

Secure Data Commons

Data Analyst User Guide

www.its.dot.gov/index.htm
Draft Report — April 13, 2020
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 | Added the Export Functionality |
| 3 | REAN Cloud | 3.0 | 10/29/2019 | Added the Manage Workstations chapter and updated the Guide per feedback comments |
| 4 | SDC Team | 4.0 | 02/06/2020 | <ul style="list-style-type: none"> • Added description of data ingestion and curation to Chapter 1 • Added description and link to instructions for setting up and logging in with Login.gov in Chapter 2 • Added note and example command to the “Download User Data from S3 Bucket through Portal” section under Chapter 2 • Added FMI section and link to GitLab page with more sample queries to Chapter 4 Sample Queries for SDC Datasets • Added note about printing documents in Chapter 5 • Added Chapter 6 Setting Up SDC with GitLab Repositories • Added more commands to the “AWS S3 CLI Commands” section under Chapter 7 |
| 5 | SDC Team | 5.0 | 04/13/2020 | <ul style="list-style-type: none"> • Added “Connecting to Waze Data in Redshift Using Python” section to Chapter 3 Accessing and Launching Workstations • Removed Chapter 4 Sample Queries for SDC Datasets due to security concerns; transferred all content to GitLab • Added new Step 5 to Chapter 4 Exporting Datasets from the SDC to describe trusted users as well as further instructions on using the auto-export functionality • Added new Q&A about sample queries to Chapter 7 Frequently Asked Questions |

Table of Contents

| | |
|---|-----|
| Table of Contents | iii |
| Chapter 1. Introduction and Document Overview | 1 |
| Prerequisites | 2 |
| Chapter 2. Initial Setup and Validation..... | 4 |
| Accessing Secure Data Commons Portal | 4 |
| Request Access to Datasets | 9 |
| Upload User Data to S3 Bucket through Portal | 11 |
| Download User Data from S3 Bucket through Portal..... | 12 |
| Chapter 3. Accessing and Launching Workstations..... | 14 |
| Launch Workstations..... | 14 |
| Software Validation..... | 16 |
| Connecting to the Data Warehouse..... | 17 |
| Connecting to Waze Data in Redshift Using SQL Workbench | 17 |
| Connecting to Waze Data in Redshift Using Python | 18 |
| Connecting to the Hadoop Hive Metastore | 20 |
| Update Data Formatting Settings in SQL Workbench..... | 20 |
| Connecting to the SDC Hadoop Data Warehouse Using Python..... | 21 |
| Connecting to Redshift from Linux Environments | 22 |
| Accessing Jupyter Notebook and RStudio Server..... | 23 |
| Manage Workstations..... | 23 |
| Resize Workstation | 24 |
| Schedule/Extend Uptime..... | 26 |
| Stop Workstations | 28 |
| Chapter 4. Exporting Datasets from the SDC | 29 |
| Chapter 5. Setting Up SDC with GitLab Repositories..... | 36 |
| Log In to GitLab..... | 36 |
| Create a Group | 37 |

| | |
|--|----|
| Add Group Members..... | 39 |
| Create a Project | 40 |
| Add Members or Groups to a Project..... | 41 |
| Generate SSH Key | 44 |
| Chapter 6. Technical Documentation and Contact Information | 48 |
| Architecture Diagram..... | 48 |
| Workstation Details..... | 48 |
| Tools and Versions..... | 49 |
| Contact Information | 49 |
| Useful Links | 49 |
| AWS S3 CLI Commands | 49 |
| Chapter 7. Frequently Asked Questions..... | 51 |
| How can I get access to the SDC Datasets? | 51 |
| How will I understand what a particular dataset consists of? | 52 |
| How can I launch a workstation? | 52 |
| Where do I store my data?..... | 53 |
| How can I bring my own datasets/algorithm to my workstation?..... | 53 |
| How can I publish my dataset/algorithm?..... | 53 |
| Where can I find sample queries for my dataset(s)? | 57 |

List of Tables

| | |
|--|----|
| Table 1: GitLab Access Rights by Role | 36 |
| Table 2: Default Workstation Details..... | 48 |
| Table 3: List of Tools Used and Their Versions | 49 |

List of Figures

| | |
|--|----|
| Figure 1: SDC Starting Homepage..... | 4 |
| Figure 2: SDC Portal – Register/Login Page | 5 |
| Figure 3: SDC Portal – Login Page..... | 5 |
| Figure 4: Login.gov – Create New Account | 6 |
| Figure 5: SDC Portal – Synchronize Login.gov Account with SDC | 7 |
| Figure 6: SDC Portal – Landing Page After Login..... | 8 |
| Figure 7: SDC Portal – Datasets Option | 9 |
| Figure 8: SDC Portal – Request Dataset Access..... | 10 |
| Figure 9: SDC Portal – Send Data Access Request | 10 |
| Figure 10: SDC Portal – Upload Files | 11 |
| Figure 11: SDC Portal – Upload Success..... | 11 |
| Figure 12: SDC Portal – Selecting Files for Download..... | 12 |
| Figure 13: SDC Portal – Starting Workstations | 14 |
| Figure 14: SDC Portal – Run Instance Success | 14 |
| Figure 15: SDC Portal – Launch Workstations..... | 15 |
| Figure 16: SDC Workstation – Initialization and Login Screens..... | 15 |
| Figure 17: SDC Workstation – SQL Workbench Icon | 17 |
| Figure 18: SQL Workbench – Create Redshift Connection Profile..... | 17 |
| Figure 19: SQL Workbench – Hive Connection Settings | 20 |
| Figure 20: SQL Workbench – Tools → Options → Data formatting | 21 |
| Figure 21: SDC Portal – Manage Workstation | 23 |
| Figure 22: SDC Portal – Manage Workstation Options..... | 24 |
| Figure 23: SDC Portal – Resize Workstation Option | 24 |
| Figure 24: SDC Portal – Workstation Stopped for Resize Changes | 24 |
| Figure 25: SDC Portal – Resizing Options | 25 |
| Figure 26: SDC Portal – Schedule Resize..... | 26 |
| Figure 27: SDC Portal – Schedule Workstation Uptime Option | 26 |
| Figure 28: SDC Portal – Schedule Uptime | 27 |
| Figure 29: SDC Portal – Tooltip with Existing Scheduled Uptime | 27 |
| Figure 30: SDC Portal – New Tooltip with Extended Uptime Schedule..... | 27 |
| Figure 31: SDC Portal – Stop Workstation..... | 28 |
| Figure 32: SDC Portal – Request Export | 29 |
| Figure 33: SDC Portal – Request Export Form..... | 30 |
| Figure 34: SDC Portal – Approval Form Fields | 31 |
| Figure 35: SDC Portal – Acceptable Use Policy..... | 32 |
| Figure 36: SDC Portal – Auto-Export Status Request Form | 33 |
| Figure 37: AWS CLI – Set Auto-Export Role..... | 34 |
| Figure 38: GitLab Login Prompt - LDAP | 37 |
| Figure 39: GitLab – Create a Group Option | 37 |
| Figure 40: GitLab – New Group Form..... | 38 |
| Figure 41: GitLab – Group Creation Success | 38 |

| | |
|--|----|
| Figure 42: GitLab – Members Left Navigation Menu | 39 |
| Figure 43: GitLab – Add a New Group Member | 40 |
| Figure 44: GitLab – Create a Project | 40 |
| Figure 45: GitLab – New Project Form | 41 |
| Figure 46: GitLab – Project Creation Success | 41 |
| Figure 47: GitLab – Select Project to Add Members..... | 42 |
| Figure 48: GitLab – Settings, Members Left Navigation Menu | 42 |
| Figure 49: GitLab – Add Member to Project | 43 |
| Figure 50: GitLab – Add Group to Project | 44 |
| Figure 51: GitLab – Create Folder for Credentials | 45 |
| Figure 52: GitLab – View Directory Objects..... | 45 |
| Figure 53: SDC Architecture Overview..... | 48 |
| Figure 54: SDC Portal – Datasets Tab | 51 |
| Figure 55: SDC Portal – Dataset Access Request..... | 52 |
| Figure 56: SDC Portal – Workstations Tab | 52 |
| Figure 57: SDC Workstation – Login | 53 |
| Figure 58: SDC Portal – Publish Button | 54 |
| Figure 59: SDC Portal – Publish Dataset Request Form | 55 |
| Figure 60: SDC Portal – Publish Algorithm Request Form..... | 56 |

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, Libre Office, 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.

During a project's onboarding phase, Data Providers work with the SDC support team to describe their project's data (for example, the type of data, frequency of new data, data file formats, etc., every project's data is unique). Then, Data Providers upload data files to designated "S3 Ingestion Buckets" (a secure, scalable object storage service provided by the SDC platform through Amazon Web Services). We call this the "Data Lake."

As new data arrives to the SDC, policies and procedures established by the Data Provider then govern who, when, and how Data Analysts can access the data. It is common that once new data arrives in the SDC, automated processes “ingest” and “curate” the data, making the data available in other forms. For example, some data may be loaded into our Data Warehouse tools (Redshift or Hadoop databases), whereas other data may be transformed from its source format into other easier-to-use formats, or filtered through a process to identify corrupt, invalid, or duplicate data. Exactly which automated processes a project’s data undergoes is determined during the project onboarding processes.

Once data has been ingested and curated, it is then available to Data Analysts through the tools listed above (note that we are always adding new Data Analyst tools based on request). A typical Data Analyst workflow may be:

- Use a tool to develop SQL queries to see a subset of the larger data set.
- Compare the data subset against models to draw unique insights (for example, develop programs utilizing the analytic capabilities of R or Python).
- Use powerful tools that present data and insights in graphical format (some use the power of Python and GeoPandas, others have developed graphical outputs in R Studio, whereas others use Libre Office as an open source alternative to Microsoft Excel). The SDC support team has worked with yet other Data Analysts to install proprietary licensed software to enhance their analytical capabilities.

Finally, there are capabilities by which Data Analysts can export their work out of the SDC, subject to the data use agreements and approval of the Data Providers.

This document provides guidance for the **Data Analyst** role. A similar guide will be prepared for the Data Providers. This 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
- Setting Up GitLab with 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>.

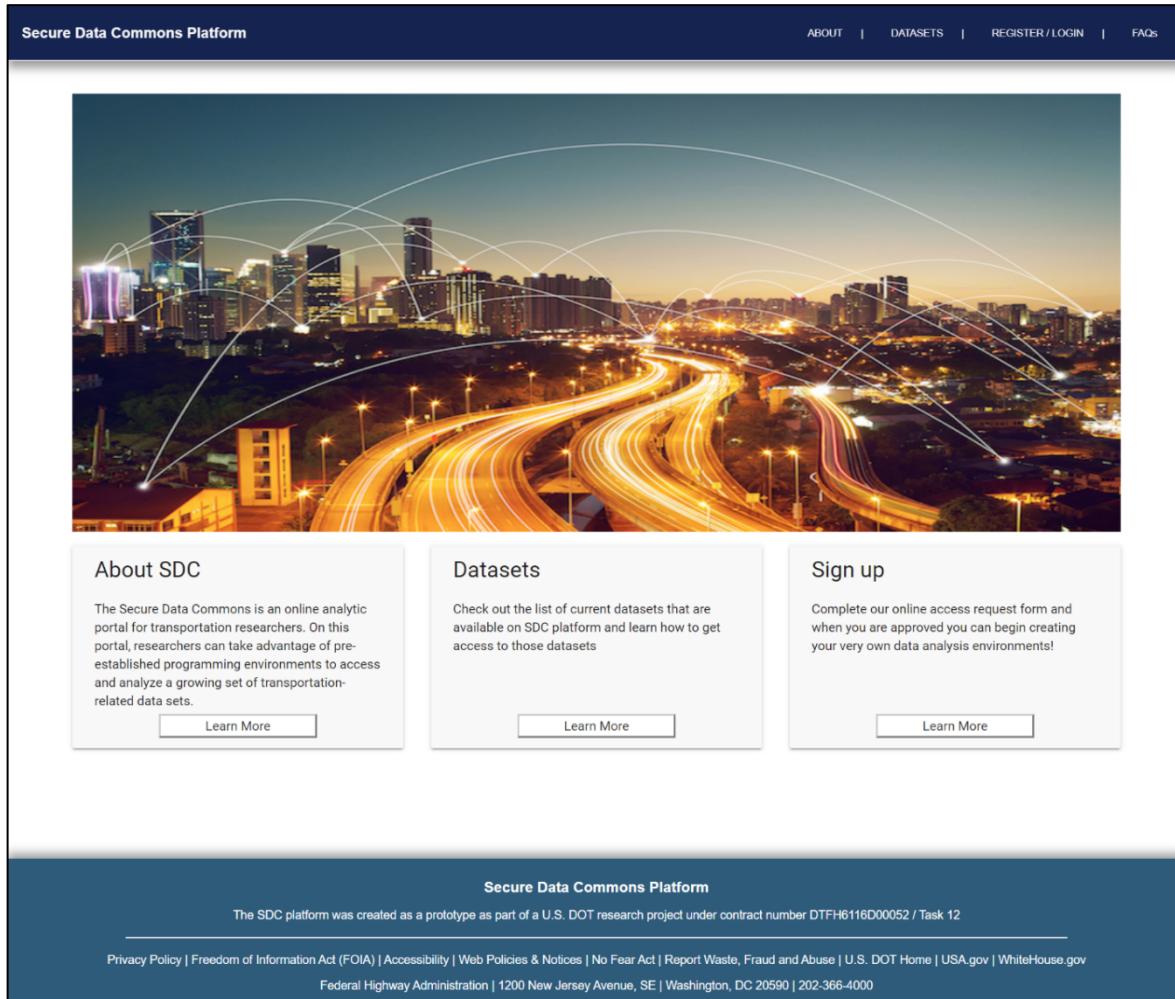


Figure 1: SDC Starting Homepage

Chapter 2. Initial Setup and Validation

Select Register/Login from the top menu to display the Access Request Form and Privacy Policy links, as well as the email address sign-in prompt. The Access Request Form takes you to a PDF form in which you can provide the support team with contact info, a rationale for your access request, your preferred workstation type, and sign-off for the SDC data use agreements.

If you have a DOT email address (@dot.gov), enter it into Email Address and then select Sign In to access the portal. If you do not have a DOT email address, continue to Page 6.

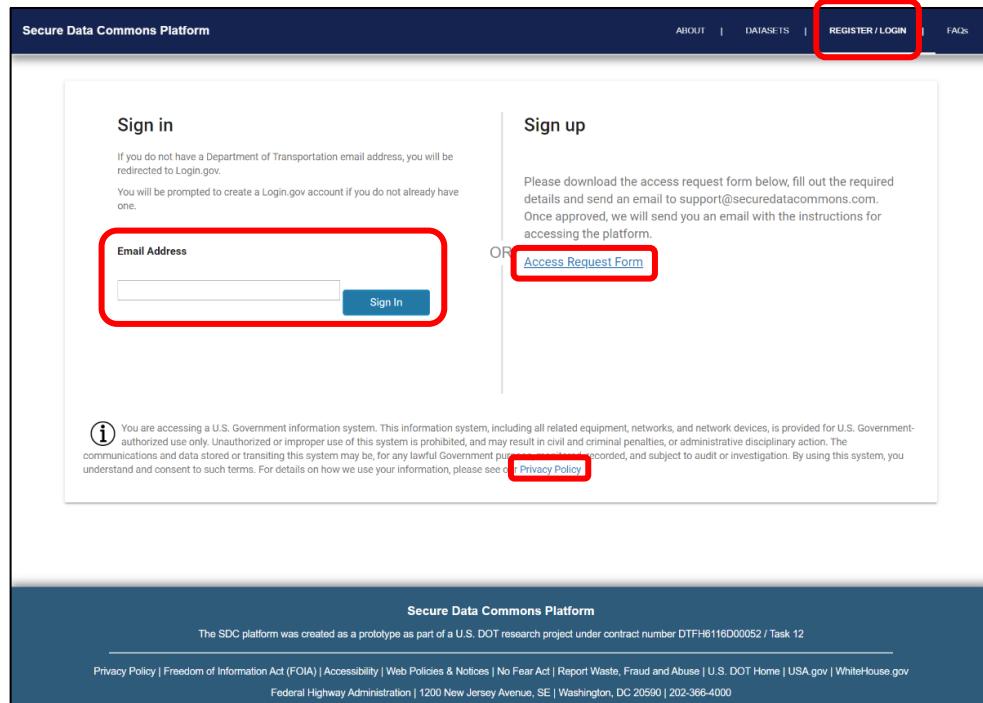


Figure 2: SDC Portal – Register/Login Page

After signing in, you will be redirected to the Secure Data Commons platform login page:

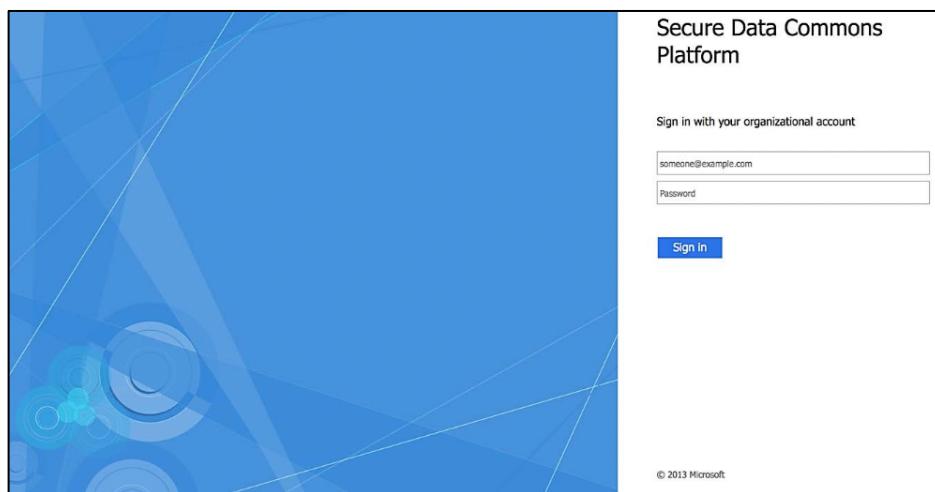


Figure 3: SDC Portal – Login Page

If you do not have a DOT email address, you will need to create a new Login.gov account or log in with your existing one so that it can be linked with your SDC credentials. Linking a Login.gov account with your SDC credentials provides extra security by adding protection you can configure as two layers of authentication. Authentication methods will consist of phone text or call; an app on your phone, tablet, or computer; a security key; a government employee ID; and pre-generated backup codes.

Enter the email address you would like to use as your Login.gov account into the portal's Email Address sign in (Figure 2 above) and then select Sign In to be redirected to the Login.gov website, where you can create a new account. If you already have an existing Login.gov account, enter your credentials and then select Sign in.

If you do not have an existing Login.gov account, select "Create an account" on the page you are redirected to:

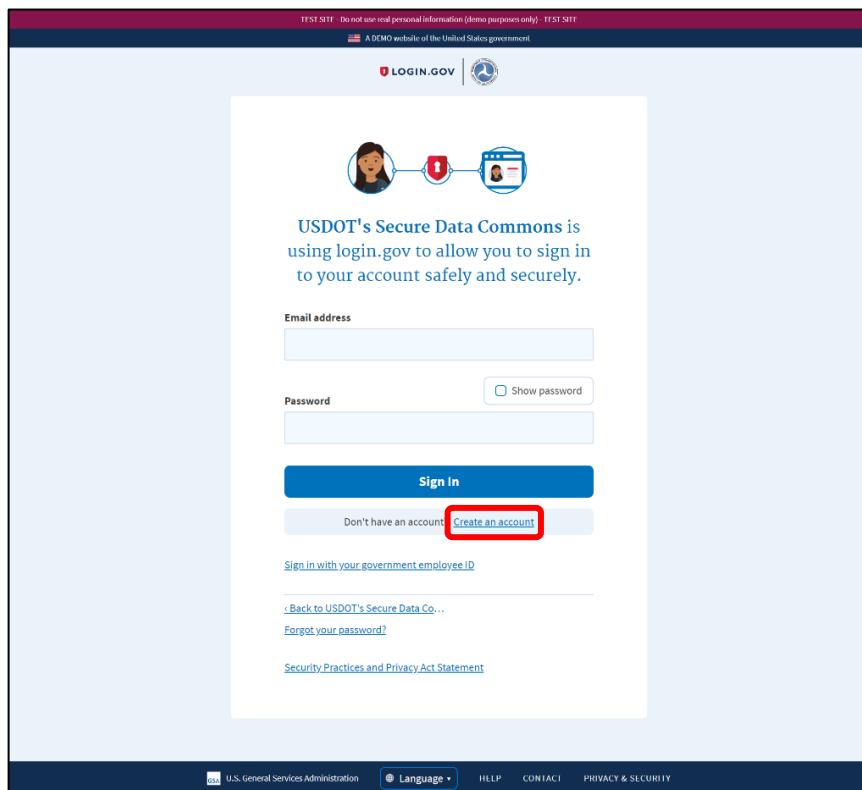


Figure 4: Login.gov – Create New Account

For further instructions on setting up your new Login.gov account and configuring its secure authentication methods, refer to <https://www.login.gov/help/creating-an-account/how-to-create-an-account/>.

Chapter 2. Initial Setup and Validation

After you have finished setting up all authentication methods or signed in with an existing Login.gov account, you will be redirected to a one-time sign-in form on the portal. Enter your SDC ADFS credentials (provided in your welcome email) and then select Sign in.

The screenshot shows the 'One Time Sign In' page of the Secure Data Commons Platform. At the top, there is a dark blue header bar with the text 'Secure Data Commons Platform' on the left and navigation links for 'HOME', 'DATASETS', 'WORKSTATIONS', 'FAQ', and 'LOGOUT' on the right. Below the header, the main content area has a white background. It features a title 'One Time Sign In' at the top center. Below the title are two input fields: 'Username' containing '@securedatacommons.com' and 'Password' containing a series of dots. There is also a 'Show Password' link. A large blue 'Sign in' button is centered below the password field. To the right of the input fields, there is a note: 'Please sign in using your internal Secure Data Commons credentials. Accounts only need to be linked once.' Below this note, another note says: 'If you do not remember your Secure Data Commons credentials, please contact the Support Team.' At the bottom of the page, there is a footer bar with the text 'Secure Data Commons Platform' and a small note: 'The SDC platform was created as a prototype as part of a U.S. DOT research project under contract number DTFH6116D00052 / Task 12'. The footer also includes links for 'Privacy Policy', 'Freedom of Information Act (FOIA)', 'Accessibility', 'Web Policies & Notices', 'No Fear Act', 'Report Waste, Fraud and Abuse', 'U.S. DOT Home', 'USA.gov', and 'WhiteHouse.gov'. It also lists the 'Federal Highway Administration | 1200 New Jersey Avenue, SE | Washington, DC 20590 | 202-368-4000'.

Figure 5: SDC Portal – Synchronize Login.gov Account with SDC

If you are accessing the portal for the first time, you will be prompted to change your password after entering the credentials provided in the welcome email.

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

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

Chapter 2. Initial Setup and Validation

The screenshot shows the homepage of the Secure Data Commons Platform (SDC). At the top, there is a dark blue header bar with the text "Secure Data Commons Platform" on the left and navigation links for "HOME", "DATASETS", "WORKSTATIONS", "FAQ", and "LOGOUT" on the right. Below the header, there is a large white content area. At the top of this area, there is a welcome message: "Welcome to the Secure Data Commons Platform (SDC), a sandbox environment in AWS managed by U.S. Department of transportation. The SDC platform provides a collaborative environment for traffic engineers, researchers, data scientists and anyone who is interested in carrying out research and analysis on different datasets related to traffic, weather, crashes etc." Below this message, there is a section titled "Ready to Get Started? Follow the steps below or directly jump to any of the sections that you wish to access. SDC Users will have the following access in the platform:" followed by four numbered steps:

- 1 Curated and published datasets**
Check out the list of currently available datasets (curated, published), data dictionaries owned by you and others. [Learn more](#) about what curated and published datasets mean
- 2 Access to workstations with programming tools**
Access the workstations assigned to you, if you have requested as part of the registration process
- 3 Bring your own datasets / algorithms**
[Learn more](#) about how you can bring your algorithms and code to the SDC platform
- 4 Publish your datasets / algorithms**
[Learn more](#) about how you can publish your research analysis and algorithms with other researchers

At the bottom of the content area, there is a dark blue footer bar with the text "Secure Data Commons Platform" and a note: "The SDC platform was created as a prototype as part of a U.S. DOT research project under contract number DTFH6116D00052 / Task 12". Below this, there are links to "Privacy Policy", "Freedom of Information Act (FOIA)", "Accessibility", "Web Policies & Notices", "No Fear Act", "Report Waste, Fraud and Abuse", "U.S. DOT Home", "USA.gov", and "WhiteHouse.gov". There is also a line of text: "Federal Highway Administration | 1200 New Jersey Avenue, SE | Washington, DC 20590 | 202-366-4000".

Figure 6: SDC Portal – 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.

The screenshot shows the SDC Portal homepage. At the top, there is a dark blue header with the text "Secure Data Commons Platform". Below the header, there is a navigation bar with links: "HOME", "DATASETS" (which is highlighted with a red box), "WORKSTATIONS", "FAQ", and "LOGOUT". The main content area has a white background. It starts with a welcome message: "Welcome to the Secure Data Commons Platform (SDC), a sandbox environment in AWS managed by U.S. Department of transportation. The SDC platform provides a collaborative environment for traffic engineers, researchers, data scientists and anyone who is interested in carrying out research and analysis on different datasets related to traffic, weather, crashes etc." Below this, there is a section titled "Ready to Get Started? Follow the steps below or directly jump to any of the sections that you wish to access. SDC Users will have the following access in the platform:" followed by four numbered steps:

- 1 Curated and published datasets
Check out the list of currently available datasets (curated, published), data dictionaries owned by you and others. [Learn more about what curated and published datasets mean](#)
- 2 Access to workstations with programming tools
Access the [workstations assigned to you](#), if you have requested as part of the registration process
- 3 Bring your own datasets / algorithms
[Learn more about how you can bring your algorithms and code to the SDC platform](#)
- 4 Publish your datasets / algorithms
[Learn more about how you can publish your research analysis and algorithms with other researchers](#)

At the bottom of the page, there is a dark blue footer with the text "Secure Data Commons Platform" and "The SDC platform was created as a prototype as part of a U.S. DOT research project under contract number DTFH6116D00052 / Task 12". Below this, there are links to "Privacy Policy", "Freedom of Information Act (FOIA)", "Accessibility", "Web Policies & Notices", "No Fear Act", "Report Waste, Fraud and Abuse", "U.S. DOT Home", "USA.gov", and "WhiteHouse.gov". The footer also includes the address "Federal Highway Administration | 1200 New Jersey Avenue, SE | Washington, DC 20590 | 202-366-4000".

Figure 7: SDC Portal – Datasets Option

Chapter 2. Initial Setup and Validation

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

The screenshot shows the 'SDC Datasets' page. It displays a table of datasets with columns: Name, Category, Description, Geographic Scope, Start / End for Data Availability, Owner, and Request Access. The 'Request Access' column contains three blue 'Request' buttons, which are all highlighted with a red box. The datasets listed are WAZE, FSD, and CVP.

| 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 | <button>Request</button> |
| FSD | Raw | Contains raw FTA sample data | All states in US | March 2019 | SDC platform | <button>Request</button> |
| CVP | Curated | Contains CVP evaluation datasets | All states in US | March 2017 to Present | SDC platform | <button>Request</button> |

Figure 8: SDC Portal – Request Dataset Access

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

The screenshot shows the 'SDC Data Access Request Form'. It includes sections for Basis for Access (radio buttons for 'Yes' and 'No'), Geographic Extent of Access (text input field), and a checkbox for accepting terms and conditions. At the bottom, there are 'SEND REQUEST' and 'CANCEL' buttons. The 'SEND REQUEST' button is highlighted with a red box.

Figure 9: SDC Portal – 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 who want to share data with other users from their project team 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.)

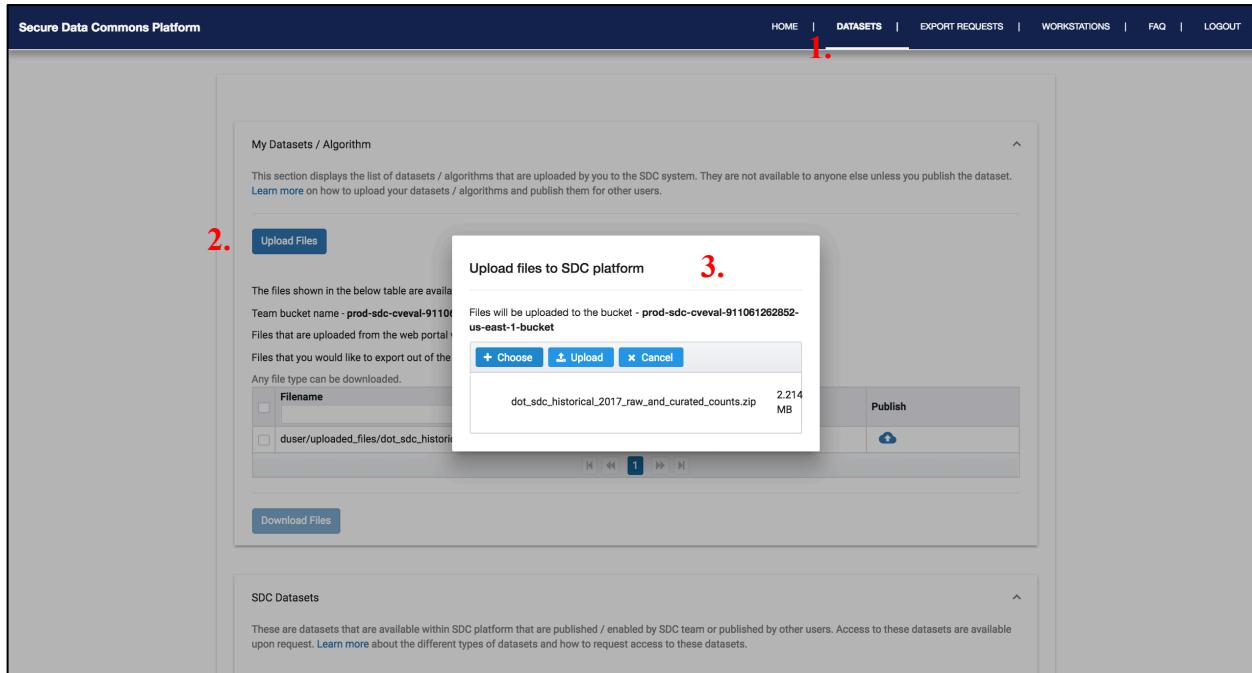


Figure 10: SDC Portal – Upload Files

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

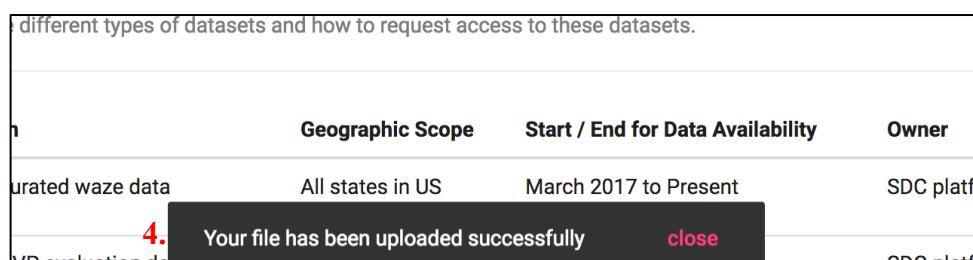


Figure 11: SDC Portal – Upload Success

5. Files that are uploaded from the web portal will be saved in the folder – `username/uploaded_files`
6. To make these files available to members of the project team, users then need to move files from `username/uploaded_files` to the project team S3 bucket using command-line tools (see AWS S3 CLI Commands for an overview).

7. Users would be able to access only the files that are under the **uploaded_files** and **export_requests** folders.

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](#).

The screenshot shows the SDC Portal interface. At the top, there's a navigation bar with links for HOME, DATASETS, EXPORT REQUESTS, WORKSTATIONS, FAQ, and LOGOUT. Below the navigation bar, the title 'My Datasets / Algorithm' is displayed. A callout '1.' points to the introductory text: 'This section displays the list of datasets / algorithms that are uploaded by you to the SDC system. They are not available to anyone else unless you publish the dataset. Learn more on how to upload your datasets / algorithms and publish them for other users.' Another callout '2.' points to a table listing files. The table has columns for 'Filename', 'Export', and 'Publish'. A file named 'duser/uploaded_files/dot_sdc_historical_2017_raw_and_curated_counts.zip' is selected, indicated by a checked checkbox in the 'Filename' column. Callout '3.' points to the 'Download Files' button at the bottom of the table.

Figure 12: SDC Portal – 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.
- Files must be copied to and from S3 buckets using the SDC workstation. **NOTE:** Files stored in S3 buckets are not visible through Windows Explorer but can be copied to their SDC workstation.
 - o For example: To copy the file shown above in Figure 12 that is hosted in an S3 bucket to your SDC workstation, use the command:

```
aws s3 cp s3://prod-sdc-cveval-911061262852-us-east-1-
bucket/duser/uploaded_files/dot_sdc_historical_2017_ra
w_and_curated_counts.zip
dot_sdc_historical_2017_raw_and_curated_counts.zip
```

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.

| # | 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 | CPU: 2 Memory(GiB): 4 | Start Launch Manage |
| 2 | Programming Environment #2 | Jupyter Notebook, R Studio Server | CPU: 2 Memory(GiB): 4 | Start Launch Manage |

Figure 13: SDC Portal – 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.

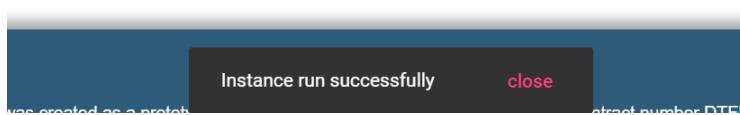


Figure 14: SDC Portal – Run Instance Success

4. Now click Launch for the workstation.

The screenshot shows the 'My Workstations' section of the SDC Portal. It lists two workstations:

| # | 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 | <button>Start</button> <button>Launch</button> <button>Manage</button> |
| 2 | Programming Environment #2 | Jupyter Notebook, R Studio Server | CPUs: 2 Memory(GiB): 4 | <button>Start</button> <button>Launch</button> <button>Manage</button> |

A red box highlights the 'Launch' button for the second workstation.

Figure 15: SDC Portal – Launch Workstations

5. 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.

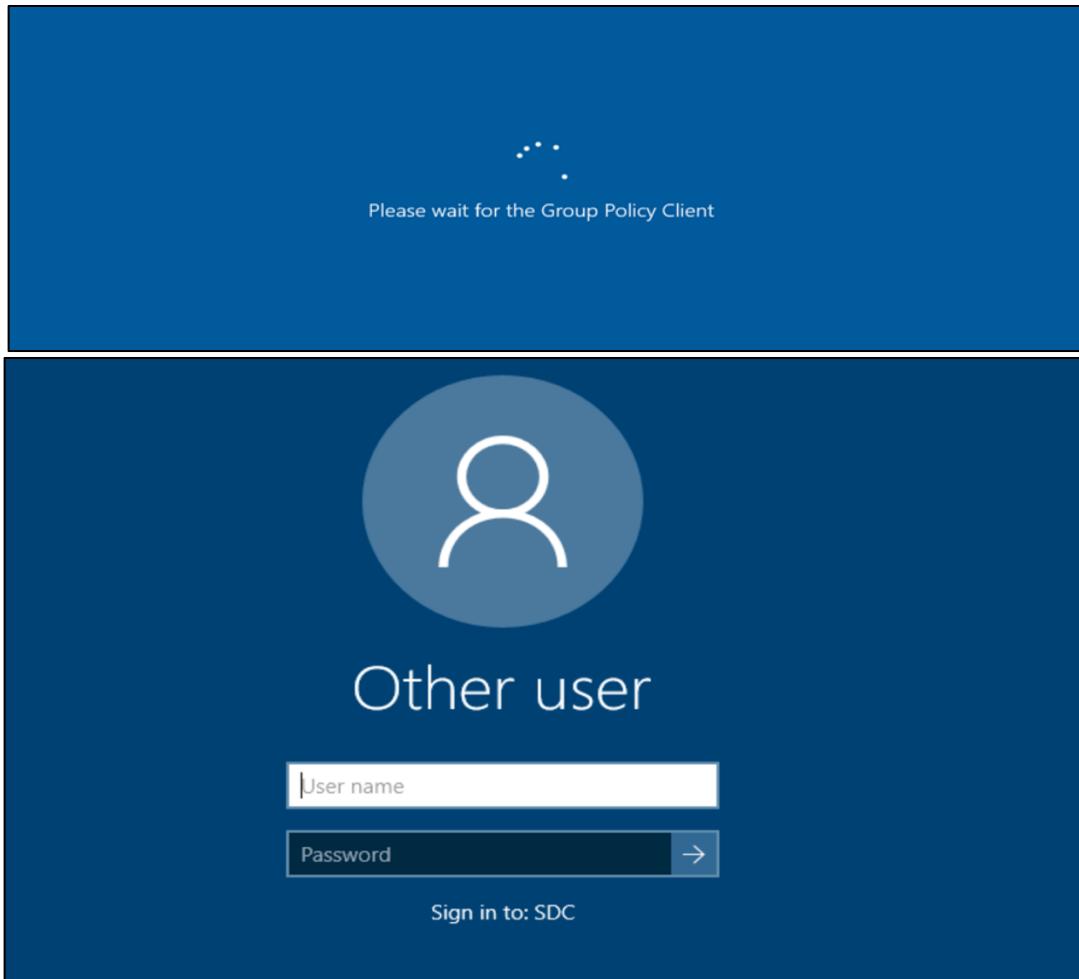


Figure 16: SDC Workstation – Initialization and Login Screens

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

Connecting to the Data Warehouse

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

Connecting to Waze Data in Redshift Using SQL Workbench

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

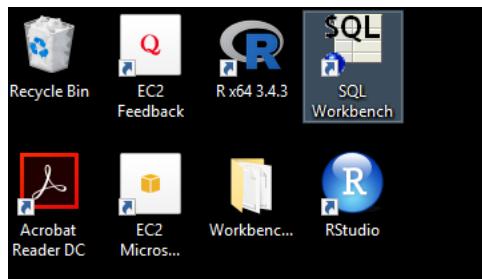


Figure 17: SDC Workstation – SQL Workbench Icon

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.

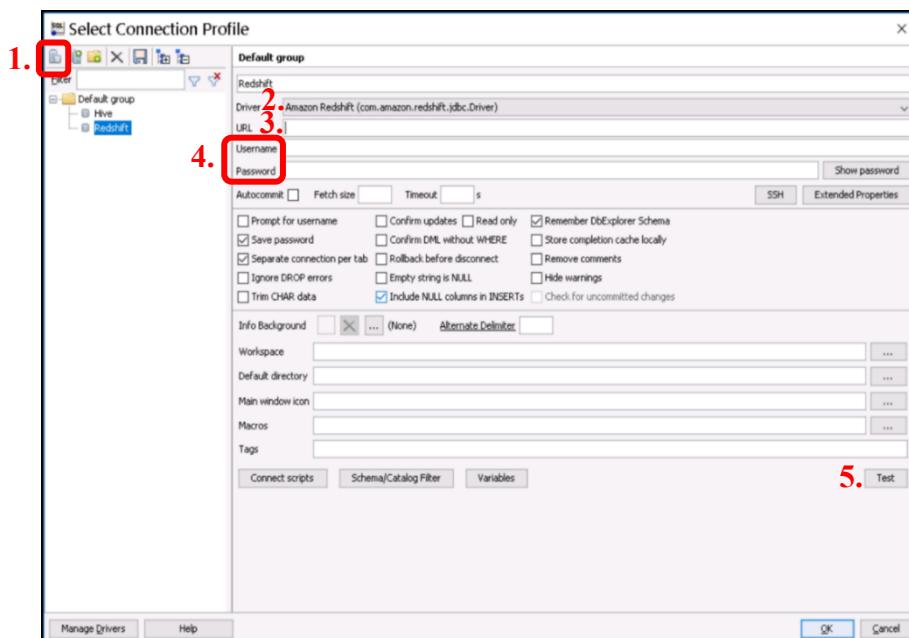


Figure 18: SQL Workbench – Create Redshift Connection Profile

Connecting to Waze Data in Redshift Using Python

NOTE: When you are granted access to Waze data, the SDC support team creates a new Redshift user for you, assigns it with a Redshift password, and emails you with information on the Redshift host you will connect to. This email message provides the redshiftHost, userName, and userPassword values shown below. Your Redshift credentials are only used for connecting to Redshift and NOT for accessing the portal, which uses your separate SDC credentials.

Important: The default version of Python installed on the SDC Windows Workstations is v2.7.4. There are two required Python modules that must be installed prior to attempting to connect to Redshift with Python using the example code below. To install these modules, open a Windows Command Prompt, and enter the following two commands:

```
C:\Users\username> pip install psycopg2  
C:\Users\username> pip install numpy
```

The above "pip install ..." command(s) only need to be run ONCE on the SDC Windows Workstation. Once the Python modules are installed, they remain available, even across reboots of the workstation.

To test Python connectivity to Redshift, open the IDLE python editor and execute the following:

```
from __future__ import print_function  
  
import psycopg2  
import numpy  
  
dbName = 'dot_sdc_redshift_db'  
redshiftHost = '[host address]'  
redshiftPort = 5439  
  
userName = '[username]'  
userPassword = '[password]'  
  
# query = 'select * from dw_waze.alert limit 10;'  
query = "select * from dw_waze.alert where  
alert_type='ACCIDENT' and city = 'Severance, CO'"  
  
conn = psycopg2.connect(  
    dbname=dbName,  
    host=redshiftHost,  
    port=redshiftPort,
```

```
user=userName,  
password=userPassword)  
  
cursor = conn.cursor()  
  
cursor.execute(query)  
result = cursor.fetchall()  
  
result = numpy.array(result)  
  
# print(result)  
for r in result:  
    print (r[1], r[8], r[18], r[19], r[22], sep='\t')
```

For further information and examples, refer to the internal SDC GitLab collaboration site.

Connecting to the Hadoop 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.
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.

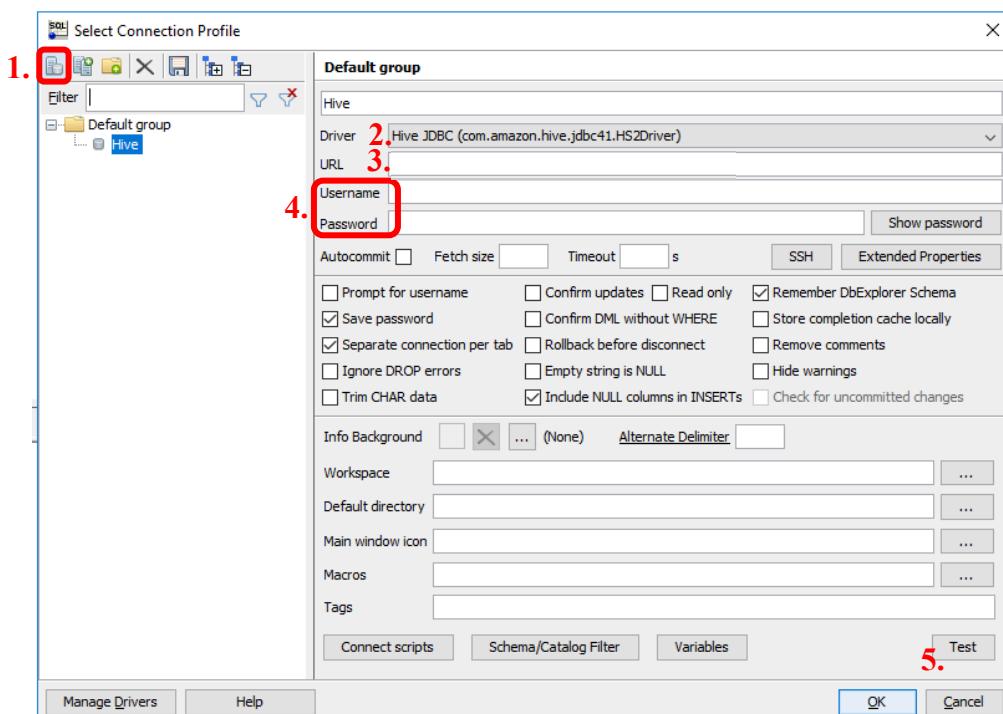


Figure 19: SQL Workbench – Hive Connection Settings

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.



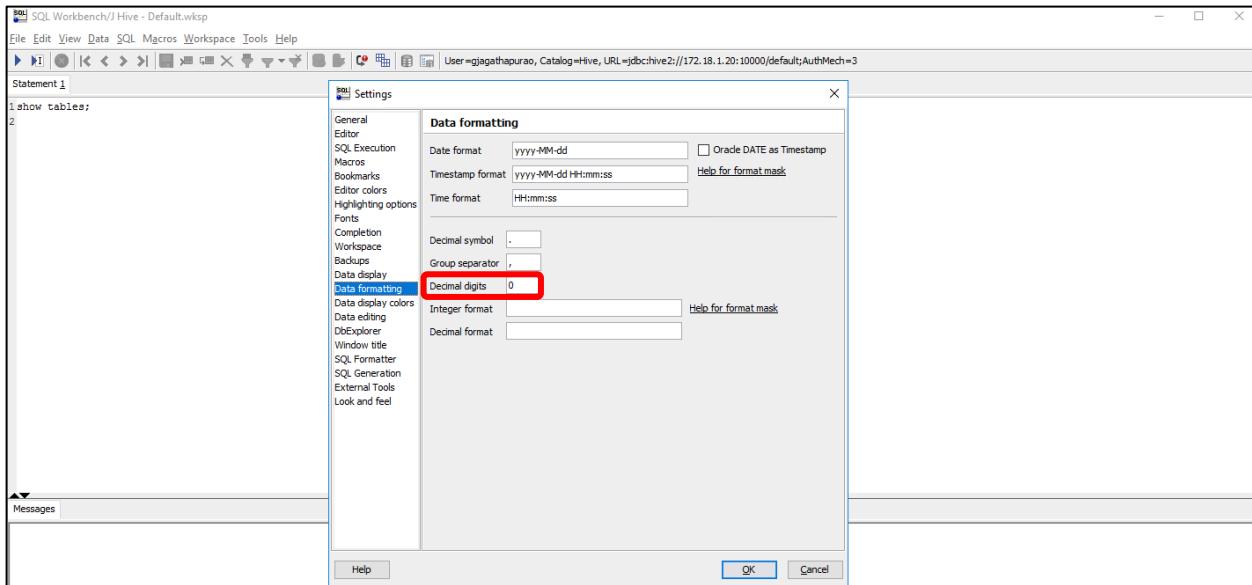


Figure 20: SQL Workbench – Tools → Options → Data formatting

Connecting to the SDC Hadoop Data Warehouse Using Python

Important: The default version of Python installed on the SDC Windows Workstations is v2.7.4. There are two required Python modules that must be installed prior to attempting to connect to Hadoop/Hive with Python using the example code below. To install these modules, open a Windows Command Prompt, and enter the following two commands:

```
C:\Users\username> pip install impyla  
C:\Users\username> pip install numpy
```

The above "pip install ..." command(s) only need to be run ONCE on the SDC Windows Workstation. Once the Python modules are installed, they remain available, even across reboots of the workstation.

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

```
from __future__ import print_function  
  
from impala.dbapi import connect  
import numpy  
  
conn = connect(  
    host='[host address]',  
    port=10000,  
    auth_mechanism='PLAIN',
```

```
user='[your_username]' ,password='[your_password]')

cursor = conn.cursor()

cursor.execute('SHOW TABLES')

result = cursor.fetchall()
result = numpy.array(result)

# print(result)
for r in result:
    print (r)
```

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

Connecting to Redshift from Linux Environments

Credentials to access the Waze Redshift database are communicated from the SDC Administrator (support@securedatacommons.com).

- 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)"
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 | CPU: 2 Memory(GiB): 4 | <button>Start</button> <button>Launch</button> Manage |
| 2 | Programming Environment #2 | Jupyter Notebook, R Studio Server | CPU: 2 Memory(GiB): 4 | <button>Start</button> <button>Launch</button> Manage |

Figure 21: SDC Portal – Manage Workstation

A dialogue window appears with two checkbox options:

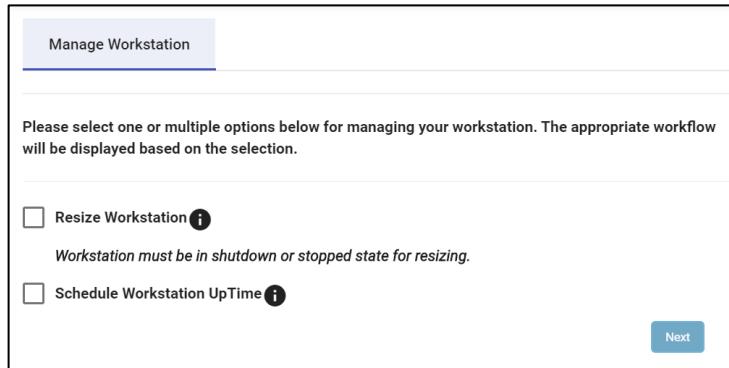


Figure 22: SDC Portal – Manage Workstation Options

Selecting each option renders the appropriate tabs in the dialogue window. The ⓘ 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.

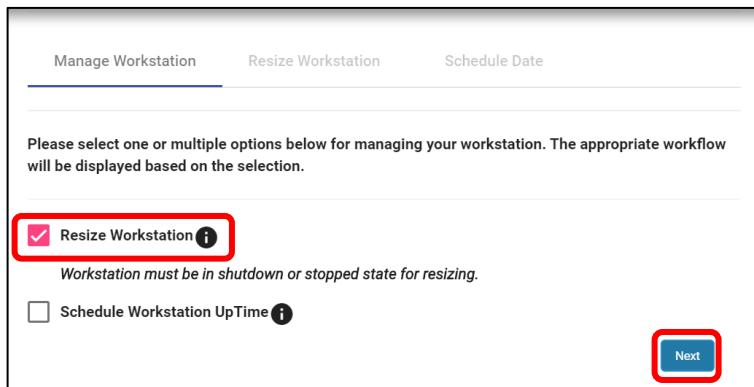


Figure 23: SDC Portal – Resize Workstation Option

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

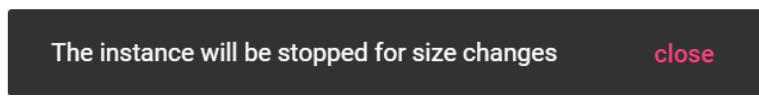


Figure 24: SDC Portal – Workstation Stopped for Resize Changes

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 selected and highlighted in blue), and 'Schedule Date'. Below the tabs, a message states: 'Your current workstation configuration is 2 CPU and 4 GB RAM'. It instructs users to 'Please choose the appropriate instance type by selecting the desired CPU and RAM' and provides a link for 'Pricing related details click here.' A note below 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.' Under the heading 'Select desired CPU / Memory', there are dropdown menus for 'CPU' (set to 4) and 'RAM' (set to 8 GB). To the right is a blue 'Submit' button. Below this, step 4. is highlighted with a red box around the checked checkbox labeled 'Please start my workstation after resizing to the new configuration'. Further down, step 6. is highlighted with a red box around the 'c5.xlarge' entry in the 'Recommended list' table. The table has columns: Instance Name, CPU, Memory, and COST. The entry for 'c5.xlarge' is: Instance Name c5.xlarge, CPU 4, Memory 8 GB, COST \$0.35 per hour. At the bottom, there are 'Cancel' and 'Next' buttons, with 'Next' being blue and bold.

| Instance Name | CPU | Memory | COST |
|---------------|-----|--------|-----------------|
| 6. c5.xlarge | 4 | 8 GB | \$0.35 per hour |

Figure 25: SDC Portal – 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 the 'Select schedule' page with three tabs: 'Manage Workstation', 'Resize Workstation', and 'Schedule Date'. The 'Schedule Date' tab is selected. A message at the top says: '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 is a section titled 'Workspace resize schedule' with 'From date' set to '10/5/2019' and 'To date' set to '10/11/2019'. At the bottom are 'Cancel' and 'Submit' buttons, with 'Submit' being highlighted by a red box.

Figure 26: SDC Portal – 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 the 'Schedule Workstation UpTime' page. It has 'Manage Workstation' and 'Schedule Workstation UpTime' tabs. A message says: 'Please select one or multiple options below for managing your workstation. The appropriate workflow will be displayed based on the selection.' There are two checkboxes: 'Resize Workstation' (unchecked) and 'Schedule Workstation UpTime' (checked). A note below says: 'Workstation must be in shutdown or stopped state for resizing.' At the bottom is a 'Next' button, which is highlighted by a red box.

Figure 27: SDC Portal – Schedule Workstation Uptime Option

Chapter 3. Accessing and Launching Workstations

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.

Manage Workstation Schedule Workstation UpTime

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.

Schedule Workstation UpTime

From date: 10/4/2019 To date: 10/8/2019

Cancel Submit

Figure 28: SDC Portal – 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.

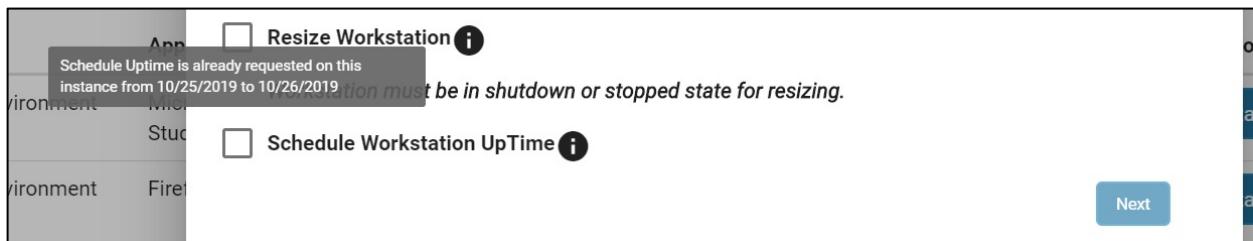


Figure 29: SDC Portal – Tooltip with Existing Scheduled Uptime

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.

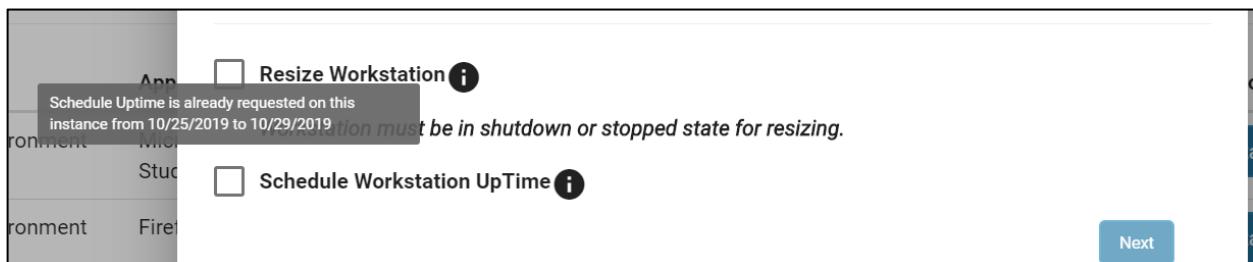
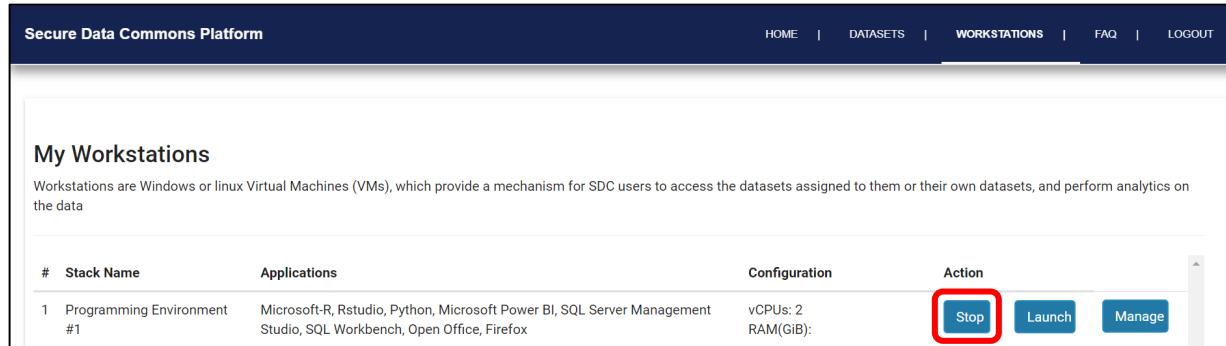


Figure 30: SDC Portal – New Tooltip with Extended Uptime Schedule

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.



The screenshot shows the 'My Workstations' section of the SDC Portal. At the top, there is a navigation bar with links for HOME, DATASETS, WORKSTATIONS (which is the active tab), FAQ, and LOGOUT. Below the navigation bar, the title 'My Workstations' is displayed, followed by a brief description: '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'. A table lists the workstations:

| # | 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 | vCPUs: 2 RAM(GiB): | Stop (button highlighted with a red box) Launch Manage |

Figure 31: SDC Portal – Stop Workstation

Chapter 4. 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. Please note that if Data Analysts want to print out a hard copy of a document, they will need to export it from the SDC workstation to their local machine.

The screenshot shows a table titled "Request Export" with the following data:

| Filename | Export | Publish |
|--|-------------------------------------|-------------------------------------|
| export_requests/test3.zip | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| export_requests/s3cmd2.zip | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| export_requests/t.zip | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| export_requests/demo9.zip | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| sbapat/uploaded_files/architecture.png | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| export_requests/s3cmd1.zip | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| export_requests/wazeAnalysis.zip | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| export_requests/dynamodbanalysis.zip | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |

Figure 32: SDC Portal – Request Export

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 Project Approval Form Trusted 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

CANCEL **NEXT**

Figure 33: SDC Portal – Request Export Form

3. 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> | | |
| Name or short description of your derived dataset * | | |
| waze derived dataset | | |
| Anchor dataset of interest or data provider * | | |
| WAZE | | |
| Specific sub-datasets or data types used * | | |
| JAM | | |
| Additional datasources * | | |
| datasources | | |
| High level description of derived dataset * | | |
| high level desc | | |
| Detailed description of the derived dataset * | | |
| detailed desc | | |
| Tags | | |
| tags | | |
| Justification of Export * | | |
| Justification | | |

Figure 34: SDC Portal – 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 |
|--|---------------|----------------|
| <p>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.</p> <p>This status is acquired per Project + Data Provider + Sub-Dataset/Data Type.</p> <p>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.</p> <p>Do you wish to request Trusted Status from the Data Provider?</p> <p><input type="radio"/> Yes <input checked="" type="radio"/> No</p> <p>Acceptable Use Policy</p> <p>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.</p> <p>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 use of the WAZE CV Pilot data</p> <p><input type="radio"/> Accept <input type="radio"/> Decline</p> | | |

Figure 35: SDC Portal – Acceptable Use Policy

5. If the user is a trusted user, he/she will be prompted with the option for requesting the auto-export approved status from the Data Provider. This will allow the analyst to be able to export data automatically, as opposed to filling out the request form and then going through the approval process before exporting. The analyst enters the derived dataset name of the data type they want to export into ‘Specify the derived dataset name’ and then the justification for requesting permissions into ‘Reason for requesting permission.’ The form will still be submitted if the user declines the option to submit a request for auto-export approved status by selecting ‘No’ under ‘Do you wish to request Auto-Export Status from the Data Provider?’ This form will be unavailable if the user does not have trusted status yet.

Auto-Export is a feature that allows data analysts to export their derived data or analysis products automatically. It avoids the need to manually request exports many times for a similar type of export (as a trusted user). Data analysts must request access to this feature for specific types of exports from data providers. At any point after approval, the data provider may revoke this permission. The trusted status associated with access to the auto-export functionality is given per Project, Data Provider, and Per Sub-Dataset/Datatype.

Note - You currently do not have Auto-Export Status from this Data Provider. We will notify you when the Data Provider has made a decision.

Do you wish to request Auto-Export Status from the Data Provider?

Yes No

Specify the derived dataset name *

Monthly_Reports

Reason for requesting permission *

Daily automated machine

Submit

Figure 36: SDC Portal – Auto-Export Status Request Form

6. Upon successful submission, the request will be sent to appropriate Data Providers. Data Providers will be responsible for accepting or rejecting the export requests.

7. Once Data Providers approve the request, Data Analysts will be able to download the dataset out of SDC through portal.
8. Log into the SDC Portal to start and launch the Guacamole workstation. Refer to the instructions in [Chapter 3 Accessing and Launching Workstations](#) if needed.
9. Log into the Guacamole workstation.
10. Perform work and save any derived data that needs to be exported to the local machine.
For example, the ‘Documents’ folder would be a good place.
11. Open PowerShell and navigate to the directory that the derived data is located at. **NOTE:** If you are unable to find PowerShell, right-click on the Windows menu, select Search, and enter “powershell.”
 - a. Use the command `cd` to change the directory, for example, `cd Documents`
12. Copy the file into your team bucket `auto_export` directory’s correct sub-directory that you were approved to auto-export. For example:
 - a. If you are a member of your team approved to auto-export ‘Monthly_Reports,’ then you would copy the file to your team bucket’s ‘`auto_export/Monthly_Reports/`’ directory
 - b. The following command would be run in the PowerShell window:
`aws s3 cp deriveddatatoexport.file s3://team-bucket/auto_export/Monthly_Reports/`
 - i. See the [AWS S3 CLI Commands](#) section for more information on the `aws s3 cp` command.
13. After the copy has completed, the SDC will run processes to auto-export the derived data to your team bucket’s auto-export S3 bucket. This process should only take a few seconds.
14. To see if your derived data has been exported, sign into the AWS CLI on your local machine using the SDC-provided authentication scripts. Make sure to set the role as your auto-export role for the specific derived datatype you are exporting. You will ONLY be able to export the data of the role you have assigned. You must change your role to export other types of derived data if you are approved for multiple datatypes.
 - a. If you need access to these authentication scripts, please contact the [SDC support desk](#).

```
PS C:\Users\...> py -3 .\samlapi_formauth_adfs3_windows.py
Username: [REDACTED]@securedatacommons.com
Password:

Please choose the role you would like to assume:
[ 0 ]: arn:aws:iam::911061262852:role/DOT-WYDOTSpeed_ReportsAutoExport
[ 1 ]: arn:aws:iam::911061262852:role/DOT-WYDOTBSMAutoExport
[ 2 ]: arn:aws:iam::911061262852:role/DOT-WYDOTUsers
[ 3 ]: arn:aws:iam::911061262852:role/DOT-Developers
Selection:
```

Figure 37: AWS CLI – Set Auto-Export Role

15. Using the AWS CLI, you may list the content of your derived datatype's directory from the team auto-export bucket and download any files located in that directory to your local machine.
- a. Your permissions are very limited, and the commands must be very precise for them to work.
 - b. List files: `aws s3 ls s3://team-auto-export-bucket/datatype/`
 - c. Download file: `aws s3 cp s3://team-auto-export-bucket/datatype/deriveddata.file localmachinelocation`

Chapter 5. Setting Up SDC with GitLab Repositories

The following sections provide instructions for SDC users on setting up GitLab repositories for committing and sharing their code and any other documents.

GitLab provides an organizational structure of users and repositories managed respectively via “groups” and “projects” to promote collaboration among teams of SDC users. Groups contain members with different access levels to their project repositories based on the assigned roles:

Table 1: GitLab Access Rights by Role

| Role | Read-Only Access | Read/Write Access |
|------------|------------------|-------------------|
| Guest | ✓ | |
| Reporter | ✓ | |
| Developer | | ✓ |
| Maintainer | | ✓ |

GitLab Features:

- Share code among a [subset of your team, your entire team, or any number of other teams on the SDC](#) and provide comments to improve each other’s code
- Simplified collaboration with shared read/write access
- Support for multiple repositories to be kept in a logical structure that separates them from each other and groups by organization.
- Version control and source code management including:
 - File locking to help prevent conflicts
 - Users can work from their local copy of the code
 - Users can branch code to make changes and then quickly merge it after approval

Log In to GitLab

1. Open a web browser, then navigate to <http://scm.securedatacommons.internal/>, which is the URL of the SDC internal GitLab implementation.

2. In the Sign in page, the LDAP login option is selected by default. Enter your SDC credentials as the LDAP username and password, then click Sign in (in the LDAP Username field, do not enter the “@securedatacommons.com” portion of your SDC username).

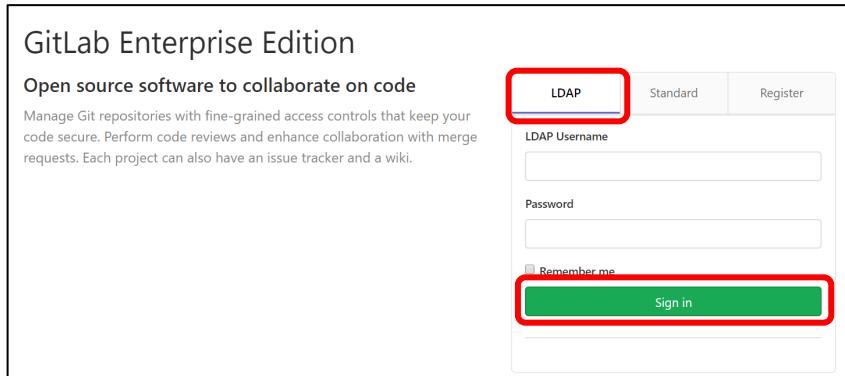


Figure 38: GitLab Login Prompt - LDAP

3. The Welcome to GitLab web page is displayed.

Create a Group

1. In the Welcome to GitLab web page, select Create a group.

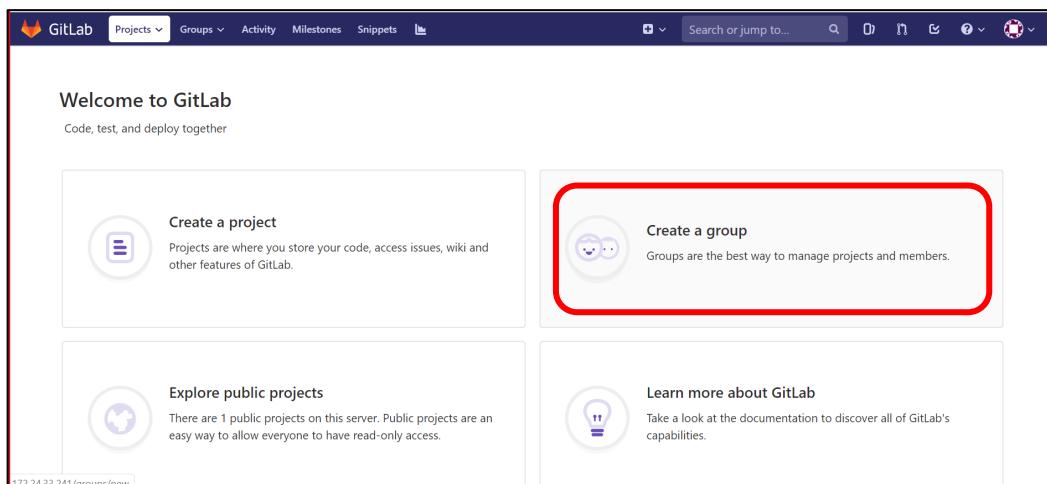


Figure 39: GitLab – Create a Group Option

2. The New group form appears. Enter the name for the new group.
3. Select the Private option under Visibility level to only allow members in the group to see the group and its projects.
4. Select Create group.

The screenshot shows the 'New group' form in GitLab. The 'Group name' field is filled with 'SDC test group'. The 'Visibility level' dropdown is set to 'Private', which is highlighted with a red box. The 'Create group' button at the bottom is also highlighted with a red box. The rest of the form fields like Group URL and Group description are empty or have placeholder text.

Figure 40: GitLab – New Group Form

5. A success message appears indicating that the group was created. You are now ready to add members to the group.

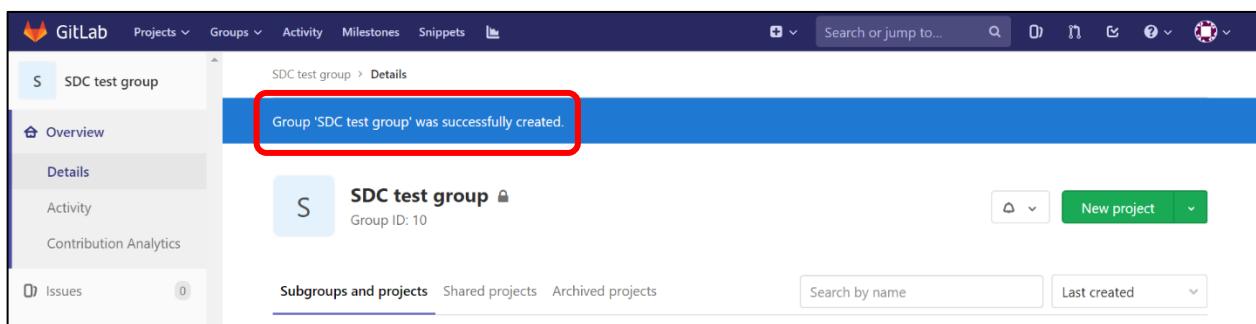


Figure 41: GitLab – Group Creation Success

Add Group Members

Users assigned with the “Developer” or “Maintainer” role can add other GitLab members in the SDC to their groups.

NOTE: SDC users are listed as a GitLab member in the Members panel only after they have logged into GitLab. If you do not see your SDC colleague listed, please ask them to log in to GitLab.

1. In the newly created project web page, click Members from the left navigation menu.

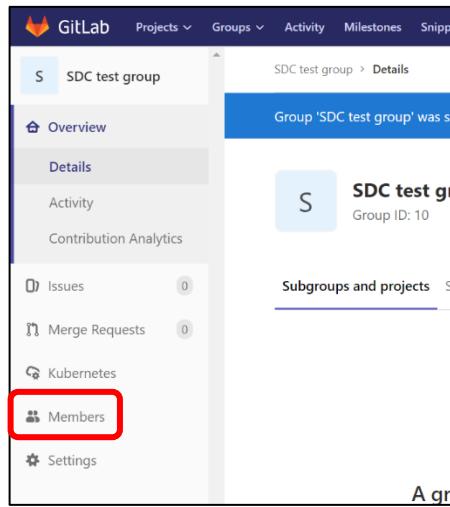


Figure 42: GitLab – Members Left Navigation Menu

2. In the “Search for members” text box, start entering the name of the SDC user you want to add; then select the member’s name from the autocomplete list.
3. Select a role from the next drop-down menu. You can select the “Read more” link to view a comprehensive list of the different available permission levels and their access rights. A summary of the roles’ access types is also shown in Table 1 at the beginning of this chapter.
4. Select Add to group.

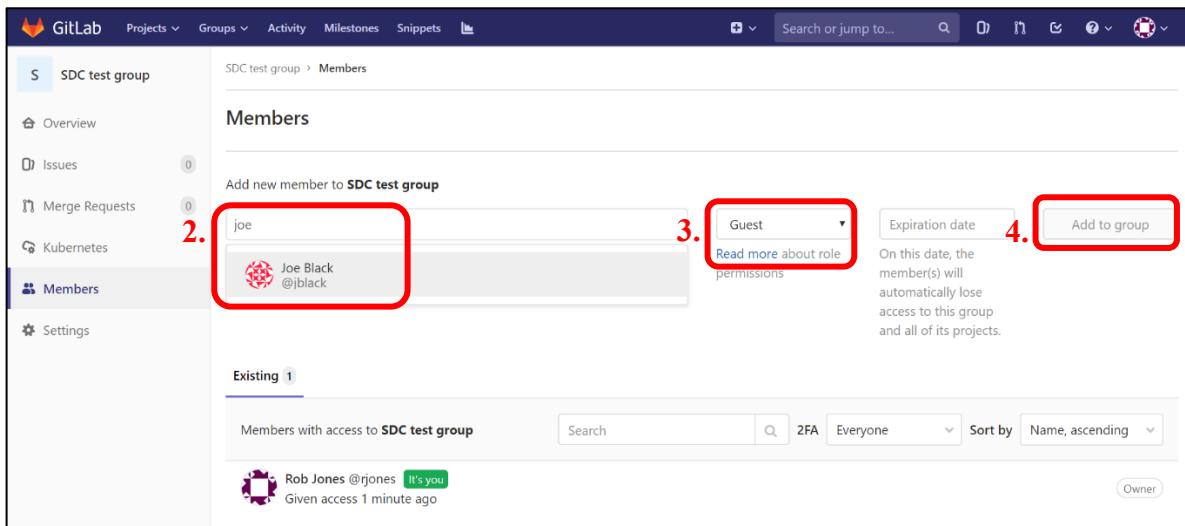


Figure 43: GitLab – Add a New Group Member

Create a Project

1. Select the GitLab icon on the upper left to return to the home page.
2. Select Create a project.

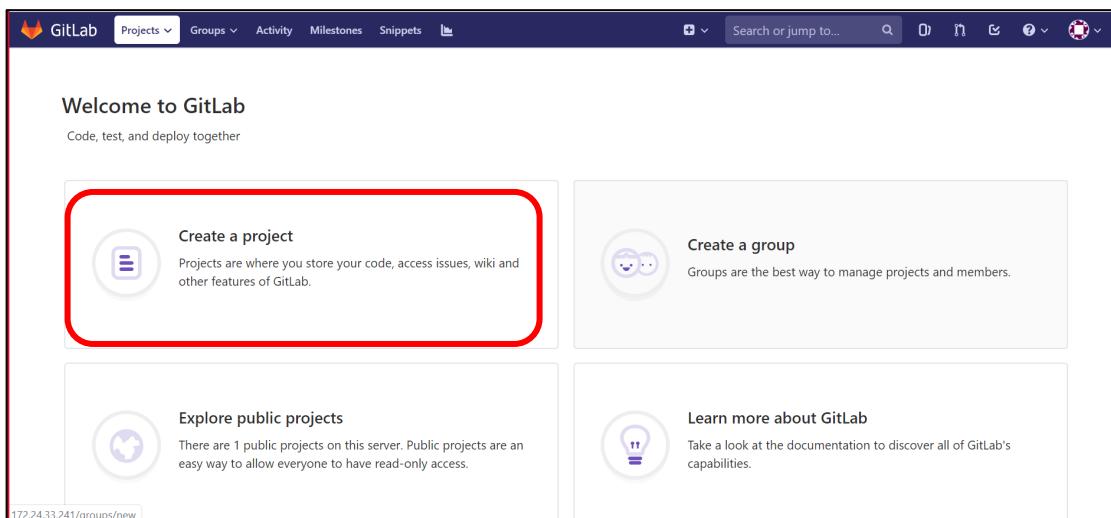


Figure 44: GitLab – Create a Project

3. The New project form appears. Enter a name for the project.
4. Select the Private option under Visibility level, so that the project is only visible to users who are granted explicit access to the project.
5. Select the Initialize repository with README checkbox to automatically generate a readme file that contains the project's repository details.
6. Select Create project.

The screenshot shows the 'New project' form in GitLab. Step 3 highlights the 'Project name' field with the value 'My awesome project'. Step 4 highlights the 'Visibility Level' section where 'Private' is selected. Step 5 highlights the 'Initialize repository with a README' checkbox. Step 6 highlights the 'Create project' button at the bottom.

Figure 45: GitLab – New Project Form

7. A success message appears indicating that the project repository was created. You are now ready to add members to the project.

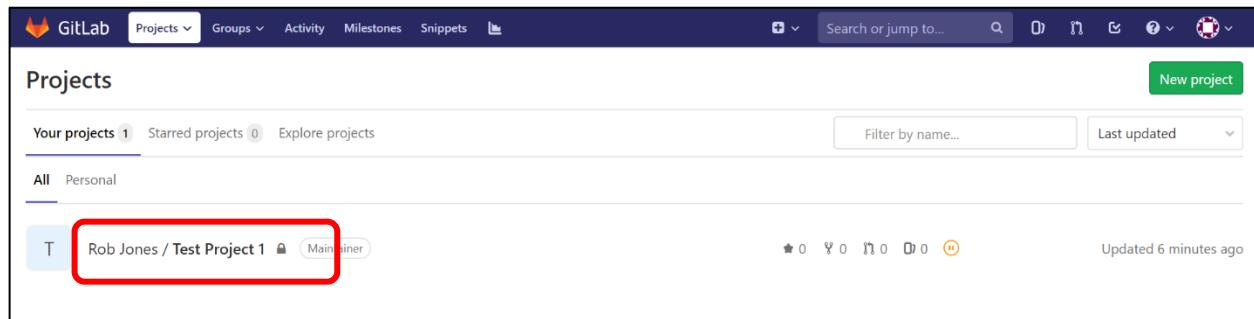


Figure 46: GitLab – Project Creation Success

Add Members or Groups to a Project

1. Select the GitLab icon on the upper left to return to the home page.

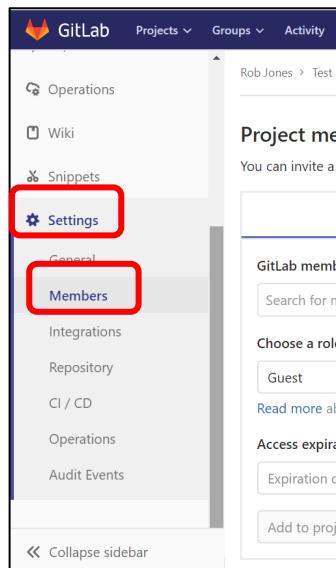
2. A list of your projects appears. Select the name of the project you want to add members /groups to.



The screenshot shows the GitLab interface with the 'Projects' tab selected. At the top, there are navigation links for 'Projects', 'Groups', 'Activity', 'Milestones', and 'Snippets'. Below the header is a search bar with placeholder text 'Search or jump to...'. To the right of the search bar are several icons: a plus sign for creating a new project, a magnifying glass for search, a refresh symbol, a gear for settings, and a user icon. The main area is titled 'Projects' and displays a list of projects. The first project listed is 'Rob Jones / Test Project 1', which is highlighted with a red rectangular box. Other projects visible include 'Starred projects' and 'Explore projects'. On the right side of the project list, there are filters for 'Filter by name...' and 'Last updated', along with a green 'New project' button. Below the project list, there are statistics: 0 stars, 0 forks, 0 issues, 0 merge requests, and 0 pipelines. The status 'Updated 6 minutes ago' is also shown.

Figure 47: GitLab – Select Project to Add Members

3. Click on Settings and then Members from the left navigation menu.



The screenshot shows the 'Settings' page for a project named 'Rob Jones > Test'. The left sidebar contains several options: 'Operations', 'Wiki', 'Snippets', 'Settings' (which is highlighted with a red box), 'General', 'Members' (which is also highlighted with a red box), 'Integrations', 'Repository', 'CI / CD', 'Operations', and 'Audit Events'. The main content area is titled 'Project members' and includes a message 'You can invite a member to this project'. It features a search bar for 'GitLab members', a dropdown for 'Choose a role' (set to 'Guest'), and sections for 'Access expiration' and 'Add to project'. At the bottom of the sidebar, there is a link 'Collapse sidebar'.

Figure 48: GitLab – Settings, Members Left Navigation Menu

4. You can add individual users or entire groups to the project.
 - a. Select the Invite member tab to add members individually to the project.
 - i. Start entering the project member's name or email address into GitLab member or Email address to select from the autocomplete list.
 - ii. Select a role for the project member from the drop-down menu under Choose a role permission.
 - iii. Click Add to project.

The screenshot shows the 'Members' page for 'Test Project 1' in GitLab. The left sidebar shows project navigation options like Project, Repository, Issues, Merge Requests, CI / CD, Operations, Wiki, Snippets, Settings, and Members. The 'Members' option is selected. The main area displays the 'Project members' section with a sub-section for inviting new members. A red box labeled 'a.' highlights the 'Invite member' button. A red box labeled 'i.' highlights the 'Search for members to update or invite' input field. A red box labeled 'ii.' highlights the 'Choose a role permission' dropdown menu, which is currently set to 'Guest'. A red box labeled 'iii.' highlights the 'Add to project' button at the bottom of the form. Below this, there is a section for 'Existing members and groups' showing 'Rob Jones @rjones' with the status 'It's you' and 'Given access 7 minutes ago'.

Figure 49: GitLab – Add Member to Project

- b. Select the Invite group tab to add an entire group to the project.
 - i. Start entering the group's name into Select a group to invite to select from the autocomplete list.
 - ii. Select a role to apply for the group from the drop-down menu under Max access level.
 - iii. Select Invite to grant the group access to the project repository.

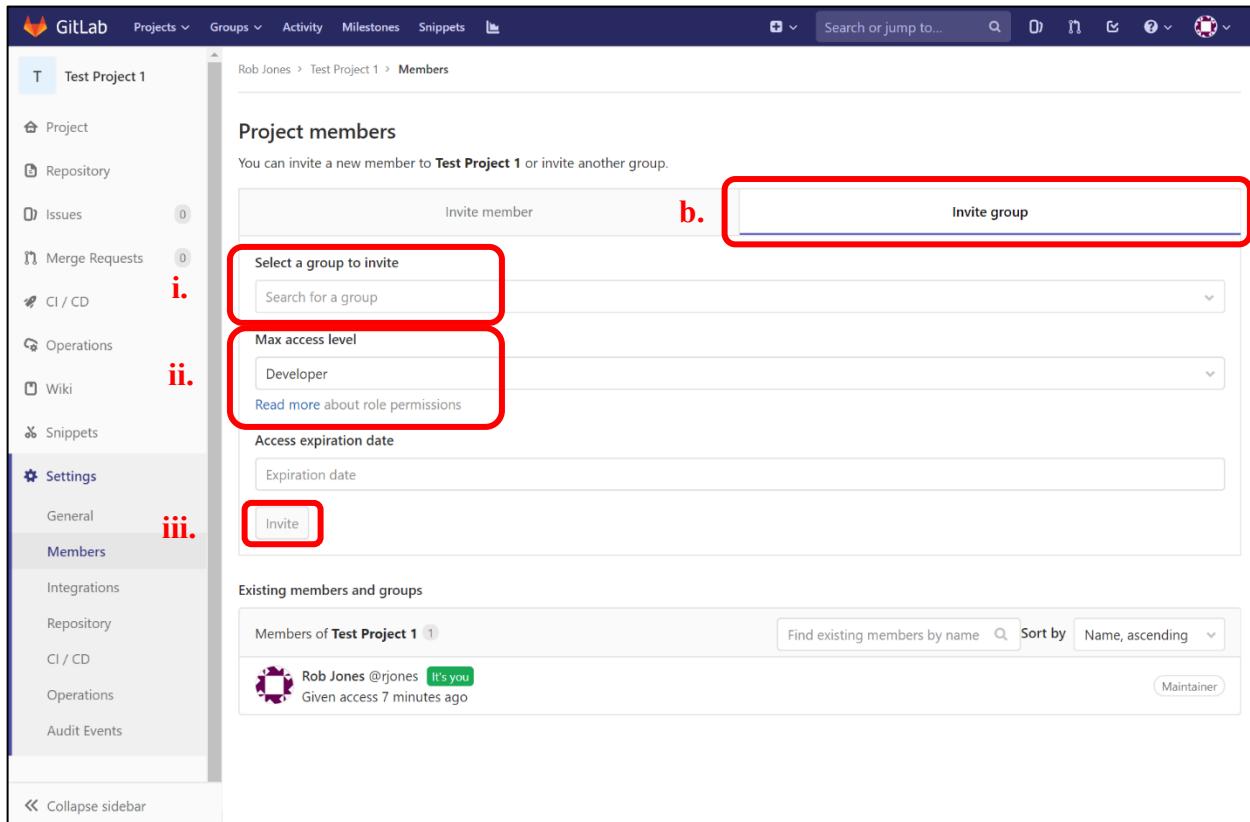


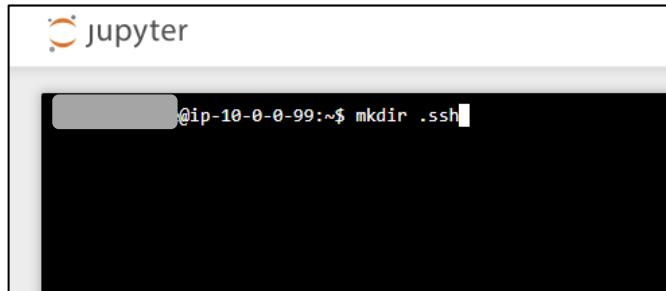
Figure 50: GitLab – Add Group to Project

Generate SSH Key

To set up git integration, you need to have an SSH key on your SDC workstation. This key authenticates your identification to GitLab so that you can push commits to the project repository.

NOTE: For Windows users, perform the following commands in Git Bash. If Git Bash is not installed on your Windows workstation, please copy the Git Bash installer, **Git-2.19.1-64-bit.exe**, from the Z:\software-distribution folder and then install it.

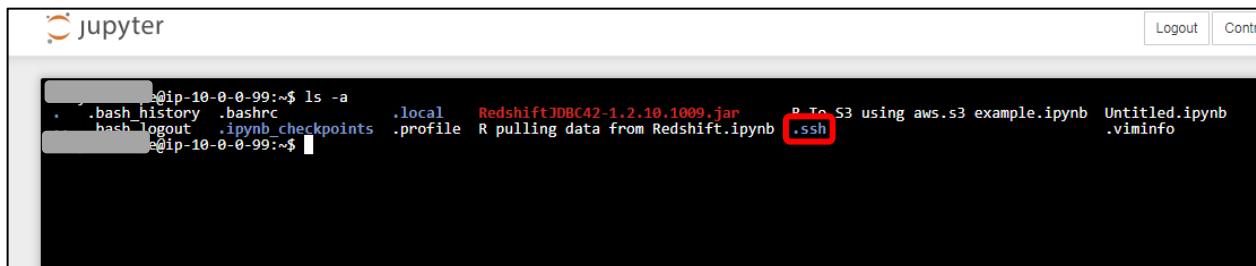
1. To store your authentication credentials for git, use a hidden folder called “.ssh.” First, check to see if the folder exists with `ls -a`. If not, create it using the command `mkdir .ssh`:



```
jupyter
@ip-10-0-0-99:~$ mkdir .ssh
```

Figure 51: GitLab – Create Folder for Credentials

2. Confirm this works by entering `ls -a` to see all objects in the current directory; you should see the .ssh listed:



```
jupyter
Logout Control
@ip-10-0-0-99:~$ ls -a
.bash_history .bashrc .local RedshiftJDBC42-1.2.10.1009.jar P_T_53_using_aws.s3.example.ipynb Untitled.ipynb
.bash_logout .ipynb_checkpoints .profile R_pulling_data_from_Redshift.ipynb .ssh .viminfo
@ip-10-0-0-99:~$
```

Figure 52: GitLab – View Directory Objects

3. Generate a new key to store in this location by typing the following command and then hitting the Enter key:

```
ssh-keygen -t rsa -b 4096 -C "SDCusername@securedatacommons.com"
```

4. It will output the following message: “*Generating public/private rsa key pair.*”
5. It will ask you where to save the file and to enter a passphrase with the following three prompts. Hit the Enter key each time (no need to enter anything else) to accept the defaults.
 - “*Enter file in which to save the key:*”
 - “*Enter passphrase (empty for no passphrase):*”
 - “*Enter same passphrase again:*”
6. After the following confirmation messages, the new SSH key is now stored in the .ssh folder.
 - “*Your identification has been saved in /c/Users/SDCusername/.ssh/id_rsa.*”
 - “*Your public key has been saved in /c/Users/SDCusername/.ssh/id_rsa.pub.*”

The default key file name is **id_rsa** and is saved under the .ssh folder (e.g., C:/Users/[SDCusername]/.ssh/ in Windows; \$HOME/.ssh/id_rsa in Linux)

7. You now need to provide the public part of this key to GitLab by copying the “.pub” part of the key. Here are the different options for copying the public SSH key.

Option 1:

On Windows, open the id_rsa.pub file (under the C:/Users/[SDCusername]/.ssh folder) in a text editor (Notepad) and copy the entire key by using Select All then Copy. Be careful not to accidentally change anything!

Option 2:

Copy your public SSH key to the clipboard by using one of the commands below depending on your Operating System:

Git Bash on Windows:

```
cat ~/.ssh/id_rsa.pub | clip
```

WSL/GNU/Linux (requires the xclip package):

```
xclip -sel clip < ~/.ssh/id_rsa.pub
```

If using Jupyter, use the **head** command. Right-click on the output and click **copy**:

```
head ~/.ssh/id_rsa.pub
```

8. From the GitLab home page, add your public SSH key to your GitLab account:
 - a. In the upper-right corner, click on your avatar image and then click Settings.
 - b. From the left navigation menu, click SSH Keys.
 - c. In the Key field, paste your public SSH key.
 - d. The Title field will be automatically populated with SDCusername@securedatacommons.com.
 - e. Click the Add key button.
9. To test whether your SSH key was added correctly, run the following command in your terminal:

```
ssh -T git@scm.securedatacommons.internal
```

10. The first time you connect to GitLab via SSH, you will be asked to verify the authenticity of the GitLab host you are connecting to. Type “yes” to add the SDC internal GitLab to the list of trusted hosts.
11. You should see the “*Welcome to GitLab, @username!*” message.

NOTE: Once added to the list of known hosts, you will not be asked to validate the authenticity of SDC internal GitLab's host again. Run the above command once more, and you should only receive a “*Welcome to GitLab, @username!*” message.

If the Welcome message does not appear, run SSH in verbose mode by replacing `-T` with `-vvvT` to debug the error:

```
ssh -vvvT git@scm.securedatacommons.internal
```

Chapter 6. Technical Documentation and Contact Information

The following sections provide technical resources for SDC users.

Architecture Diagram

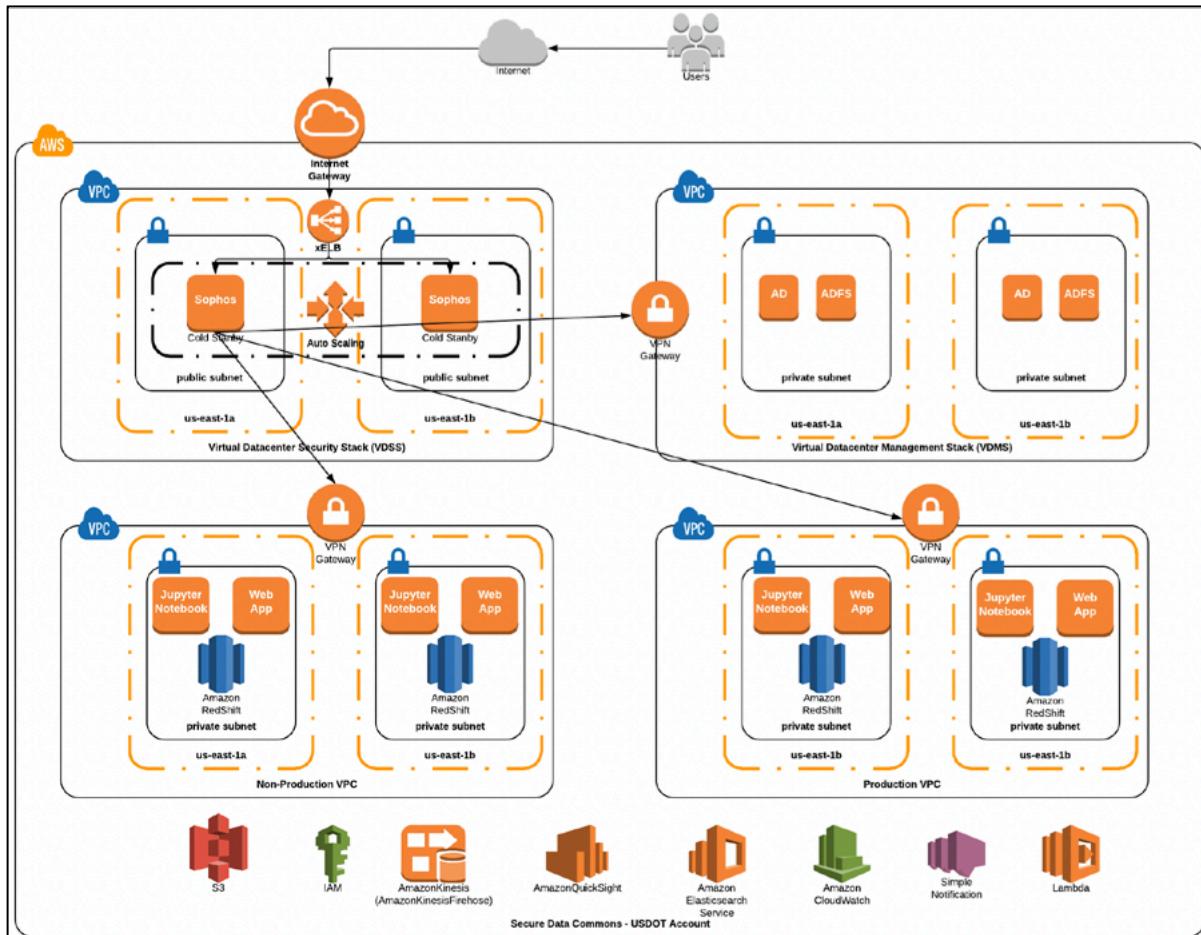


Figure 53: SDC Architecture Overview

Workstation Details

Table 2: 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. Please refer to the [Resize Workstation](#) section under Chapter 3 for instructions.

Tools and Versions

Table 3: 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

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

Useful Links

[S3:](#) Amazon 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 user’s SDC workstation (EC2 instance). The listed commands can only be used from within the user’s workstation environment.

- **List objects in a bucket** - If there are any files at the bucket level, then this 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** - This command will list all the objects/files under that folder or prefix.

```
aws s3 ls s3://<bucketName>/<prefix>/
```

- **Copy a local file to S3** - This command will copy a local file (test.txt) from your workstation over to an S3 bucket:

```
aws s3 cp test.txt s3://<bucketName>/test.txt
```

- **Copy an S3 object to a local file** - The below command copies an S3 object to your workstation.

```
aws s3 cp s3://<bucketName>/test.txt test.txt
```

Chapter 7. 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?](#)
- [Where can I find sample queries for my dataset\(s\)?](#)

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 ‘Datasets’ in the top menu.

The screenshot shows the SDC Portal interface with the 'DATASETS' tab selected. The page title is 'SDC Datasets'. A descriptive text block states: 'These are datasets that are available within SDC platform that are published / enabled by SDC team or published by other users. Access to these datasets are available upon request. Learn more about the different types of datasets and how to request access to these datasets.' Below this is a table listing two datasets:

| 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 54: SDC Portal – Datasets Tab

Chapter 7. Frequently Asked Questions

All the available datasets are listed under the ‘SDC Datasets’. To request access to a dataset, click on the ‘Request’ button. A form will pop up. Fill out the form and click on ‘Send Request’ button.

The screenshot shows the SDC Portal interface. On the left, there's a sidebar with 'My Datasets / Algorithm' and 'SDC Datasets'. The main area displays a table of datasets. One row for 'WAZE' is selected, showing its details: Name (WAZE), Type (Dataset(Curated)), Description (Contains curated waze data), Programming tools / language (---), Geographic Scope (All states in US), Start / End for Data Availability (March 2017 to Present), and Request Access (a blue 'Request' button). A modal window titled 'SDC Datasets Access Request' is open over the table. It has fields for 'Enter the list of state*' and 'Provide comma separated state name (e.g. KN, VI)'. At the bottom of the modal are two buttons: 'SEND REQUEST' (highlighted with a red box) and 'CANCEL'.

Figure 55: SDC Portal – Dataset Access Request

The request will be sent to the SDC 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 ‘Workstations’ and click on the ‘Launch’ button of any workstation you want to access.

The screenshot shows the SDC Portal interface with the 'WORKSTATIONS' tab selected. The main area is titled 'My Workstations'. It displays a table of workstations:

| # | Stack Name | Applications | Action |
|---|-------------------------------------|---|---|
| 1 | Programming Environment #1 | Microsoft R, RStudio, Python, Microsoft Power BI, SQL Server Management Studio, SQL Workbench, Open Office, Firefox | <button>Stop</button> <button>Launch</button> |
| 2 | Programming Environment (AppStream) | Firefox, Jupyter Notebook, RStudio | <button>Stop</button> <button>Launch</button> |

Figure 56: SDC Portal – Workstations Tab

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

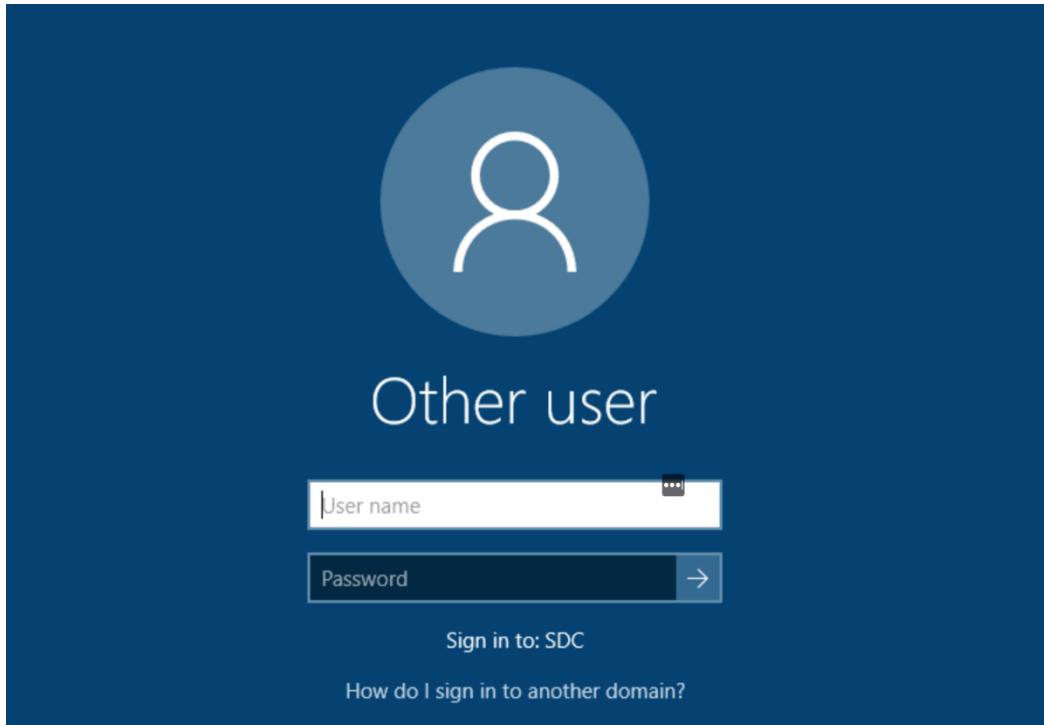


Figure 57: SDC Workstation – Login

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?

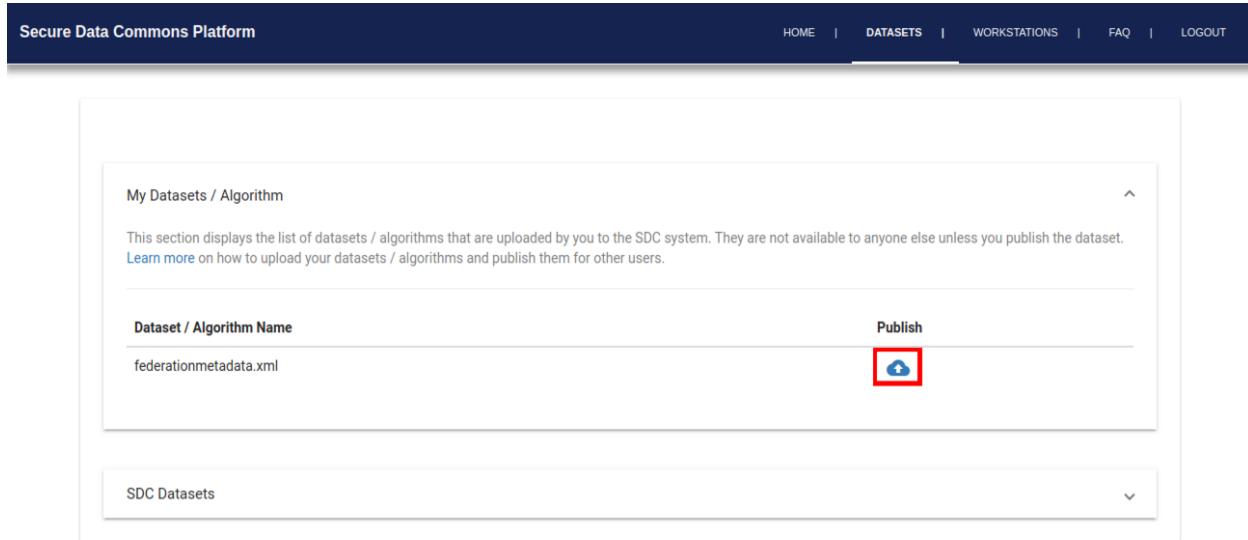
Please refer to [Upload User Data to S3 Bucket Through Portal](#) to bring your own datasets/algorithm to workstation.

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 Publish button for the dataset/algorithm you wish to publish.

Chapter 7. Frequently Asked Questions



The screenshot shows the SDC Portal interface. At the top, there is a dark blue header bar with the text "Secure Data Commons Platform" on the left and navigation links "HOME", "DATASETS", "WORKSTATIONS", "FAQ", and "LOGOUT" on the right. Below the header, the main content area has a white background. It features a section titled "My Datasets / Algorithm" with a sub-instruction: "This section displays the list of datasets / algorithms that are uploaded by you to the SDC system. They are not available to anyone else unless you publish the dataset. [Learn more](#) on how to upload your datasets / algorithms and publish them for other users." Below this, there is a table-like structure with two columns: "Dataset / Algorithm Name" and "Publish". The first row contains the name "federationmetadata.xml" and a "Publish" button, which is highlighted with a red box. At the bottom of the page, there is a section titled "SDC Datasets" with a small downward arrow icon.

Figure 58: SDC Portal – Publish Button

3. In the pop-up window, there are two options for the Type: either a **Dataset** or **Algorithm**.
 - a. From the Type drop-down menu, select **Dataset**:

The screenshot shows a web-based form titled "Publish Dataset / Algorithm Request". The "Type" field is set to "Dataset". The "Name *" field contains "appstream.py". The "Description *" field is empty. The "File/folder name *" field is empty. The "Readme / Data dictionary file name *" field is empty. The "Geographic Scope *" field is empty. The "Start/End Date for Data Availability *" field is empty. At the bottom are two buttons: "SEND REQUEST" and "CANCEL".

Figure 59: SDC Portal – Publish Dataset Request 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.

Chapter 7. Frequently Asked Questions

- 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. From the Type drop-down menu, select **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 60: SDC Portal – Publish 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.

Where can I find sample queries for my dataset(s)?

Sample queries are provided for each of the datasets on a GitLab page we have set up for code sharing: http://scm.securedatacommons.internal/sdc_code_collaboration/sdc-hadoop-sql-query-examples. Please refer to the wiki page of a specific dataset to find its related queries there.

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