join thea_bsm_partii join
thea_bsm_partii_crumbdata where thea_bsm_core.bsmid =
thea_bsm_partii.bsmid and thea_bsm_partii.partiid =
thea_bsm_partii_crumbdata.partiid limit 10;
```

Geospatial Queries

- The Data Warehouse has geospatial querying capabilities. Functions such as ST_Point, ST_Polygon, ST_Contains and others can be used in queries. For the full list of supported functions see: <https://github.com/Esri/spatial-framework-for-hadoop/wiki/UDF-Documentation>
- As an example, here is a sample query to retrieve a count of messages generated by vehicles between latitudes 40 and 41 and longitudes between -106 and -105.

```
select count(*) from wydot_bsm_core where
ST_Contains(ST_Polygon(40, -105, 41, -105, 41, -106, 40, -
106),
ST_Point(coredatalatitude, coredatalongitude));
```

Waze Data

Overview

The SDC platform ingests and curates Waze data for all 50 states of United States of America using the Waze API. Users have access only to Waze data pertaining to the specific geographic region specified in their data access policies.

Waze Alert Data

An alert is a User Generated Incident (UGI) reported by a Waze user or group of users, as defined by Waze. In SQL Workbench, the following query can be used to test access to the *alert* table.

```
select distinct alert_uuid from alert where
alert_type='<alert_type>' and pub_utc_timestamp between
<start_time_stamp> and <end_time_stamp> and state
=<state_name>;

select distinct alert_uuid from alert where
alert_type='ACCIDENT' and pub_utc_timestamp between '2018-
```

```
01-01 04:59:43.00' and '2018-01-30 04:59:43.00' and state  
='CA';
```

Waze Jam Data

Waze provides information of traffic jams and events that affect road conditions either from wazers or external sources. A traffic jam maybe associated with an alert.

```
select top 10 * from jam;  
select top 10 * from jam_point_sequence;
```

Waze Irregularity Data

Irregularities are similar to Jams, where Waze derives these events based on unusual traffic patterns. These could also be a result of an alert or a jam.

```
select top 10 * from irregularity;  
select top 10 * from irregularity_point_sequence;  
select top 10 * from irregularity_alert;  
select top 10 * from irregularity_jam;
```

Chapter 6. Exporting Datasets from the SDC

Data Analysts should be able to export the data of the system based on the compliance and data usage policies set forth by a Data Provider.

There are two different types of analysts:

1. **General Analyst:** This type of analyst must provide justification to the Data Provider for each data product that they want to export out of the SDC system. The intent is to ensure that the Data Provider has oversight of the exported data. This type of analyst can also request trusted status from the Data Provider while filling out the approval form.
2. **Trusted Analyst:** This type of analyst already has a trusted status which is provided by the Data Providers. The intent is to reduce the effort for exporting data products of analyses out of the SDC system. A trusted user has a pre-existing and approved relationship with the Data Provider.

Once the Data Analyst completes creating derived datasets, either working on the SDC datasets or combining with other datasets that they import into the system, they can export the derived datasets or share the datasets with other team members.

The following are the steps that the Data Analyst needs to follow to export the data of their analysis from the SDC system to support their research:

1. Each Data Analyst is part of a team bucket which is displayed in the Datasets section. When ready to export, Data Analysts can select the file (or files) that they want to export out of the SDC system and place them in a separate staging folder (i.e., **export_requests**) in their team bucket. Data Analysts can request for exporting a file in this folder by clicking on the export symbol for the file they want to export out of the SDC system.

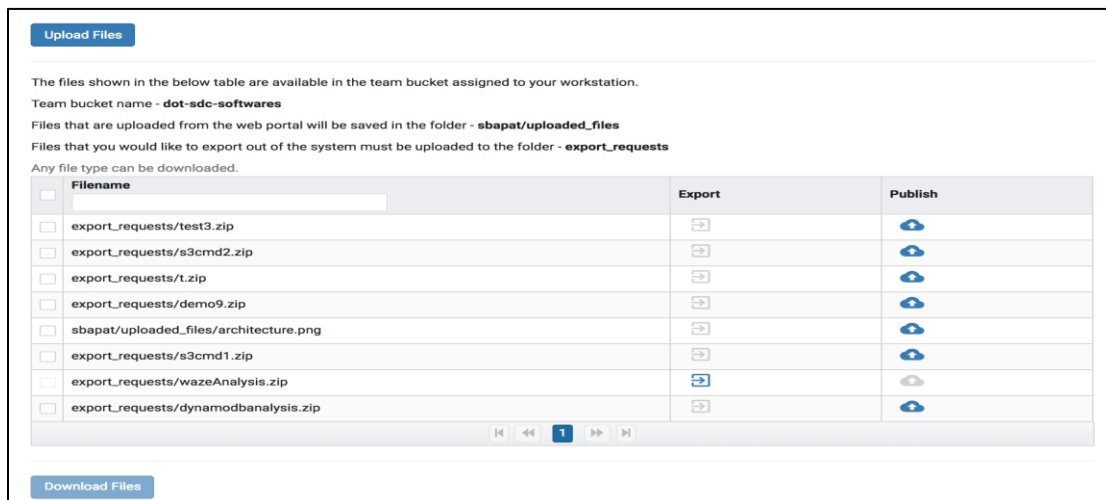


Figure 31: Request Export

U.S. Department of Transportation
Intelligent Transportation Systems Joint Program Office

2. Once the export button is selected, a dialog box for requesting the export data will be displayed. The analyst will then need to provide the details of the Project, Data Provider, and Data Type that he has used to create his own dataset and click on the NEXT button once finished.

Request to Export Data

Select ProjectApproval FormTrusted Status

Please select the Project, Data Provider, and the primary Sub-Dataset/Data Type, that you have used to create your derived dataset. This will help us to route your request to the appropriate Data Provider for approval.

Note - If your derived Dataset is created using multiple Sub-Datasets/Data Types, that are available within SDC or external Datasets/Data Sources that you have uploaded into the system, you will be provided an option to list all such Datasets/Data Types in the next section of the workflow.

Project/Dataset

WAZE

Data Provider

WAZE

Sub-Dataset/Data Type

JAM

CANCELNEXT

Figure 32: Request Export Form

- The additional information regarding the request for exporting the data must be filled out in the approval form below. These details are shared with the Data Providers, which helps them to accept or reject the request made by the Data Analysts.

| Select Project | Approval Form | Trusted Status |
|--|---------------|----------------|
| <p>Please provide additional information pertaining to your request for exporting data, by filling out the fields below. These details will be shared with the Data Provider to help them review your request and provide their decision.</p> <p>* All fields are mandatory</p> <p>Name or short description of your derived dataset *</p> <p>waze derived dataset</p> <p>Anchor dataset of interest or data provider *</p> <p>WAZE</p> <p>Specific sub-datasets or data types used *</p> <p>JAM</p> <p>Additional datasources *</p> <p>datasources</p> <p>High level description of derived dataset *</p> <p>high level desc</p> <p>Detailed description of the derived dataset *</p> <p>detailed desc</p> <p>Tags</p> <p>tags</p> <p>Justification of Export *</p> <p>Justification </p> | | |

Figure 33: Approval Form Fields

4. If the user is not a trusted user, he/she will be prompted with the option for requesting the trusted status from the Data Provider. This will allow the analyst to export the data immediately, as opposed to waiting for review and approval from the Data Provider. The user must accept the Acceptable Usage policy for the request to go through to the Data Provider. The form will not be submitted if the user declines.

Select Project

Approval Form

Trusted Status

Trusted Status is a mechanism for analysts to obtain a passport from a data provider. Obtaining this passport allows analyst to export their data immediately (for subsequent similar requests), as opposed to waiting for the review and approval of a data provider.
This status is acquired per Project + Data Provider + Sub-Dataset/Data Type.

Note - Based on the dataset and datatype selection, you currently do not have a Trusted Status from this Data Provider. We will notify the Data Provider about your request and send it for approval. Your request will be processed based on the decision from the Data Provider.

Do you wish to request Trusted Status from the Data Provider?

☐ Yes
 ☒ No

Acceptable Use Policy

The WAZE DOT is providing ongoing access to data generated by the Connected Vehicle Pilot deployment to support performance measurement and evaluation activities to a select group of explicitly approved individuals. The CV Pilot is an ongoing research activity and includes access to rapidly evolving data sets and products. WAZE DOT makes no claims, promises or guarantees about the accuracy, completeness, or adequacy of the contents of data and expressly disclaims liability for errors and omissions in the data.

Conducting research activities on WAZE DOT CV pilot data and resources is restricted to authorized individuals for the purpose for which access was granted. Further users of the WAZE CV Pilot data

☐ Accept
 ☐ Decline

Figure 34: Acceptable Use Policy

5. Upon successful submission, the request will be sent to appropriate Data Providers. Data Providers will be responsible for accepting or rejecting the export requests.
6. Once Data Providers approve the request, Data Analysts will be able to download the dataset out of SDC through portal.

Chapter 7. Technical Documentation and Contact Information

The following sections provide technical resources for SDC users.

Architecture Diagram

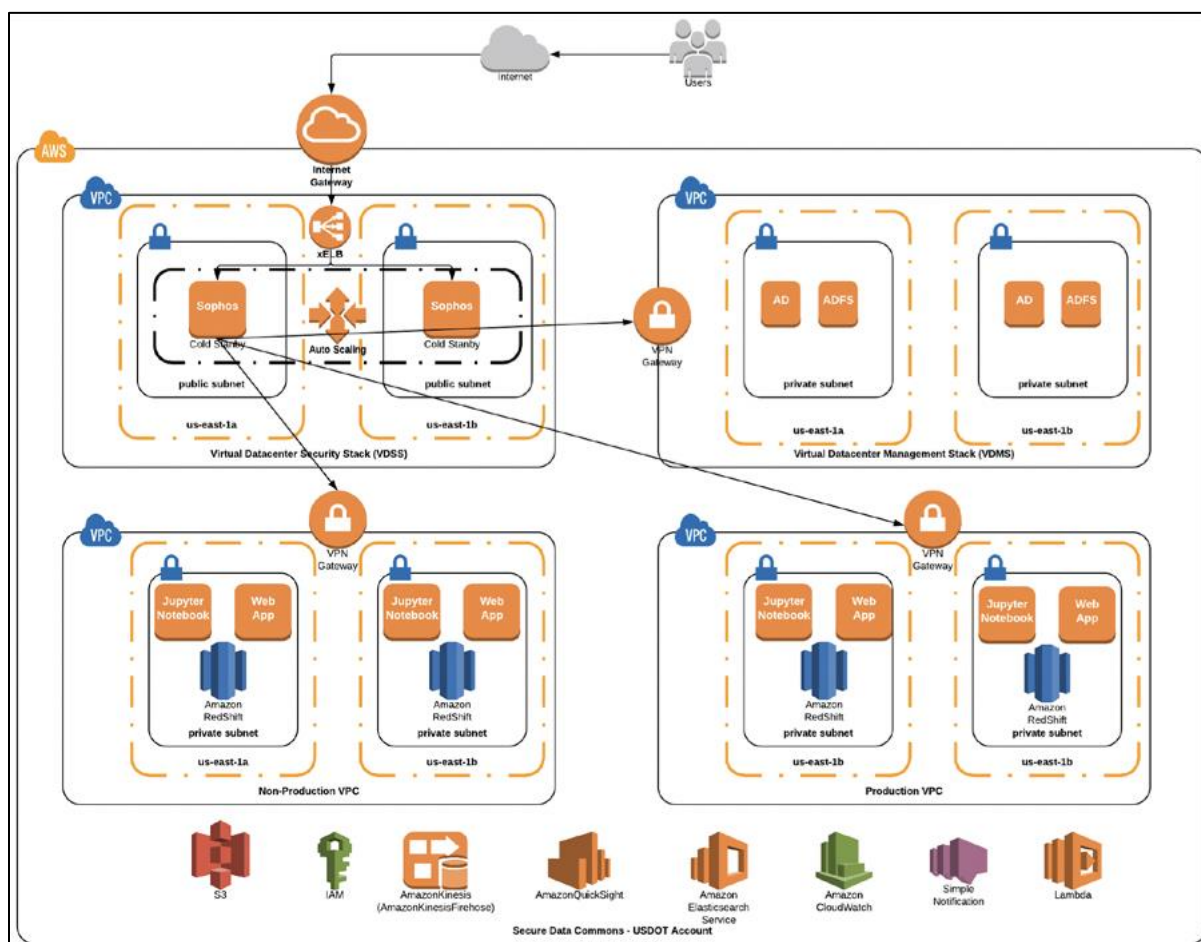


Figure 35: SDC Architecture Overview

Workstation Details

Table 1: Default Workstation Details

| Workstation Type | Type |
|---------------------|-----------|
| Linux Workstation | t2.medium |
| Windows Workstation | t2.medium |

Note: Workstation size and type can be increased upon user request.

Tools and Versions

Table 2: List of Tools Used and Their Versions

| Tool Name | Version | Workstation |
|----------------|-----------|----------------|
| Java | 1.8.0_162 | Linux, Windows |
| Python | 2.7.14 | Linux, Windows |
| SQLWorkbench/J | Build 125 | Windows |
| R | 3.4.3 | Linux, Windows |
| RStudio | 1.1.423.0 | Linux, Windows |
| Libre Office | 5.3.6.1 | Windows |
| Visual Studio | 1.20.1.0 | Windows |
| AWS CLI | 1.14.46 | Windows |
| 7Zip | 18.01 | Windows |
| PuTTY | 0.70 | Windows |
| Firefox | 59.0.2.0 | Windows |

Contact Information

SDC support team can be reached out at support@securedatacommons.com

Useful Links

- S3:** Simple Storage Service, a place to store data.
- Jupyter:** An interactive, browser-based programming environment, mostly used for Python scripts but can also run R or other languages and can weave formatted text in with code and results of code into one ‘notebook’ file.
- Redshift:** A database system, which can be queried with SQL.
- Hive:** The Apache Hive™ data warehouse software facilitates reading, writing, and managing large datasets residing in distributed storage and queried using SQL syntax.

AWS S3 CLI Commands

The following are a list of helpful commands to work with S3 from the terminal. The AWS_S3_CLI_Cheat_Sheet.pdf file is also available on the desktop of all workstations with useful commands. In these commands, ‘local’ refers to the EC2 instance being run inside SDC:

- **List objects in a bucket** - If there are any files at the bucket level then the below command will return the list of files. If there are only folders/prefixes under the bucket then it will return the top level folder/prefix names of that bucket
aws s3 ls s3://<bucketName>
- **List objects under a folder/prefix** - The below command will list all the objects/files under that folder or prefix.
aws s3 ls s3://<bucketName>/<prefix>/

U.S. Department of Transportation
Intelligent Transportation Systems Joint Program Office

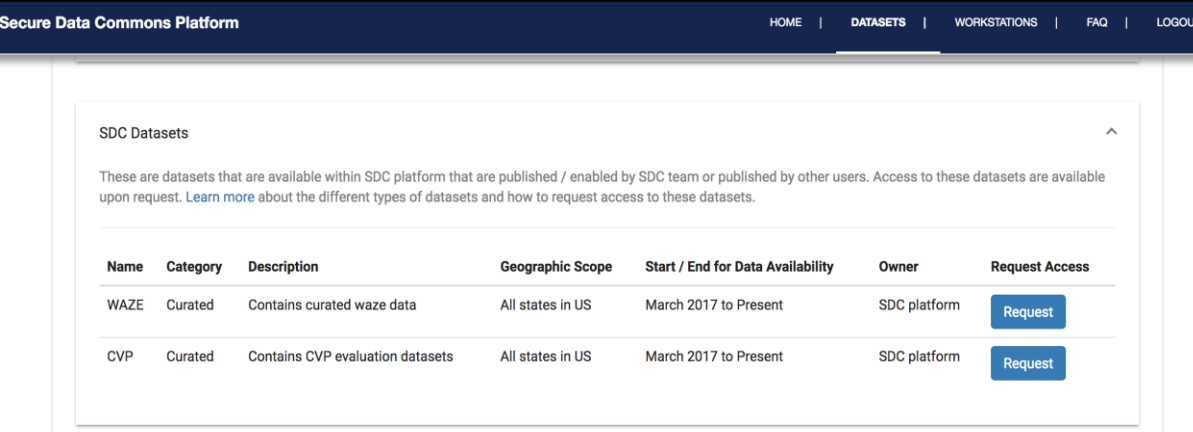
- **Copy a local file to an s3 bucket** - This command will copy the file at the root of the bucket .
`aws s3 cp nohup.out s3://<bucketName>/`
- **Copy a local file to a specific folder/prefix in an s3 bucket** - To create a new output folder and copy to that destination, simply add the desired path to the end, following bucket/.
`aws s3 cp <filename> s3://<bucketName>/<prefix>/`
- **Get a file from S3 to the ec2 instance** - This command will help you get a file from S3 to the EC2 instance:
`aws s3 cp s3://<bucketName>/<prefix>/<fileName> .`

Chapter 8. Frequently Asked Questions

- [How can I get access to the SDC Datasets?](#)
- [How will I understand what a particular dataset consists of?](#)
- [How can I launch a workstation?](#)
- [Where do I store my data?](#)
- [How can I bring my own datasets/algorithm to my workstation?](#)
- [How can I publish my dataset/algorithm?](#)

How can I get access to the SDC Datasets?

These are datasets that are available within SDC platform that are published / enabled by the SDC team or published by other users. Access to these datasets are available upon request. Once you are logged in, go to 'Dataset' in the top menu.



| Name | Category | Description | Geographic Scope | Start / End for Data Availability | Owner | Request Access |
|------|----------|----------------------------------|------------------|-----------------------------------|--------------|-------------------------|
| WAZE | Curated | Contains curated waze data | All states in US | March 2017 to Present | SDC platform | Request |
| CVP | Curated | Contains CVP evaluation datasets | All states in US | March 2017 to Present | SDC platform | Request |

Figure 36: Datasets Tab

Click on ‘SDC Dataset’ to view all the available datasets.

To access a dataset, you need to click on the ‘Request’ button.

A form will pop up. Fill in that form and click on ‘Send Request’ button.

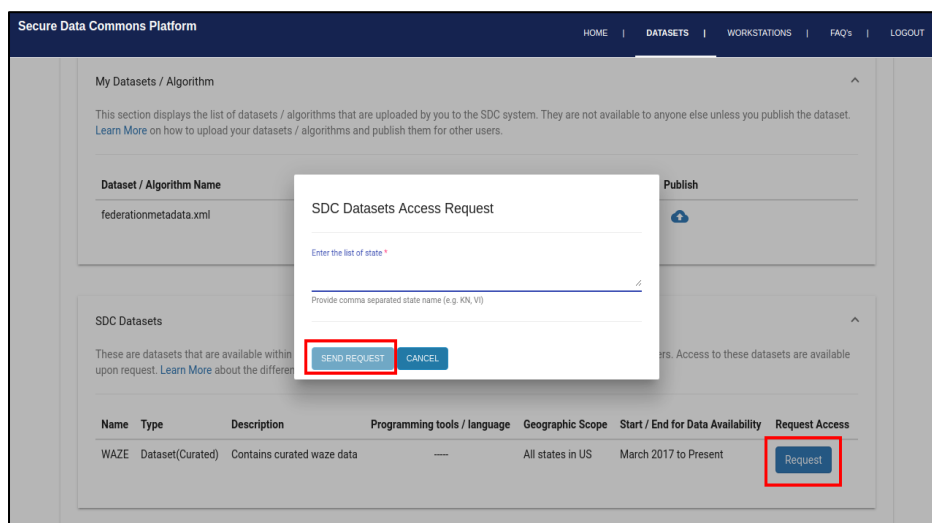


Figure 37: Dataset Access Request

The request will be sent to support team and access to the requested dataset will be given upon approval.

How will I understand what a particular dataset consists of?

Click on Name of Dataset, you can see README of that particular dataset below it.

How can I launch a workstation?

Click on ‘Workstation’ and click on the ‘Launch’ button of any workstation you want to access.

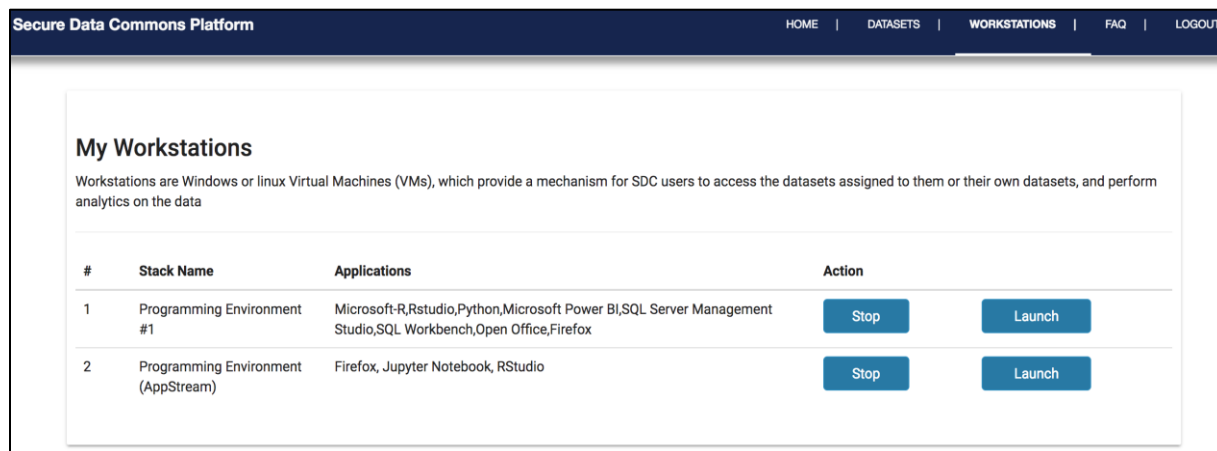


Figure 38: Workstations Tab

For Programming Environment #1, you will be prompted with username and password to log in to the Windows workstation.

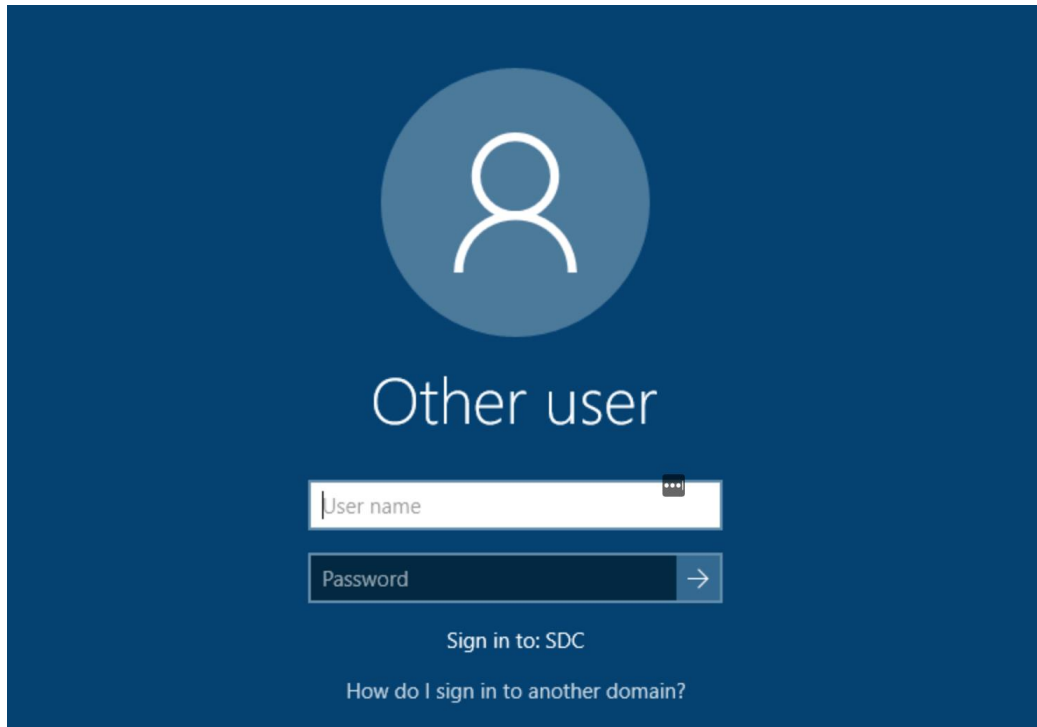


Figure 39: Login for Programming Environment

Where do I store my data?

You can store your data in your team/individual bucket. Please refer to [Upload User Data to S3 Bucket Through Portal](#)

How can I bring my own datasets/algorithm to my workstation?

Please refer to [Upload User Data to S3 Bucket Through Portal](#) to bring your own datasets/algorithm to workstation.

How can I publish my dataset/algorithm?

Follow the below steps to publish your datasets / algorithms and share with other SDC Users.

1. Navigate to the Datasets page.

2. Click on the upload icon under publish for the dataset/algorithm you wish to publish.

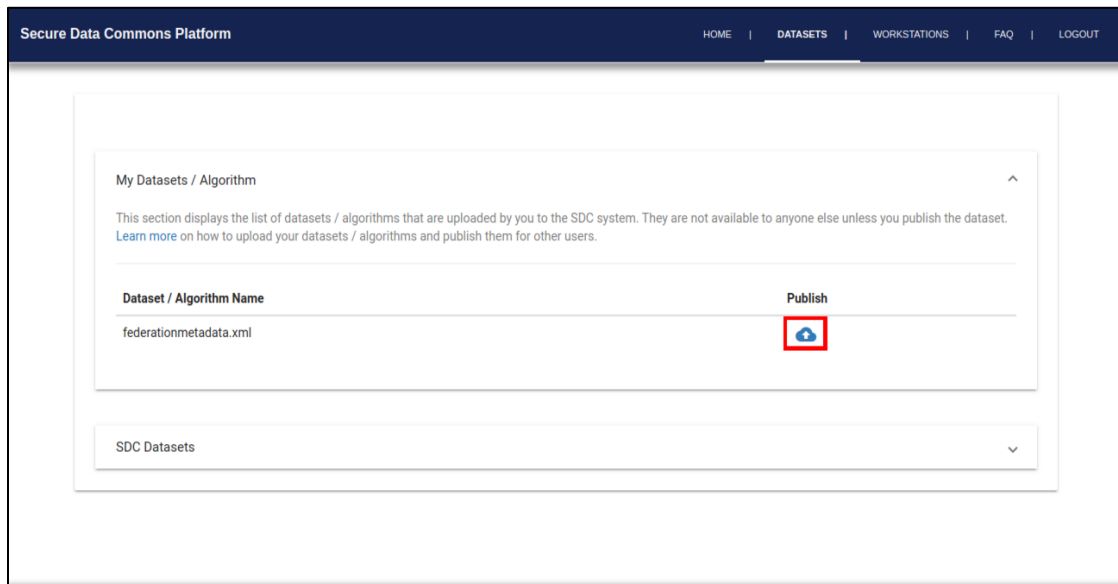


Figure 40: Upload Button

3. In the pop-up window, you will have two options, either a Dataset or Algorithm.
 - a. Select the drop-down value as Dataset and fill in the required values.

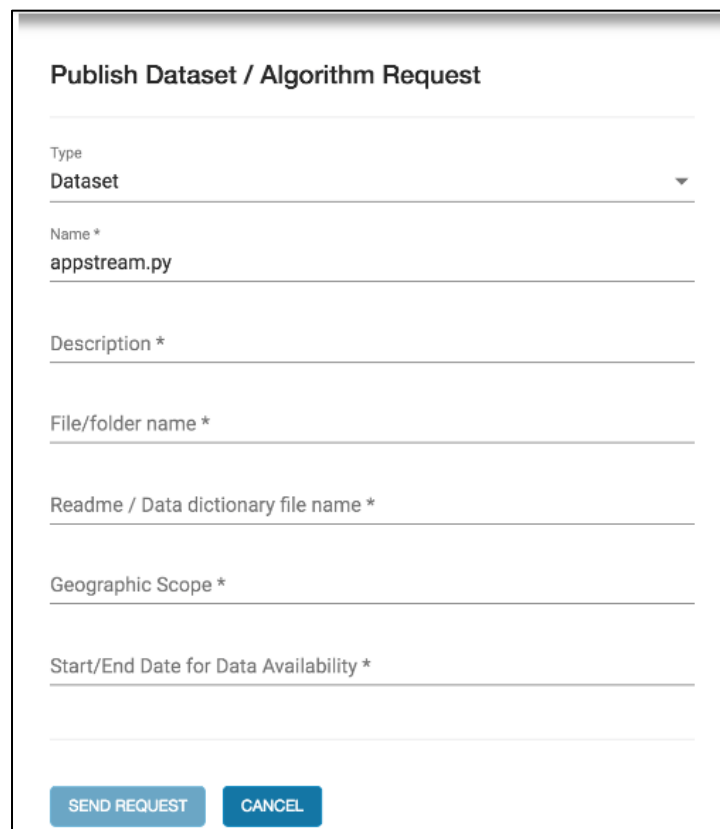
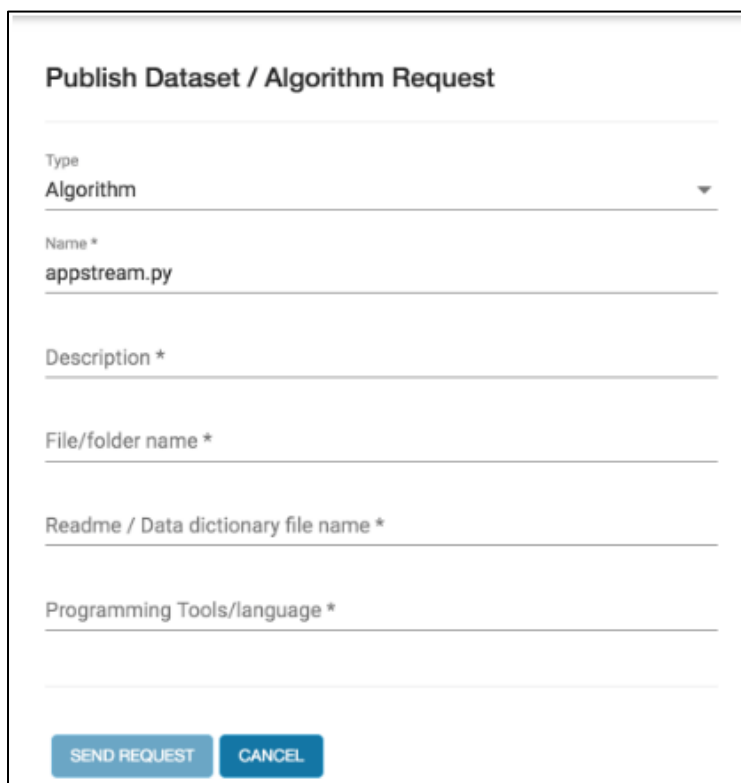
The screenshot shows a 'Publish Dataset / Algorithm Request' form. It has several input fields: 'Type' (a dropdown menu currently showing 'Dataset'), 'Name *' (containing 'appstream.py'), 'Description *', 'File/folder name *', 'Readme / Data dictionary file name *', 'Geographic Scope *', and 'Start/End Date for Data Availability *'. At the bottom of the form are two buttons: 'SEND REQUEST' and 'CANCEL'.

Figure 41: Publish Dataset Form

- i. Name - Name of the dataset, which you wish to call it. Users will see your dataset with this name under SDC Datasets section.
 - ii. Description - Provide a short description so users can get an idea about your dataset.
 - iii. File/folder name - Name of the file or folder where your dataset resides in your S3 Bucket. We need this information, so the support team can publish this dataset and make it available to other users.
 - iv. Readme / Data dictionary file name - This file should provide detailed instructions about your dataset, how it was created or any relevant information that helps user to understand and use the dataset. Save this file in your home folder relative to the dataset file/folder name.
 - v. Geographic scope - Indicate the geographic scope for your dataset whether it belongs to a specific state, region, country etc.
 - vi. Start/End Date for data availability - Provide the start and end dates of the data that belongs in your dataset. For example, your dataset may contain data from March 2017 to August 2017.
- b. Select the drop-down value as Algorithm.



Publish Dataset / Algorithm Request

Type
Algorithm

Name *
appstream.py

Description *

File/folder name *

Readme / Data dictionary file name *

Programming Tools/language *

SEND REQUEST CANCEL

Figure 42: Algorithm Request Form

- i. Name - Enter the name for your algorithm. Users will see your algorithm with this name under SDC Datasets section
- ii. Description - Provide a short description about your algorithm

- iii. File/Folder name - Name of the file or folder where your algorithm resides in your S3 bucket. We need this information, so SDC support team can publish this algorithm and make it available to other users
- iv. Readme / Data dictionary file name - This file should provide detailed instructions about your algorithm, how it was created or any relevant information that helps user to understand and use the algorithm. Save this file in your home folder relative to the algorithm file/folder name
- v. Programming Tools/language - Provide the details of programming tools and/or languages that were used to create this algorithm, so users can leverage the same to run your program.

U.S. Department of Transportation
ITS Joint Program Office-HOIT
1200 New Jersey Avenue, SE Washington, DC 20590

Toll-Free “Help Line” 866-367-7487

www.its.dot.gov

FHWA-JPO-18-XXX



U.S. Department of Transportation