

Objective

Operationalize a data stack. The components would have the following utilities –

1. ELT (data loader)
2. Orchestrator for transforms
3. Data Store
4. Dashboards/ Visualizer

Here is a simplistic diagram for project data flow, setting up each of them will require multiple AWS services for instantiation



Data Source

Salesforce Opportunity daily slice will be converted to sales coverage trend metrics

Execution Plan

These are the details of execution that will be carried out for this project:

Local/ Sandbox/ Test

1. Salesforce source access
 - a. Get access to Salesforce Sandbox Environment, setup API access credentials
2. Install ELT Service Airbyte (local mac)
3. Configure Salesforce connection in Airbyte
4. Create Dockerized container environment for python code runs, and code deployment
5. Docker push to ECR

Production/ AWS Cloud

1. Salesforce Prod access
 - a. Get access to Salesforce Sandbox Environment, setup API access credentials
2. Install ELT Service Airbyte on AWS EC2
3. Instantiate a Redshift cluster on AWS
4. Instantiate Airflow instance (managed Airflow) - AWS MWAA
5. Docker push to ECR for Airflow jobs/ Redshift transforms
6. Use Tableau to read/ consume Redshift transforms on Salesforce opportunity trend metrics

Pre-requisites

1. Local Machine

All workings detailed are performed on a Macbook M1 machine. If you are working on Linux machines, or Macbook Intel chip machines, the commands and settings can be replicated verbatim. For Windows machines, the local system workings/ commands may not work as mentioned here.

2. Docker

o Context -

- Docker Desktop is an easy-to-install application for your Mac, Linux, or Windows environment that enables you to build and share containerized applications and microservices. It provides a simple interface that enables you to manage your containers, applications, and images directly from your machine without having to use the CLI to perform core actions.
- We will be running Python code in the Docker containers. The python code will be used for data transforms on Redshift. Using Docker containers allows us to replicate airflow host settings through exposing environment variables.

o Signup/ Access –

- start at <https://hub.docker.com/>
- signup with email
- opt for single-user free tier \$0 (*to begin with*)

o Download & Install –

- <https://www.docker.com/products/docker-desktop/>

o Docker Desktop should be running when trying to test/ execute commands locally.

TEST/ Check#1: On Mac Menu, you should be able to see the Docker icon, IF the docker desktop is active/ running



TEST/ Check#2: On Mac Terminal : run to see all running containers

docker ps --all

if this command runs fine (i.e. no error exceptions, even with no running containers) - it indicates docker is ON

This is what we see in terminal for Docker dependent commands, if Docker isn't running

```
[djenadewi@jenia-mac deng-jobs % make env-up
docker build -t deng-jobs:local
Cannot connect to the Docker daemon at unix:///var/run/docker.sock. Is the docker daemon running?
make: *** [env-up] Error 1
djenadewi@jenia-mac deng-jobs % ]
```

- When pushed to production environment, the docker images are pushed into AWS – ECR (Elastic Container Repository) for workflow read and pipeline execution
2. Python & Poetry
- Install Python, using a consistent package manager (like brew, conda, pip , etc)
 - Although latest python versions are advised, we will use Python 3.9 / 3.10 since I find these versions most stable and lesser tweaks to deal with package dependencies
 - With higher python versions managing python dependencies becomes a routine and this is where using consistent package manager, or poetry like dependency managers becomes key
 - Consistent use of package manager is key, and so is maintaining package version dependencies
 - Recommend to use poetry for version and package dependency management
 - Poetry is a tool for **dependency management and packaging** in Python. It allows you to declare the libraries your project depends on and it will manage (install/update) them for you. Poetry offers a lockfile to ensure repeatable installs, and can build your project for distribution.
 - Create the project using poetry i.e., start from here <https://python-poetry.org/docs/basic-usage/>
 - Initialize poetry – it results in a **pyproject.toml** file
 - To add any package say pandas use \$ poetry add pandas
 - this poetry add commands reflects in entry in pyproject.toml and poetry.lock file is the one which controls the registry versions of the installed packages for the project

3. Email & Distribution Alias

- We will need a mailbox for instantiating Airflow. Ideally we need 2 email ids. The first one is a full inbox to be used by airflow for its smtp configuration. This mailbox mimics any human user, and has the same features as a regular mailbox. (Note: singup for a gmail account, and ensure to enable IMAP settings). The second is a creation of a group email alias, the group email has the formerly created mailbox as its member. The idea is to add as many users as members of the group, and the airflow email account is a member of the distribution.

Implementation Log

1. Salesforce Access (Sandbox environment)

We get Salesforce credentials as username and password. However for programmatic ingestion, we will need api access methods, which follows OAuth2 authentication. Here are the steps -

- a. Get access to Salesforce (username & password)

- b. Salesforce Sandbox environment domain eg: <https://test.salesforce.com>
 - c. You will also need the CLIENT_ID for the Salesforce credentials (which looks something like 3MVG96LA2t1yu9WK_2g3tejYpVb0YWaH0dj8gu00Ao7Tv.RHqpRWTitD0gLl5gDy6..ZsBZ.WQXte1XntAaoH)
 - d. Can create a connected app in Salesforce for all sfdc interactions
 - i. details outlined here
 1. https://docs.datawatch.com/swarm/desktop/Generating_a_Client_ID_and_ClientSecret_Key_for_Salesforce_Connections.htm
 - ii. Also refer these:
 1. <https://community.boomi.com/s/article/howtoconnecttosalesforcerestapiwithoAuth20>
 2. <https://medium.com/@bpmmendis94/obtain-access-refresh-tokens-from-salesforce-rest-api-a324fe4cccd9b>
 - e. So summarily: inputs for SFDC access are username, password, and CLIENT_ID (from sfdc- admin or connected app setup)
 - i. Using these credentials we will need CLIENT_ID, CLIENT_SECRET and REFRESH_TOKEN for use in subsequent steps
 - ii. I had to use postman to post the api methods and code, and finally landed on the refresh token screen as this (P.S. you can use any rest api handler, I used postman by signup with a free account on <https://www.postman.com/>)

Overview [POST https://test.salesforce.com](https://test.salesforce.com) + [...](#) No Environment

https://test.salesforce.com/services/oauth2/token?code=aPrxmCocI5Str3ckbLStugIPUR7OrgNNy whole URL

Save [Edit](#)

POST https://test.salesforce.com/services/oauth2/token?code=aPrxmCocI5Str3ckbLStugIPUR7OrgNNy whole URL

Params • Authorization Headers (7) Body Pre-request Script Tests Settings Cookies

Body Cookies (3) Headers (10) Test Results Status: 200 OK Time: 664 ms Size: 1.89 KB Save Response

Pretty Raw Preview Visualize JSON

```
1 "access_token": "00DD5000000msU!AcQcAQEYKHz95z4Lvhv4scXA...4cnyjD5CBPu0U7Agoya1KmUx0mN4wRvcI4UYAKfbDpbk1Kv1JYBdFB8wtFqeQzH011vj",
2 "signature": "D5aJxutuQKqeeFuHwIPQoEXFksgePYUyOf4KG3aPxss=",
3 "scope": "full",
4 "id_token": "eyJraiwQ0i1YmgILC0eXai0i1Kv10i1Cjhbc0i0jS1uI1n19.
5 eyJhdF9oYXNlOi1tM5X1pRwB1RjdfNETfMK1RsMfDyIsIn1N1Yi16mh0dBz018vdGv2dC5zYwx1c2ZvcmN1LmNvbS9pZC8wMEREUzAwMDAwMG1zVwMyQukvMDA1MGgwMDAwME10en3JQUF5IiwiYXvKIjo1M01WrxK2TEy0df5t1X1S1yZn20w2pzcFzT11YuGuwgZ423wUf0fVN1R2L1J3cXBsv1RpdeRPZ0xsNwEEtYulpzQ1ov1FyGUxw500wFvSClsIm1zcy16im0oHbz018vdGv2dC5zYwx1c2ZvcmN1LmNvbS1Im4C16MT2NQ0MD12YwiaWF01joxNjY1NDQwMTQzIQ.
6 Qw8R1WiegdsDcARLkMq3BCADGGYODC_Gx1upqJ15n0w1mHy2WMSHwGhLhMeYUKvkDjshd0SpPV1TgNJAzzBdz1TrwqEA3ZhFavZ1Wb0Bu1nRoFu923dTiam2R2YHg04x2k0h0p9xp1PxVxJjfkmk25zIkgRzrXuIZAuuz1rJ2aaHxDXFuccBGDBSK9fxggf1BpN0BaZ8_MFPErAUuo2t0xihbBwnJwdN9-17K0ksQ0tfifl_1Ms9@0t6Mfuah0J4a8f19qAWKngE151D1JK81z1230npijYtB1zuVP06tVK4fG212TJ5lVRxAcou515tHdNWd7aA1b1vA8j1qj6mEqTTp18oncsbe1jV0NqjwENyD3BiYffKKRAK0Sxx0fK02GMhrc5jap3ayhe22EIA9px1zdw41yu7uw-fQCjAtV1zhuUo0ep0pozjw4010PrlqJAT1GfIEExaQvShIg9Qf0U0QjPka8tWuQzdB0q3fmoknmVnHjmP_",
7 "instance_url": "https://sumologic--uat.sandbox.my.salesforce.com",
8 "id": "https://test.salesforce.com/id/00DD5000000msUcB2AI/0050h00000ItzrNAAR",
9 "token_type": "Bearer",
10 "issued_at": "1665440143795"
```

So, the credentials I had generated were for CLIENT_ID, SECRET_KEY and REFRESH_TOKEN.

- f. We will need these credentials in the next step to ingest Salesforce Opportunities object into a data store

2. Airbyte Installation

The installation is for local macbook as well as EC2 machines. Security and networking aspects associated with EC2 machines is beyond the scope of this project, but a necessary aspect for the system to be operational. Typically bastions are configured to secure access to EC2, along with dedicated keys along with user name, password, IP authentication methods.

- a. An ELT service like ELT makes it faster to ingest any such custom apps / sources. An alternate (and bit more time consuming) approach would be to design for an api ingestion to consume from any api endpoint and load into a data store. However, my intent is to pump the raw data faster through some means into the data store, so that using/ processing of the raw data can be done using SQL (which is my comfort skill, and easier to understand from business stakeholder aspect too), so I chose to deploy ELT services. Similar other open-source ELT service is meltano (<https://docs.meltano.com/getting-started/installation>). Use of Airbyte abstracts a lot of the heavylifting required for custom apis. Also if we see the connector catalog, we find pretty good coverage with Airbyte for a lot of the apps and integrations.
 - i. There are also quite a few licensed ELT providers where environment setup can be done in few minutes. Stitch¹, Fivetran² are 2 names that are top of mind recalls in this area.
 - b. For local installation, its nicely detailed and simple, following instructions on <https://airbytehq.github.io/deploying-airbyte/local-deployment>
 - c. Here are the commands to use:

```
git clone https://github.com/airbytehq/airbyte.git  
cd airbyte  
docker-compose up
```

```
Last login: Tue Nov 1 19:18:40 on ttys002
bjenna@bjenna-mac ~ % cd Documents/Github
bjenna@bjenna-mac GitHub % git clone https://github.com/airbytehq/airbyte.git
Cloning into 'airbyte'...
remote: Enumerating objects: 510997, done.
remote: Counting objects: 100% (4710/4710), done.
remote: Compressing objects: 100% (1373/1373), done.
remote: Total 510997 (delta 2234), reused 4468 (delta 2134), pack-reused 506287
Receiving objects: 100% (510997/510997) 597.79 MiB | 1.91 MiB/s, done.
Resolving deltas: 100% (281923/281923), done.
Updating files: 100% (11589/11589), done.
bjenna@bjenna-mac GitHub % cd airbyte
bjenna@bjenna-mac airbyte % docker-compose up
[WARN 0000] The "RUN_DATABASE_MIGRATION_ON_STARTUP" variable is not set. Defaulting to a blank string.
[WARN 0000] The "NEW_RELIC_APP_NAME" variable is not set. Defaulting to a blank string.
[WARN 0000] The "SECRET_PERSISTENCE" variable is not set. Defaulting to a blank string.
[WARN 0000] The "JOB_ERROR_REPORTING_SENTRY_DSN" variable is not set. Defaulting to a blank string.
[WARN 0000] The "WORKER_ENVIRONMENT" variable is not set. Defaulting to a blank string.
[WARN 0000] The "GITHUB_STORE_BRANCH" variable is not set. Defaulting to a blank string.
[WARN 0000] The "DEPLOYMENT_MODE" variable is not set. Defaulting to a blank string.
[WARN 0000] The "LOG_CONNECTOR_MESSAGES" variable is not set. Defaulting to a blank string.
[WARN 0000] The "SECRET_PERSISTENCE" variable is not set. Defaulting to a blank string.
[WARN 0000] The "JOB_ERROR_REPORTING_SENTRY_DSN" variable is not set. Defaulting to a blank string.
[WARN 0000] The "DEPLOYMENT_MODE" variable is not set. Defaulting to a blank string.
[WARN 0000] The "REMOTE_CONNECTOR_CATALOG_URL" variable is not set. Defaulting to a blank string.
[WARN 0000] The "TEMPORAL_HISTORY_RETENTION_IN_DAYS" variable is not set. Defaulting to a blank string.
[WARN 0000] The "UPDATE_DEFINITIONS_CRON_ENABLED" variable is not set. Defaulting to a blank string.
[+] Running 1/1
  1. worker Pulling
    ↗ 371f42f94d Waiting
    ↗ dde8ecef1337 Waiting
    ↗ 2c7446c91a79 Waiting
    ↗ 3501992f2158 Waiting
    ↗ ccd6da350bcb Waiting
    ↗ 9e259bac12296 Waiting
    ↗ 368865401d4 Waiting
    ↗ 73984f31054 Waiting
    ↗ 7cc59c95faed Waiting
    ↗ 4f4fb790e754 Waiting
    ↗ 284d4a9ba227 Waiting
  1. init Pulling
    ↗ 90cd3b7c325 Waiting
    ↗ 6c33fd4dbf481 Waiting
    ↗ 3501992f2158 Waiting
  1. airbyte-cron Pulling
    ↗ 7fa5051a0bf9e Waiting
  1. airbyte-proxy Pulling
    ↗ d5d519d6fc13 Already exists
```

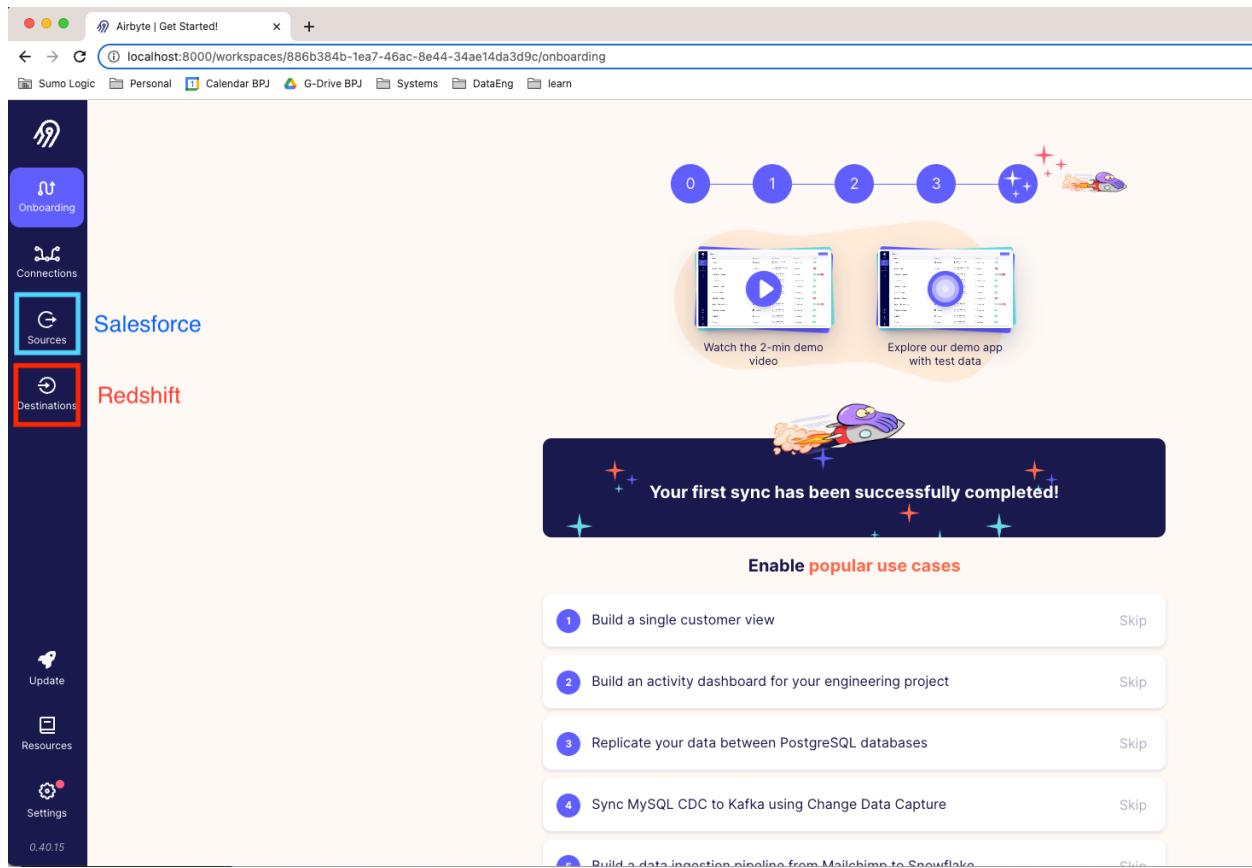
¹ <https://www.stitchdata.com/>

² <https://www.fivetran.com/>

d. For EC2 installation i.e. to productionalize Airbyte, the steps are listed <https://airbytehq.github.io/deploying-airbyte/on-aws-ec2>. Beyond bringing up the application, providing SSH keys for secure access, securing the EC2 host behind a bastion, having RDS (Postgres/MySQL) datastore as the application backend, and enabling cloudwatch logs for application monitoring are some of the additional steps that need to be done in the AWS environment.

e. Airbyte can be accessed at <http://localhost:8000>

In a nutshell, The links marked as “Sources” and “Destinations” are linked/ made operational through “Connection”



3. Add Salesforce Source to Airbyte

a. For trend analysis, I will be extracting **Opportunity** object alone. The credentials generated from #1 are provided on a “New Source”. The credentials are provided here...

Salesforce

Prereq

- Sale
- Ded
- (For

Step 1:

1. Log
2. On t
3. In th
4. For l
5. Clie
6. Scrc
7. Scrc
8. On t
9. Use
10. Fill c

4. Redshift Instance

- The core of this data stack is the analytic data store. It can be called Data Mart or Data Warehouse depending on the scope of use cases it supports. For this task, it will host and process only Opportunity object from Salesforce (so can be referred to as a Marketing Data Mart).
- Access to AWS, IAM role provisioning, firewall rules, MFA enforcement and having rights on the account are necessary aspects which I wont be going into as much detail. Instead I will be focusing on the Redshift instantiation following the instructions <https://docs.aws.amazon.com/redshift/latest/gsg/bring-own-data.html>
- While the instance is being instantiated and initialized, need to have a bit futuristic vision of data access w.r.t. roles and permissions on the database object organizing.

- d. So the redshift instance will be created resulting in super-admin user which will be used to create admin user for all subsequent role grants and user administration.
 - e. For redshift user creation the command would be as:
- ```
create user xyz in group read_only_users password 'abc';
```
- f. For creating a Redshift cluster,
    - i. signin to AWS,
    - ii. type Redshift in search bar,
    - iii. create cluster
    - iv. provide cluster configuration details – identifier, node type, nodes
      - 1. this is result of a planning/sizing exercise
      - 2. chose dc2.large and 2 nodes for this assignment
    - v. provide these values:
      - 1. database name: db\_prod
      - 2. database port: 5439
      - 3. admin user: admin
      - 4. admin password: awsuser

This is how ready-to-use Redshift instance looks like

The screenshot shows the 'Properties' tab of a Redshift cluster named 'redshift-cluster-1'. The top section displays cluster metadata: Identifier (redshift-cluster-1), Namespace (10467e54-fb64-4e16-9703-8b192a6359d5), Status (Available), Node type (dc2.large), and Number of nodes (2). The Endpoint URL is highlighted with a red box. Below this, the JDBC URL and ODBC URL are listed. The bottom section, titled 'Database configurations', lists database parameters: Database name (dev), Port (5439), Admin user name (admin), Parameter group (deng-redshift-parameter), SSH ingestion setting (cluster public key), Encryption (Disabled), AWS KMS key ID (-), Audit logging (Enabled), Log export type (S3 bucket), S3 bucket name (sumo-dw-redshift-logs), S3 key prefix (redshift/), and Last successful delivery (November 08, 2022, 14:06 UTC-08:00). A red box highlights the 'Database name' row in the database configurations table.

| Cluster identifier<br>redshift-cluster-1                  | Status<br><span style="color: green;">Available</span> | Node type<br>dc2.large | Endpoint<br><a href="#">redshift-cluster-1.c6audqj62kaf.us-west-2.redshift.amazonaws.com</a>                                                              |
|-----------------------------------------------------------|--------------------------------------------------------|------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------|
| Cluster namespace<br>10467e54-fb64-4e16-9703-8b192a6359d5 | Date created<br>September 29, 2022, 13:11 (UTC-07:00)  | Number of nodes<br>2   | JDBC URL<br><a href="#">jdbc:redshift://redshift-cluster-1.c6audqj62kaf.us-west-2.redshift.amazonaws.com:5439/dev</a>                                     |
|                                                           | Storage used<br>30.47% (97.51 of 320 GB used)          |                        | ODBC URL<br><a href="#">Driver={Amazon Redshift (x64)};Server=redshift-cluster-1.c6audqj62kaf.us-west-2.redshift.amazonaws.com;Port=5439;Database=dev</a> |

Cluster performance | Query monitoring | Schedules | Maintenance | **Properties**

**Database configurations**

| Database name<br>dev     | Parameter group<br>Defines database parameter and query queues for all the databases.<br><a href="#">deng-redshift-parameter</a> | Encryption<br>Disabled | Audit logging<br>Enabled                                         |
|--------------------------|----------------------------------------------------------------------------------------------------------------------------------|------------------------|------------------------------------------------------------------|
| Port<br>5439             | SSH ingestion setting (cluster public key)<br><a href="#">ssh-rsa AAAAB3NzaC1yc2EAAAQAB...</a>                                   | AWS KMS key ID<br>-    | Log export type<br>S3 bucket                                     |
| Admin user name<br>admin |                                                                                                                                  |                        | S3 bucket name<br>sumo-dw-redshift-logs                          |
|                          |                                                                                                                                  |                        | S3 key prefix<br>redshift/                                       |
|                          |                                                                                                                                  |                        | Last successful delivery<br>November 08, 2022, 14:06 (UTC-08:00) |

[Edit admin credentials](#) [Rotate encryption keys](#) [Edit ▾](#)

**Network and security settings**

|                                           |                                  |                                                                                                     |                                                                           |
|-------------------------------------------|----------------------------------|-----------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------|
| Virtual private cloud (VPC)<br>[REDACTED] | Availability Zone<br>us-west-2b  | VPC security group<br>Specify which instances and devices can connect to the cluster.<br>[REDACTED] | Publicly accessible<br>Allow connections from outside the VPC.<br>Enabled |
| Subnet<br>amazonmmwaa-subnet-group        | Enhanced VPC routing<br>Disabled |                                                                                                     |                                                                           |
| Endpoint URL<br>[REDACTED]                |                                  |                                                                                                     |                                                                           |

**Cluster permissions**

Create an IAM role as the default for this cluster that has the [AmazonRedshiftAllCommandsFullAccess](#) policy attached. This policy includes permissions to run SQL commands to COPY, UNLOAD, and query data with Amazon Redshift. The policy also grants permissions to run SELECT statements for related services, such as Amazon S3, Amazon CloudWatch logs, Amazon SageMaker, and AWS Glue.

**Associated IAM roles (1) Info**

Create, associate, or remove an IAM role. You can associate up to 50 IAM roles. You also choose an IAM role and set it as the default for this cluster.

| IAM roles           | Status  | Role type |
|---------------------|---------|-----------|
| s3_to_redshift_role | In-sync | --        |

- vi. I have ignored detailed around IAM role creation and managing access. One key aspect is configuring rules with security groups around inbound SQL clients. So the idea is, Redshift instance is not a public data store. We need to control who can access what.

1. Redshift roles and groups controls the user's permissions and access levels within Redshift
2. Modify security group rules to allow inbound traffic only from Redshift recognized users. Here's how inbound SQL clients can be limited only to the IP's bound in the security group of the Redshift cluster

**Security Groups (1/1) Info**

Filter security groups

Security group ID: sg-07e532d23d0518cb7

Actions Export security groups to CSV Create security group

| Name     | Security group ID | Security group name | VPC ID     | Description                | Owner        | Inbound rules count   | Outbound rules   |
|----------|-------------------|---------------------|------------|----------------------------|--------------|-----------------------|------------------|
| RedShift | [REDACTED]        | default             | [REDACTED] | default VPC security gr... | 412926778423 | 10 Permission entries | 1 Permission ent |

**Inbound rules (10)**

Filter security group rules

| Name                   | Security group rule... | IP version | Type     | Protocol | Port range | Source                 | Description        |
|------------------------|------------------------|------------|----------|----------|------------|------------------------|--------------------|
| Nicholas Zhu Home...   | sgr-0ca12fde68a05b6    | IPv4       | Redshift | TCP      | 5439       | 67.180.196.98/32       | nzhu Home IP       |
| EC2 Airbyte Secur...   | sgr-01f7a1503bdf60302  | IPv4       | Redshift | TCP      | 5439       | 52.39.94.219/32        | EC2 Airbyte        |
| RWC Office IP (20+...) | sgr-02664c24c58a2f60b  | IPv4       | Redshift | TCP      | 5439       | 50.200.63.82/32        | Kay RWC office ip  |
| Tableau us-west-2b...  | sgr-06192a4bd4ead7bcc  | IPv4       | Redshift | TCP      | 5439       | 34.214.85.54/32        | Tableau us-west-2l |
| -                      | sgr-0a46008219d7e9...  | -          | Redshift | TCP      | 5439       | sg-0de4f40762cbd2b...  | -                  |
| binay home IP          | sgr-0383be6b88c283e... | IPv4       | Redshift | TCP      | 5439       | 99.21.67.216/32        | binay home ip      |
| -                      | sgr-02f69b9218ae095f   | -          | H1TIPS   | TCP      | 443        | sg-07e532d23d0518cb... | -                  |
| Nicholas Zhu add IP... | sgr-042150179dd23...   | IPv4       | Redshift | TCP      | 5439       | 24.23.251.78/32        | nzhu 2nd Home IP   |
| Tableau us-west-2b...  | sgr-0a13243e5f849cbfc  | IPv4       | Redshift | TCP      | 5439       | 34.214.85.244/32       | Tableau us-west-2l |
| Kay Home IP (14-oc...) | sgr-0aaadcc8820781b7a  | IPv4       | Redshift | TCP      | 5439       | 104.56.116.63/32       | Kay Home IP        |

- vii. Since it is an analytic data store, we need monitors with Redshift queries , disk usage, CPU, etc. CloudWatch allows for these configurations. However, I will skip the monitoring piece on Redshift to pivot to project scope.

- g. To access the Redshift cluster,
- i. can use clients like DBeaver, Aginity, SQL Workbench, Datagrip, etc with the following credentials:
    1. Hostname
    2. Port
    3. Username
    4. Password
    5. Database
  - ii. For each new user, the inbound rules need to be updated along with access provision
- h. Create a Redshift user credentials (service user) to be used in Airbyte data load . So we want to load the Salesforce object to Redshift.
- i. Create a separate schema for Airbyte access with the created user.
  - ii. I have named the schema **ab\_sfdc** , destination of sfdc raw data in Redshift

## 5. Airbyte Redshift Destination

The credentials created in #4 above are used to setup the Airbyte Redshift destination.

**Destination Settings**

**Destination type**

Redshift BETA

Beta connectors are in development but stable and reliable and support is provided. See our [documentation](#) for more details.

**Destination name** ⓘ sfdc\_prod

**Host** ⓘ redshift-cluster-1.c6audql62kaf.us-west-2.redshift.amazonaws.com

**Port** ⓘ 5439

**Username** ⓘ user\_airbyte\_elt

**Password** ⓘ \*\*\*\*\* Edit

**Database** ⓘ db\_prod

**Default Schema** ⓘ ab\_sfdc

**JDBC URL Params** ⓘ Optional

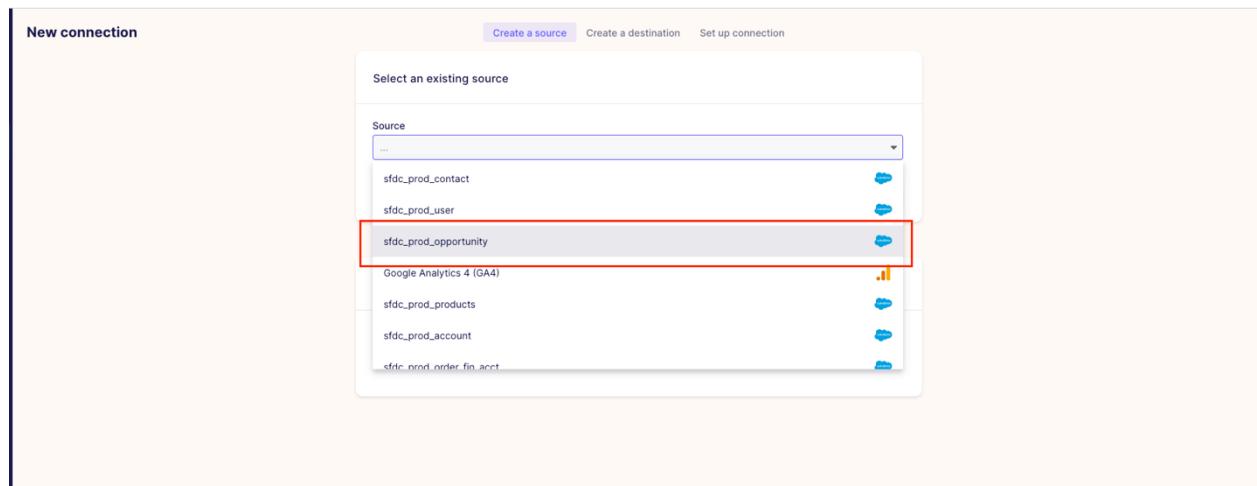
**Uploading Method** ⓘ Standard

**NOTE:** It's important to have users isolated from each other. So in Redshift we create a separate schema for each source, and a dedicated user for each schema. So in this example `user_airbyte_elt` is for sfdc source, all objects of which get loaded into `ab_sfdc` schema in Redshift. If this detail isn't practiced, there's a risk of one source malfunction impacting others which isn't ideal.

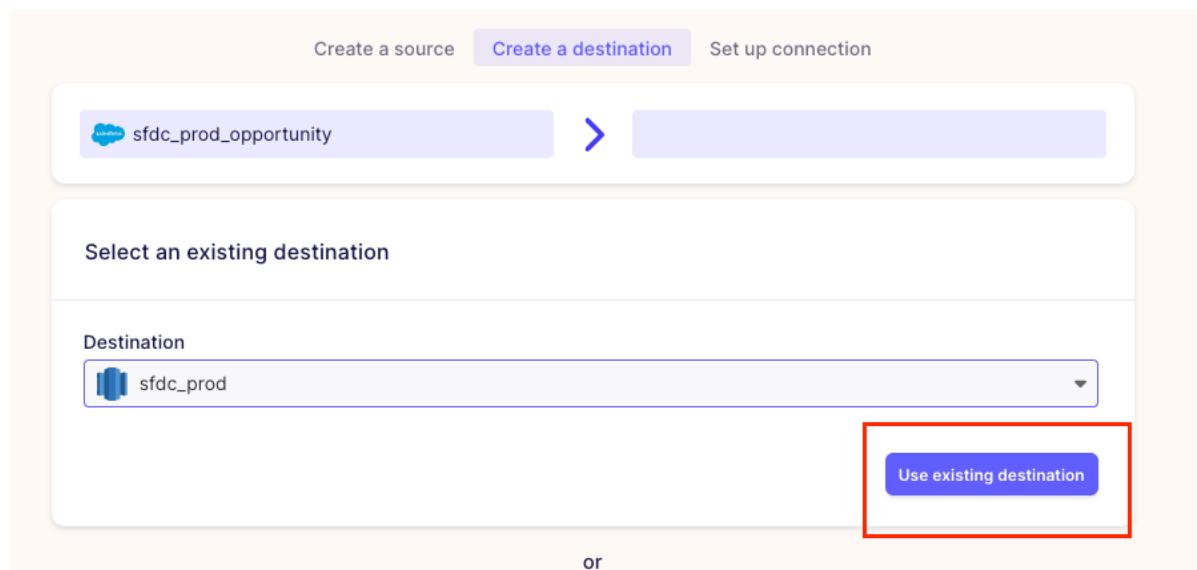
## 6. Airbyte connection: Salesforce - Redshift

So, in #3 we had created a source in Airbyte (a Salesforce source object which parses for "Opportunity" object), and in #5 we added a destination in Redshift. The "Connection" in Airbyte is what connects either of these.

Use the source created in #3,



Use the Destination in #5,



Finally, name of the connection, and schedule it. I have done a *cron* schedule at UTC 1 daily.

On this screen we can use if the loads needs be overwrite, incremental loads, etc. Typically I chose overwrites to begin with (so that full history is available), and then subsequent schedules are incremental dedupe and append loads.

The screenshot shows the 'Set up connection' page for a new connection named 'sfdc\_prod\_opportunity'. In the 'Transfer' section, the replication frequency is set to 'Cron' with the expression '0 0 1 \* \* ?' and the time zone is 'UTC'. In the 'Streams' section, the destination namespace is set to 'Mirror source structure' and the stream prefix is 'prefix'. The 'Activate the streams you want to sync' table lists two streams: 'No namespace' (Sync mode: Overwrite) and 'Opportunity' (Sync mode: Deduped + history). A 'Refresh source schema' button is visible in the top right of the table area.

| Sync                                                        | Source ⓘ     | Sync mode ⓘ          | Cursor field ⓘ                  | Primary key ⓘ   | Destination ⓘ                  |
|-------------------------------------------------------------|--------------|----------------------|---------------------------------|-----------------|--------------------------------|
| <input type="checkbox"/>                                    | Namespace    | Source   Destination |                                 |                 | Namespace Stream name          |
| <input type="checkbox"/> > <input checked="" type="radio"/> | No namespace | Describe             | Full refresh   Overwrite        | <source schema> | Describe                       |
| <input type="checkbox"/> > <input checked="" type="radio"/> | No namespace | Opportunity          | Incremental   Deduped + history | SystemModstamp  | Id <source schema> Opportunity |

Once connection is setup, sync the connection (it will eventually run on the configured schedule, this is just force running once manually, to see everything is fine)

## 7. Access the created/updated table in Redshift

Assuming a scheduled run of the connection is fine, we would expect the opportunity table in ab\_sfdc schema in Redshift. I use [DBeaver client](#) to connect to Redshift.

|    | abc_airbyte_unique_key           | abc_id              | abc_name   | abc_o__c          | abc_type     | iswon |
|----|----------------------------------|---------------------|------------|-------------------|--------------|-------|
| 1  | 0032f4140f9ca0b962ab7b966e9869d1 | 0060h000017WPbAAO   | [REDACTED] | Sam Phelan        | New Logo     | [ ]   |
| 2  | 0036bbfb4475bf62160f117da161280  | 0060h00001BxeRAAT   | [REDACTED] | Louis Frickman    | Paid POC     | [v]   |
| 3  | 003d5605c9e415f5fdbcc3886873d221 | 0060h00001E7BV7AAAN | [REDACTED] | Steve Tsang       | Upgrade      | [ ]   |
| 4  | 0040946bebc46143bb2e6b3efdb215   | 0060L00000mDpxQAW   | [REDACTED] | Ed Gibson         | New Logo     | [ ]   |
| 5  | 006cb1f16a6769e1a111b570832b3e9  | 0068a000010f0duAAJ  | [REDACTED] | Darsika Nanda     | Renewal      | [ ]   |
| 6  | 00196b6028c04075d32aaa373987c3   | 0060h00018YbaAAK    | [REDACTED] | Todd Hicks        | New Logo     | [ ]   |
| 7  | 007d4e080c3db331da9a483064300c07 | 0060h00019H8ccAAC   | [REDACTED] | Ted Saludes       | New Logo     | [ ]   |
| 8  | 008620b36ebe426ba468223d6ac902   | 0060L000001SEzQAM   | [REDACTED] | Daniel Payomi     | New Logo     | [ ]   |
| 9  | 008730a5c3e298948c0b9e9684ad8393 | 0060h000018YCHuAO   | [REDACTED] | Jim Hoppe         | New Logo     | [ ]   |
| 10 | 00aa67bcf490d4eb8e6654c7a8110df  | 0060h000018YCaQAW   | [REDACTED] | Sam Phelan        | New Logo     | [ ]   |
| 11 | 002f229269ea6404db142af7602ffe9  | 0060h00001E72h4AAB  | [REDACTED] | Hiroaki Kawamura  | Upgrade      | [v]   |
| 12 | 0000e0d751469c4e7c678e2744fe7    | 0060h00001AZDBVA5   | [REDACTED] | Gary Newgaard     | New Logo     | [ ]   |
| 13 | 005f4846309e70d2a6390d07e2cb2    | 0060L00000pcjvNAQ   | [REDACTED] | Chris Bowman      | New Logo     | [ ]   |
| 14 | 00640c6eeb0d8749a0121f931f708b5b | 0060h00001CV3hRAAT  | [REDACTED] | Peter Grammaticas | New Logo     | [ ]   |
| 15 | 00757c9d718b67283f60a21263ed179  | 0060h00001AXpGfAAI  | [REDACTED] | Erik Ruby         | New Logo     | [ ]   |
| 16 | 0039ae60fa914d529f85e1fc9de57bc  | 0060L00000pcwhYQAQ  | [REDACTED] | Paul Wilcox       | New Logo     | [ ]   |
| 17 | 007641717ea6991e2397dc8676fc5e37 | 006E000000fh565IA   | [REDACTED] | Lindsey Liranzo   | SDR          | [ ]   |
| 18 | 0068f1470aed9767a2e83d4570be4e56 | 006E000000SqzIA     | [REDACTED] | Paul Wilcox       | Legacy Dupli | [ ]   |
| 19 | 003fed90eb501f448c660250dd1590   | 006E000000gbybmIAI  | [REDACTED] | Ed Gibson         | New Logo     | [v]   |
| 20 | 00815c45f99ddc242b6114c65254ff0  | 0060L00000mAdkbQAC  | [REDACTED] | Jim Hoppe         | Upgrade      | [v]   |
| 21 | 008c4759fcbedafaca6e41e59a50b8c3 | 006E000000yENPIA4   | [REDACTED] | Lindsey Liranzo   | SDR          | [ ]   |

## 8. Redshift transforms

So we have the raw/ source sfdc opportunity object loaded into `ab_sfdc.opportunity`. Opportunity in real world is a business expansion entity, it can relate to a team or a department in a company.

Requirement:

- trend of business opportunities dollar value at various slices as type, geo, industry segment, etc.
- weekly trend generated on 1<sup>st</sup>, 8<sup>th</sup>, 15, 22<sup>nd</sup> and 29<sup>th</sup> day of each month (and not Mondays/ or Sundays)
- we have a plan / forecast/ estimate for the dollar values, and would like to have a comparative view of the metrics

To implement these requirements,

- we keep a daily snapshot of the `opportunity` table , (let's say its pipeline )
- use a `schedule` table which has list of dates i.e., 1,8,15,22 and 29
- use a `plan_closed_won` table for forecast/ plan of leads.

The query to produce weekly trend of opportunity dollar values is [https://github.com/bpjena/ms-dsc/blob/master/dsc680\\_project/redshift/sfdc\\_opportunity\\_trend.sql](https://github.com/bpjena/ms-dsc/blob/master/dsc680_project/redshift/sfdc_opportunity_trend.sql).

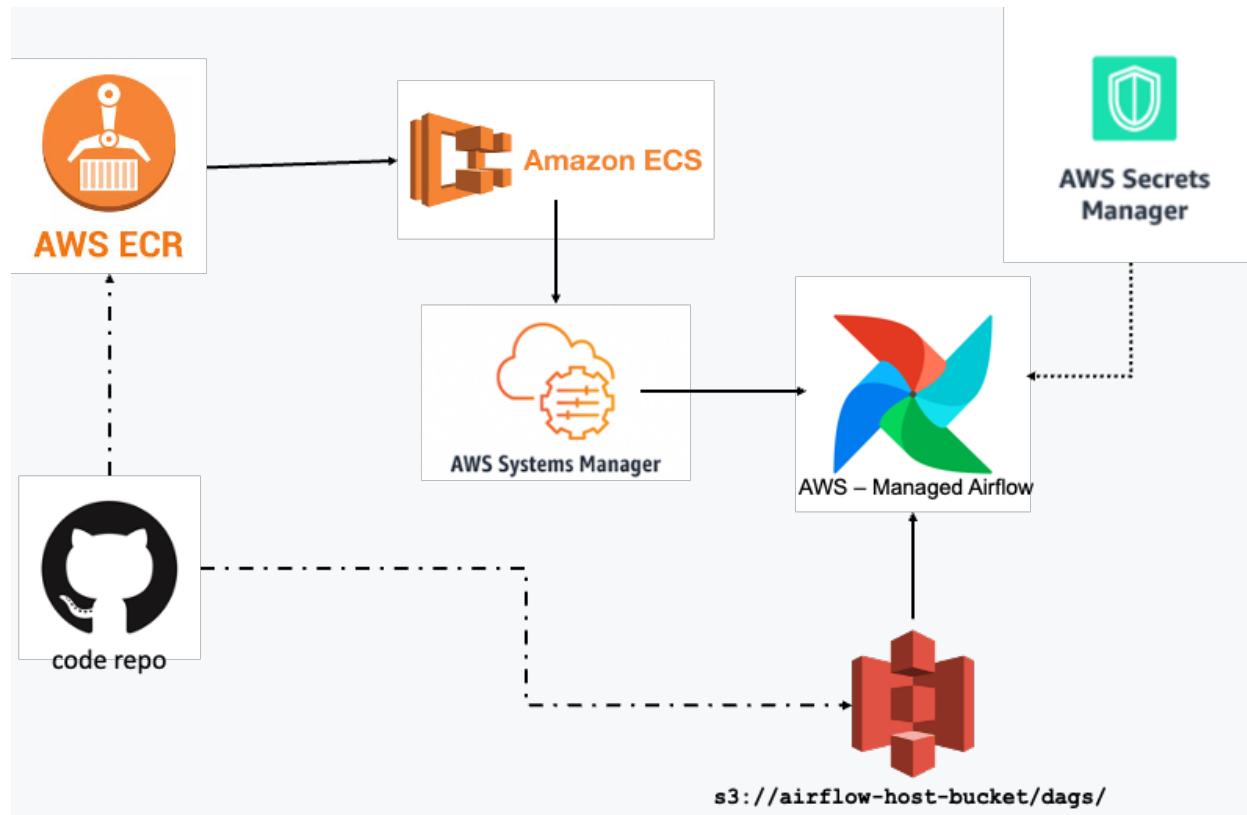
With the current week being represented as `ref_week = 0`, final dataset produces a futuristic 52-week view and similarly a historical 52-week view of the opportunity dollar values. Simply executing the query produces a dataset using DBeaver. However, to make this a data pipeline, will need to automate this dataset creation.

## 9. Airflow jobs

The intention is to be able to automate the redshift transformation activity through a scheduler. While I could have chosen a simple cron scheduler on any EC2 machine as a working solution, I choose to perform a mature deployment of open-source orchestration platform like [airflow](#) to not only be able to perform this transform automation. But also use airflow features for any data related automation like ingestion jobs, checking data quality of datasets, etc

I will be using [python redshift connector](#) as the core library to script a python job which will be scheduled using airflow. Code and related instructions to run this local machine is detailed [https://github.com/bpjena/ms-dsc/tree/master/dsc680\\_project](https://github.com/bpjena/ms-dsc/tree/master/dsc680_project). Core idea is to run code in Docker containers, and ensure the replication of container when code runs by Airflow. This method has inherent scaling benefits, since scheduling and code execution environments are decoupled and each can scale as needed.

While there are many custom Airflow installations in use, due to the open-source nature of this utility, I've chosen [Managed Airflow](#) from AWS as my airflow version of use, since a lot of airflow components and hence its maintenance/ infra is abstracted through this managed service from AWS. The code detailed here is for executing the python jobs in local docker containers. When airflow runs these, the local docker container is replicated by ECS Operator usage from Airflow (elastic container service). Here is how the deployment of the code looks like from Airflow job flow aspect.



In order to initialize the Airflow instance, refer the steps  
<https://docs.aws.amazon.com/mwaa/latest/userguide/get-started.html>

Code push/ deployment as an activity is about pushing the code into AWS ECR (elastic container registry) and copying the dag file into S3 bucket (dags folder in S3 maps to the Airflow jobs or DAGs).

## 10. Tableau Setup

Tableau Redshift connection credentials has read privilege on the transformed dataset created from the job. Create separate credentials to be used for Tableau connection to Redshift.

Views 1    Data Sources 2    Connected Metrics 0    Custom Views 0    Subscriptions 0

Select All

The screenshot shows the Tableau Data Sources page. At the top, there are counts for Views (1), Data Sources (2), Connected Metrics (0), Custom Views (0), and Subscriptions (0). Below this, a "Select All" button is visible. The main area displays two data sources, each represented by a cylinder icon and labeled "Live connection". The first source is "ODS\_SFDC\_ACCOUNT (SALESFOR..." and the second is "ODS\_SFDC\_OPPORTUNITY\_ACCTO...". Both sources have three dots at the end, indicating more options.

Here is a screenshot of the graphical plots for the transformed dataset when rendered through Tableau.



## 11. Local python code run / setup

Please refer the instructions in the [https://github.com/bpjena/ms-dsc/tree/master/dsc680\\_project/README.md](https://github.com/bpjena/ms-dsc/tree/master/dsc680_project/README.md) for local setup. To run the python code locally will need .env file within which has the Redshift access credentials or any security params that we want to pass as environment variables. AWS credentials can also be stored here in the .env file for local code runs and tests. Commands to run the job in local machine is detailed in .local file. So the job run in local docker container has the following execution logs.

```
[bjena@bjena-mac deng-jobs % make env-up
docker build -t deng-jobs:local .
[*] Building 7.4s (17/17) FINISHED
=> [internal] load build definition from Dockerfile
=> => transferring Dockerfile: 37B
=> [internal] load dockerignore
=> => transferring dockerignore: 0B
=> [internal] load metadata for docker.io/library/python:3.10-slim
=> [internal] load build context
=> => transferring context: 514.89K
=> [base 1/12] FROM docker.io/library/python:3.10-slim@sha256:13b98a4ddf43c98e3f91bd5136844ff8b0c1d2a4c913b4a430d552a43d4c7b3b
=> CACHED [base 2/12] RUN groupadd --gid 10001 deng && useradd --create-home --no-user-group --gid 10001 --uid 10001 deng
=> CACHED [base 3/12] WORKDIR /app
=> CACHED [base 4/12] COPY pyproject.toml poetry.lock ./app/
=> CACHED [base 5/12] RUN apt-get update -yqq && apt-get upgrade -yqq && apt-get install -yqq --no-install-recommends freetds-dev libkrb5-dev libsasl2-dev libssl-dev libffi-dev libpq-dev git && apt-get autoremove -yqq --purge 0.0s
=> CACHED [base 6/12] RUN curl -sSL https://raw.githubusercontent.com/python-poetry/poetry/master/install-poetry.py | python -
=> CACHED [base 7/12] RUN poetry config virtualenvs.create false
=> CACHED [base 8/12] RUN poetry install --only main --no-root
=> CACHED [base 9/12] RUN apt-get purge --auto-remove -yqq freetds-dev libkrb5-dev libsasl2-dev libssl-dev libffi-dev libpq-dev git && apt-get autoremove -yqq --purge 0.0s
=> [base 10/12] COPY . ./app
=> [base 11/12] RUN poetry install --only main
=> [base 12/12] RUN chown -R deng:deng ./home/deng /app
=> exporting to image
=> exporting layers
=> writing image sha256:67aa9388bcd7db2565b74be7e49f5cb7ece7227f51accd86ed130155b0275947
=> naming to docker.io/library/deng-jobs:local
Use 'docker scan' to run Snyk tests against images and learn how to fix them
docker run -it --rm -u deng -v /Users/bjena/.aws:/home/deng/.aws --env-file .env -w /app deng-jobs:local /bin/bash
deng@8ab5412db285:/app$ /bin/bash -c "PYTHONPATH=. python python/redshift/sfdc_open_pipe.py"
2022-11-11 06:42:49.323 [INFO] (sfdc_open_pipe.py:64): INVALID schedule...
2022-11-11 06:42:49.324 [INFO] (sfdc_open_pipe.py:118): schedule skipped 2022-11-11...
2022-11-11 06:42:49.325 [INFO] (sfdc_open_pipe.py:119): All Done ;-)
deng@8ab5412db285:/app$
```

```
[deng@8ab5412db285:/app$ /bin/bash -c "PYTHONPATH=. python python/redshift/sfdc_transforms.py"
2022-11-11 06:43:56,953 [INFO] (sfdc_transforms.py:46): starting.. sfdc.account
2022-11-11 06:43:57,509 [INFO] (sfdc_transforms.py:51): cleaned.. sfdc.account [truncate]
2022-11-11 06:44:40,153 [INFO] (sfdc_transforms.py:57): inserted.. sfdc.account
2022-11-11 06:44:40,153 [INFO] (sfdc_transforms.py:46): starting.. sfdc.opportunity
2022-11-11 06:44:40,875 [INFO] (sfdc_transforms.py:51): cleaned.. sfdc.opportunity [truncate]
2022-11-11 06:45:13,401 [INFO] (sfdc_transforms.py:57): inserted.. sfdc.opportunity
2022-11-11 06:45:13,401 [INFO] (sfdc_transforms.py:46): starting.. sfdc.pipeline
2022-11-11 06:45:14,040 [INFO] (sfdc_transforms.py:51): cleaned.. sfdc.pipeline [daily_snapshot]
2022-11-11 06:45:58,198 [INFO] (sfdc_transforms.py:57): inserted.. sfdc.pipeline
2022-11-11 06:45:58,198 [INFO] (sfdc_transforms.py:46): starting.. sfdc.campaign
2022-11-11 06:45:58,836 [INFO] (sfdc_transforms.py:51): cleaned.. sfdc.campaign [truncate]
2022-11-11 06:46:01,513 [INFO] (sfdc_transforms.py:57): inserted.. sfdc.campaign
2022-11-11 06:46:01,514 [INFO] (sfdc_transforms.py:46): starting.. sfdc.contact
2022-11-11 06:46:02,747 [INFO] (sfdc_transforms.py:51): cleaned.. sfdc.contact [truncate]
2022-11-11 06:49:13,383 [INFO] (sfdc_transforms.py:57): inserted.. sfdc.contact
2022-11-11 06:49:13,384 [INFO] (sfdc_transforms.py:46): starting.. sfdc.lead_share
2022-11-11 06:49:14,322 [INFO] (sfdc_transforms.py:51): cleaned.. sfdc.lead_share [truncate]
2022-11-11 06:49:15,727 [INFO] (sfdc_transforms.py:57): inserted.. sfdc.lead_share
2022-11-11 06:49:15,729 [INFO] (sfdc_transforms.py:46): starting.. sfdc.lead_status
2022-11-11 06:49:16,331 [INFO] (sfdc_transforms.py:51): cleaned.. sfdc.lead_status [truncate]
2022-11-11 06:49:16,865 [INFO] (sfdc_transforms.py:57): inserted.. sfdc.lead_status
2022-11-11 06:49:16,873 [INFO] (sfdc_transforms.py:58): done.
2022-11-11 06:49:16,873 [INFO] (sfdc_transforms.py:63): All Done ;-)
deng@8ab5412db285:/app$
```