# StreamSets Data Collector 5.7.1：Dockerの展開とWebSocketパイプラインの作成

本ガイドでは、StreamSets Data Collector 5.7.1をDockerコンテナに展開する手順、およびWebSocketに接続して受信したデータをローカルに保存するパイプラインの作成方法を説明します。

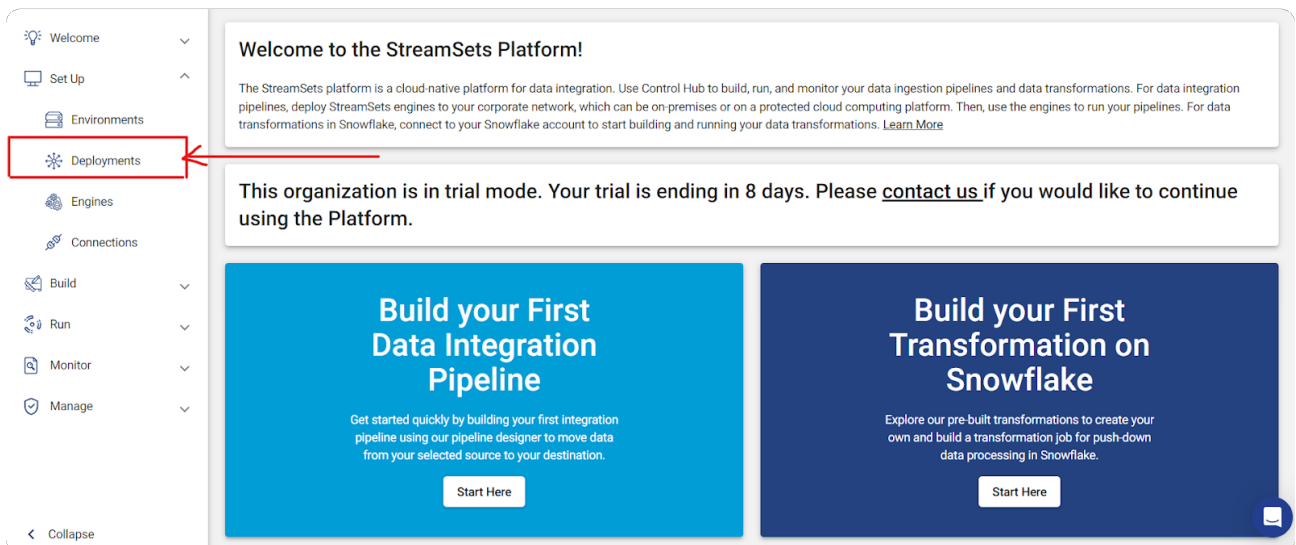## 前提条件：

- お使いのマシンにDockerがインストールされていること（Dockerをダウンロードしてインストールしてください）。
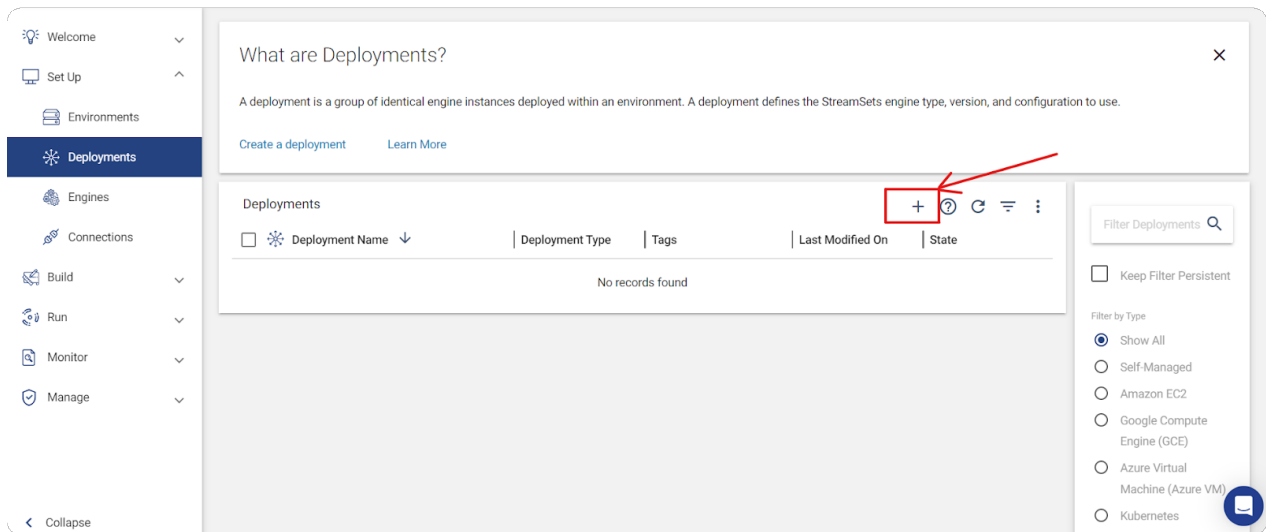
- StreamSetsのアカウント。

## StreamSets Data Collectorの展開

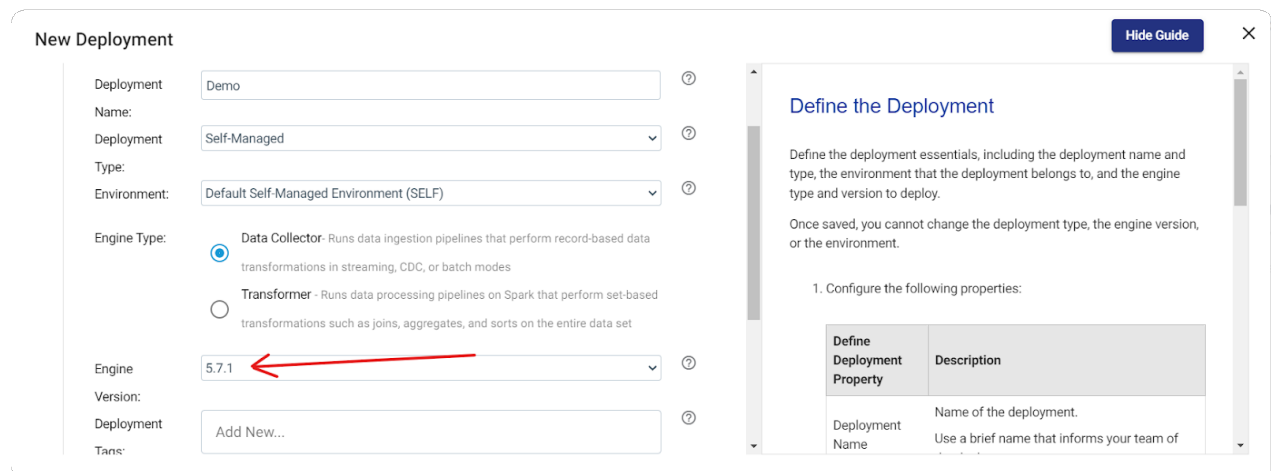### ステップ1： Data CollectorのDeployment設定
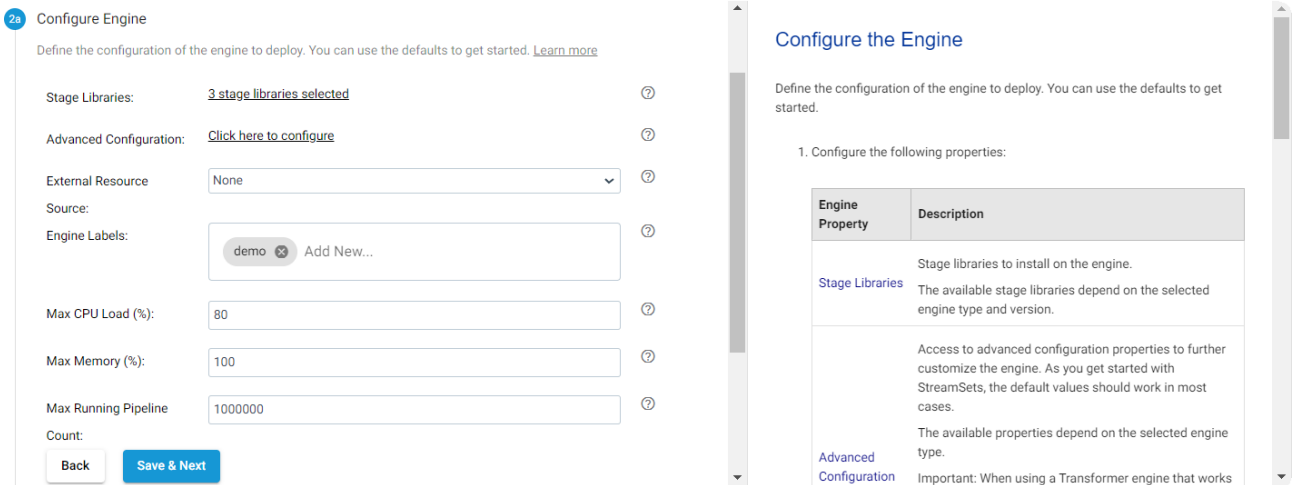
- ログイン後、サイドバーで「Deployment」を選択します。

- 以下のように「+」をクリックして新しい展開（Deployment）を開始します。



- 「Engine」で現在の安定版リリースであるバージョン 5.7.1 を選択し、ご使用の要件に従って定義された展開を完了します。「Save & Next」ボタンをクリックします。

- エンジンの設定をお使いの仕様に合わせ、「Save & Next」ボタンを選択します。



- インストールの種類として「Docker Image」を選択し、「Save & Next」ボタンをクリックします。

- 展開の共有設定をカスタマイズして他のユーザーやグループにアクセスを許可します。「Save & Next」ボタンをクリックします。



- 「Start & Generate Install Script 」ボタンをクリックします。

- Dockerが起動していることを確認してから、以下のコマンドをコピーしてください。このコマンドをwindows/Ubuntu/macのターミナルに貼り付けると、エンジンが起動します。

- コンテナがアクティブになっていることを確認します。以下のスクリーンショットの通り、Docker Desktop内で実行中のコンテナを確認してください。

## ステップ2： Data Collectorの詳細設定

WhoisXML APIのWebSocketを活用するために、データコレクターの設定をカスタマイズしてください。

- ビルドされたjarファイルはこちらからダウンロードできます。ダウンロード後、Dockerコンテナ内の既存のファイルを新しく取得したファイルで置き換えます。ファイルのシームレスな置き換えには、Docker Desktopが効率的で便利です。
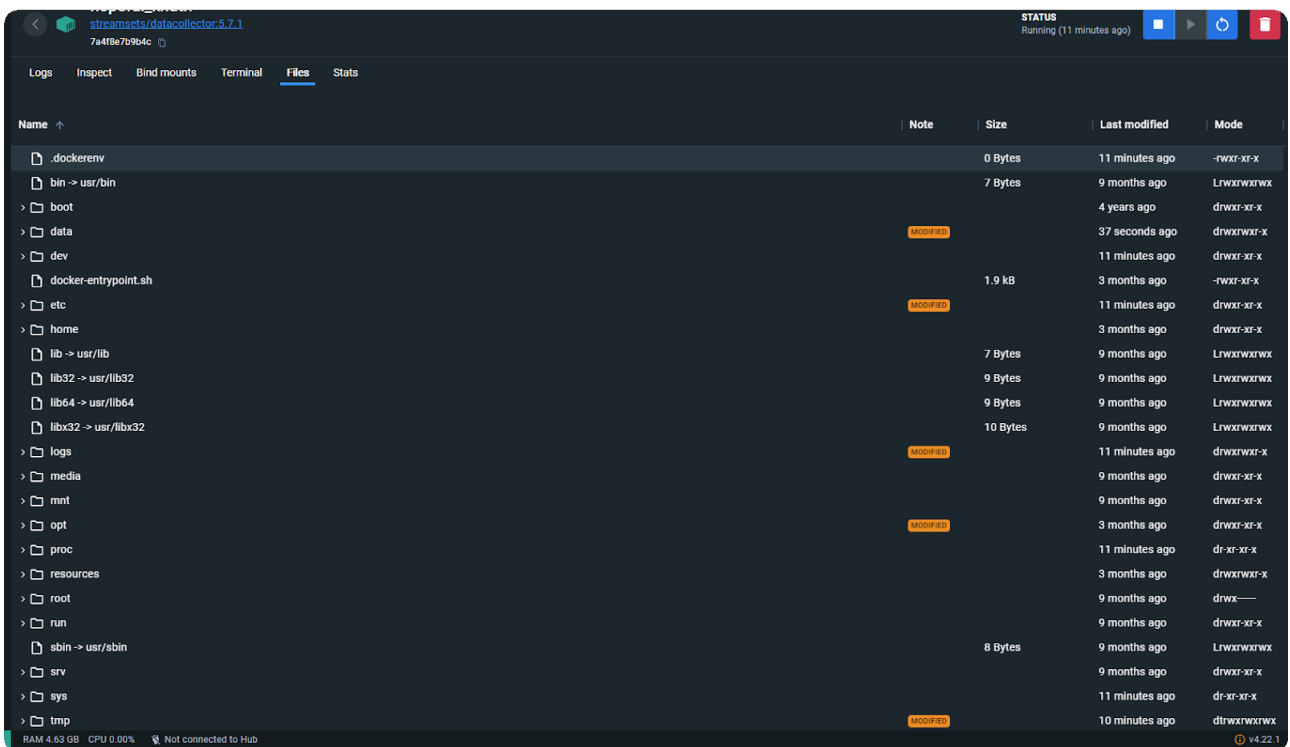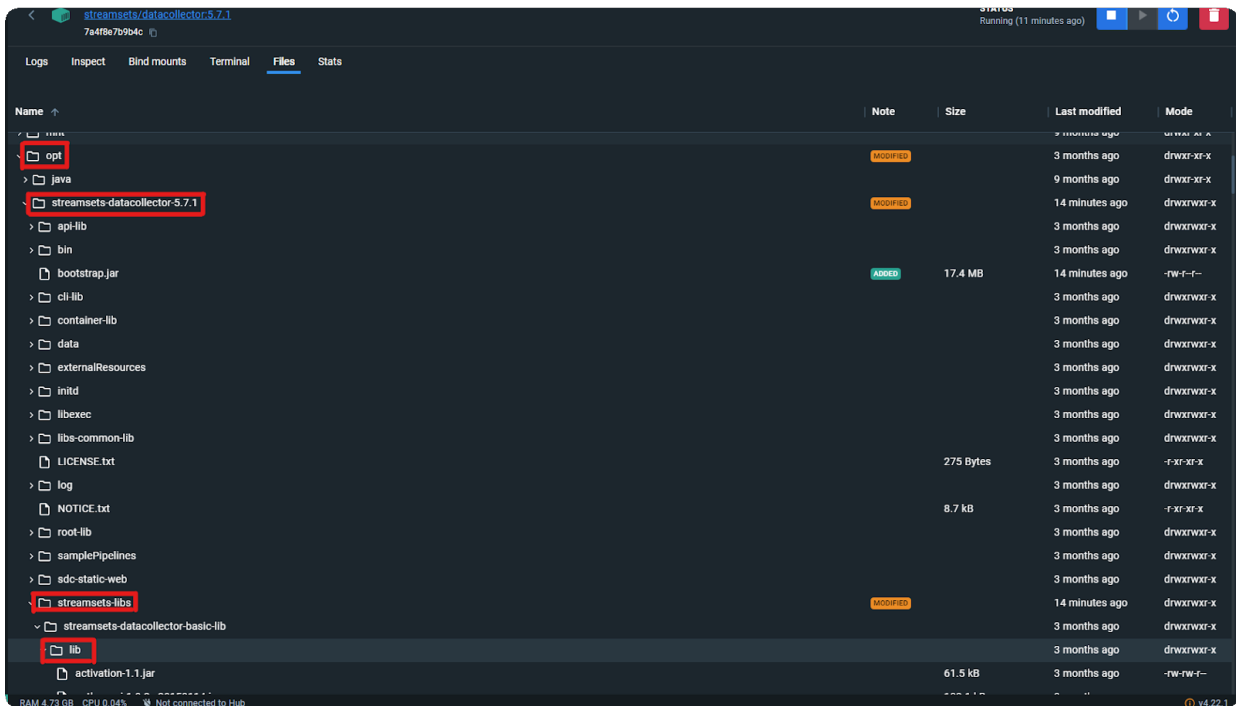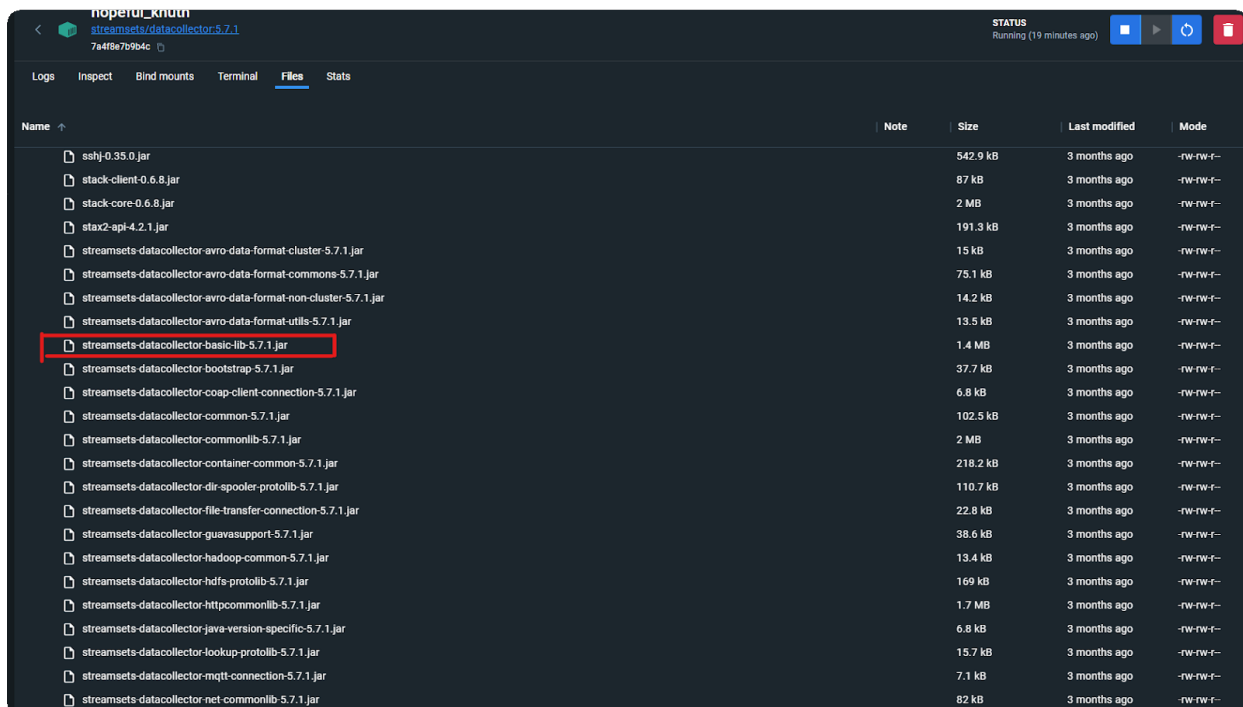


- ディレクトリ /opt/streamsets-datacollector-5.7.1/streamsets-libs/streamsets-datacollector-basic-lib/lib に移動します。

- このディレクトリに、「streamsets-datacollector-basic-lib-5.7.1.jar」というjarファイルがあるはずです。



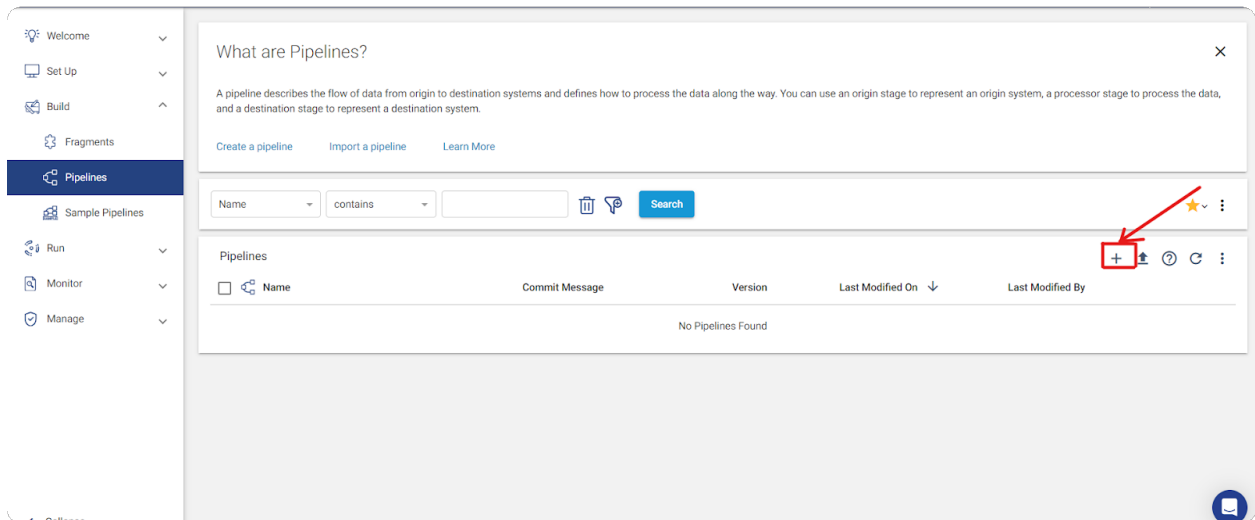- 既存のファイルをダウンロードしたファイルで置き換えます。または、ダウンロードしたファイルをこの場所にドラッグすれば、簡単に直接置き換えることができます。

- Data Collector/Containerを再起動します。

# パイプラインの作成

## ステップ3：パイプラインのセットアップ

- StreamSets UIでサイドバーから「Build」を選択し、次に「Pipelines」を選択します。以下のように「＋」ボタンをクリックしてパイプラインを開始します。



- 新しいパイプラインをお使いの仕様に合わせてカスタマイズし、「Next」ボタンをクリックして進みます。

- パイプラインの設定を調整して指定のデータコレクタを選択し、「Save & Open in Canvas」ボタンをクリックします。



- 以下のようなユーザーインターフェース（UI）が表示されます。



- 「Add Stage」ボタンをクリックし、「WebSocket」を検索して 「WebSocket Client」を選択します。

WebSocketステージを選択し、適宜設定します：

## WebSocketの設定

- Resource URLを入力

- Request Dataを入力（WHOIS APIキーが入ります）

- Max Message Length (bytes) として最小522184 を入力

**Data Formatの設定**

- Data FormatとしてJSONを選択

- Max Object Length (chars) として9999999 を入力（必要に応じて変更できます）



- UIの「Add Stage」ボタンでステージを追加し、「Local FS」ステージを選択します。

## Local FSの設定

「Local FS」ステージを選択し、要件に従って設定します。

- 希望する出力ファイルの場所を「Directory Template」に入力します。



要件に応じて「Data Format」に必要な設定を入力します。

- 「Data Format」としてJSONを選択します。



設定後、「Validate」ボタンをクリックしてパイプラインを検証し、エラーを特定して修正しま

す。パイプラインの最終状態は、以下の例のようになります。

# 最終ステップ



# パイプラインの実行

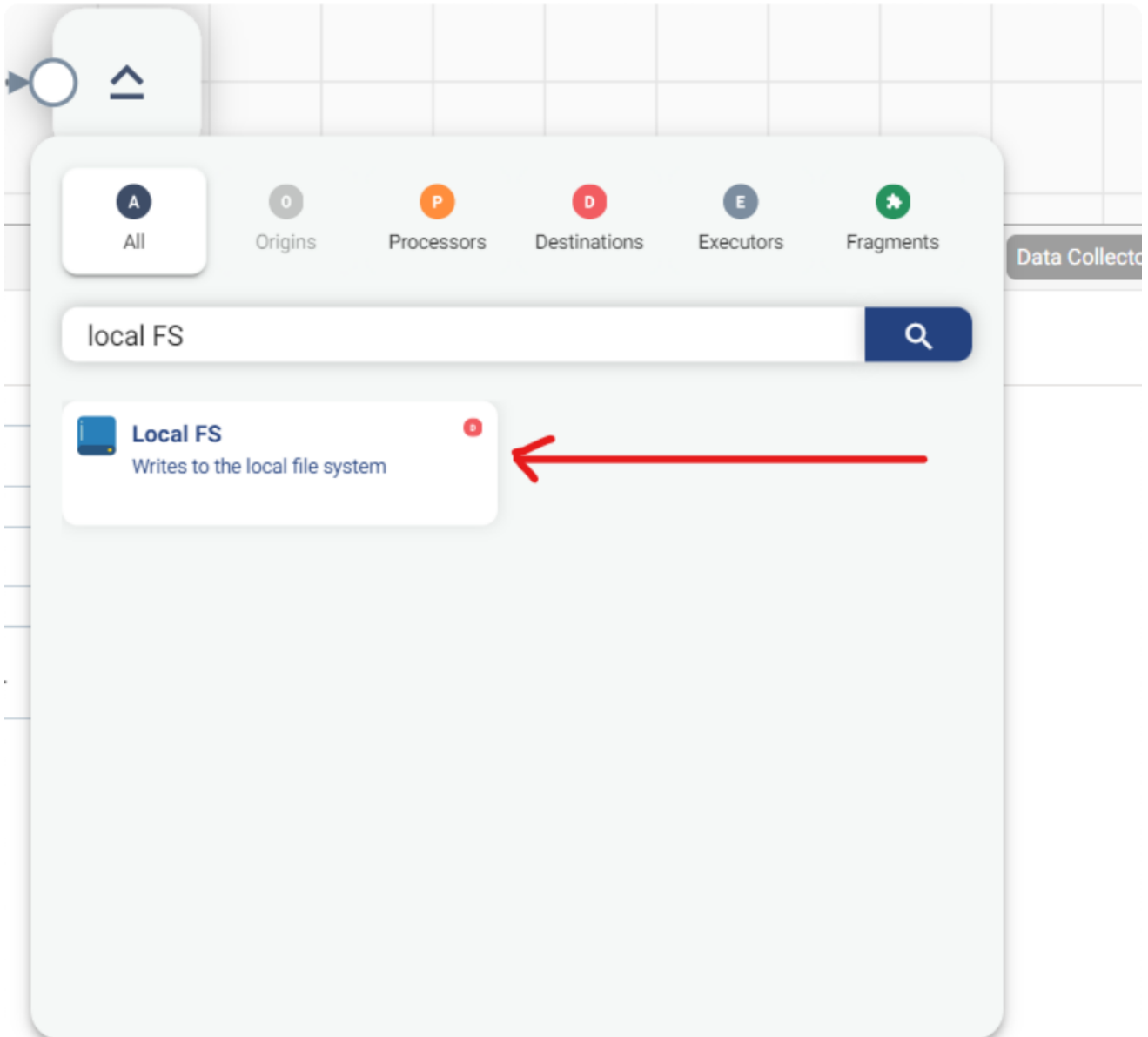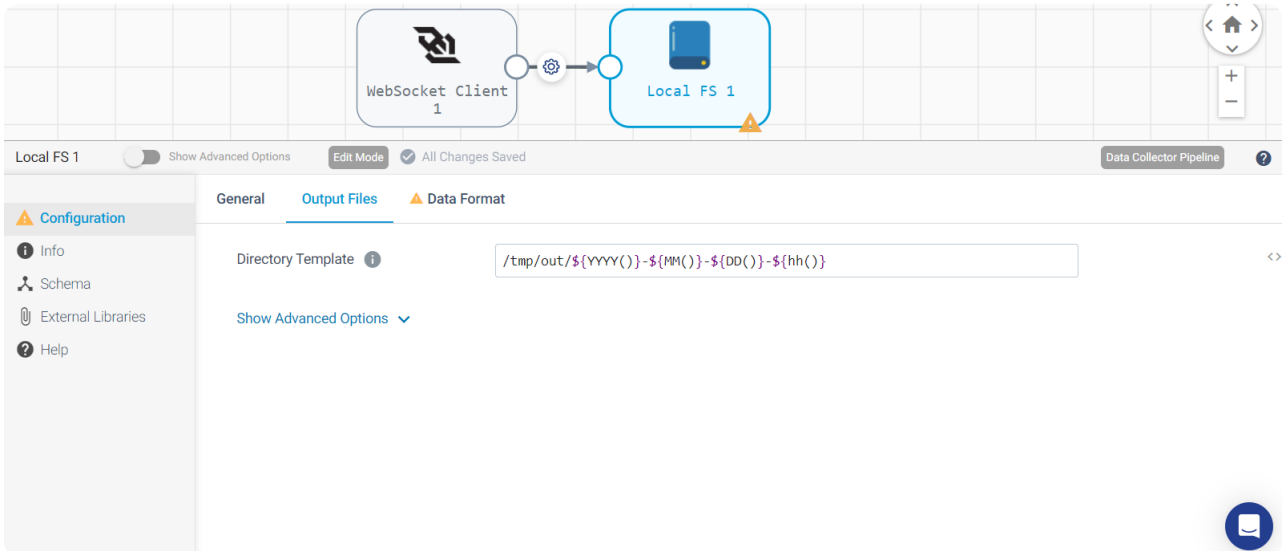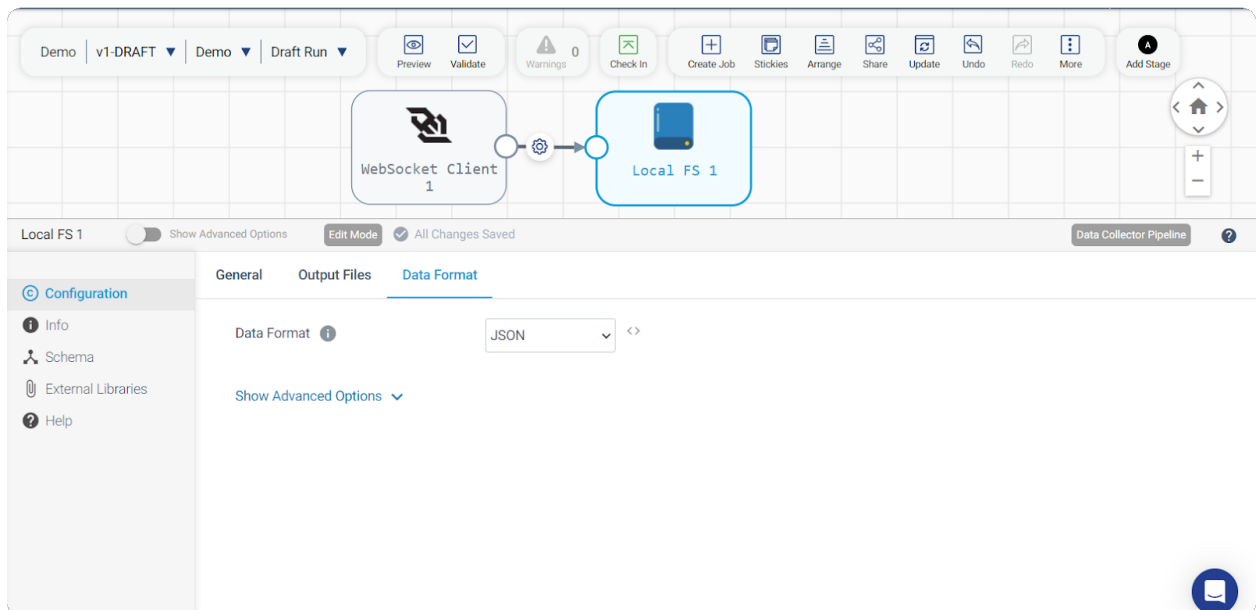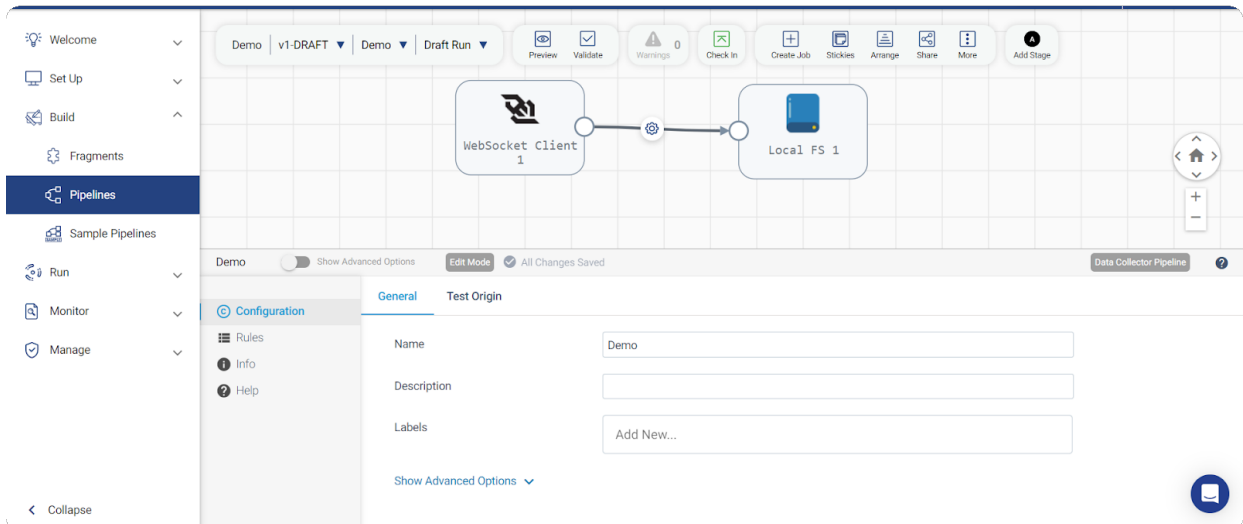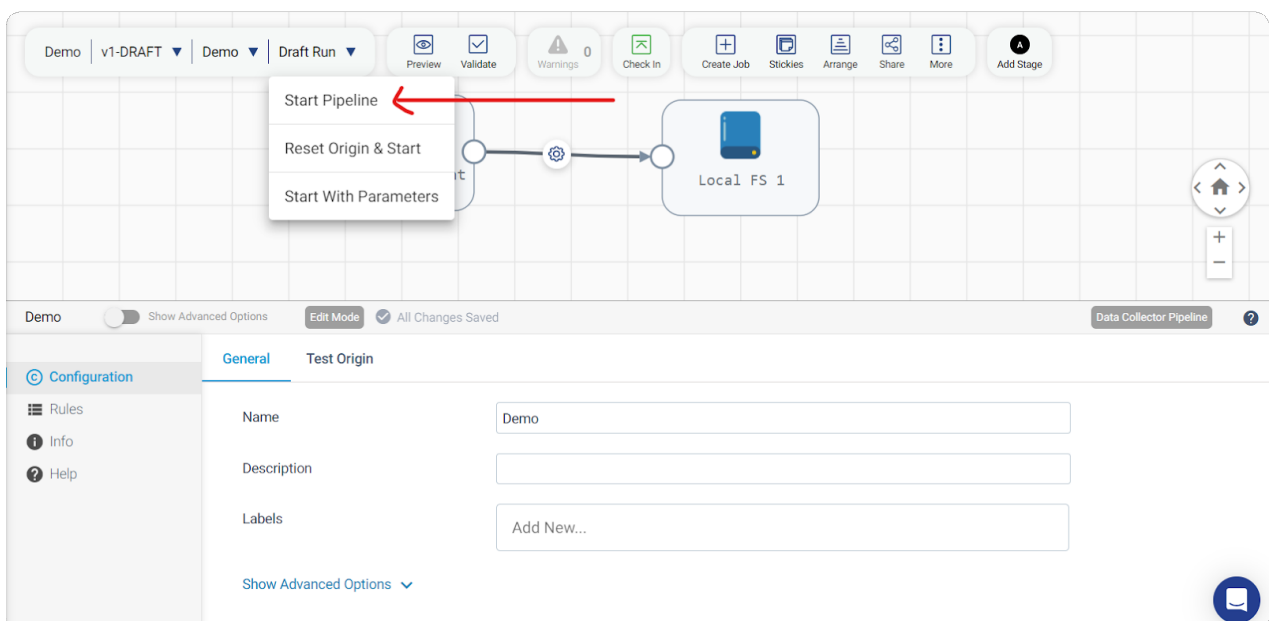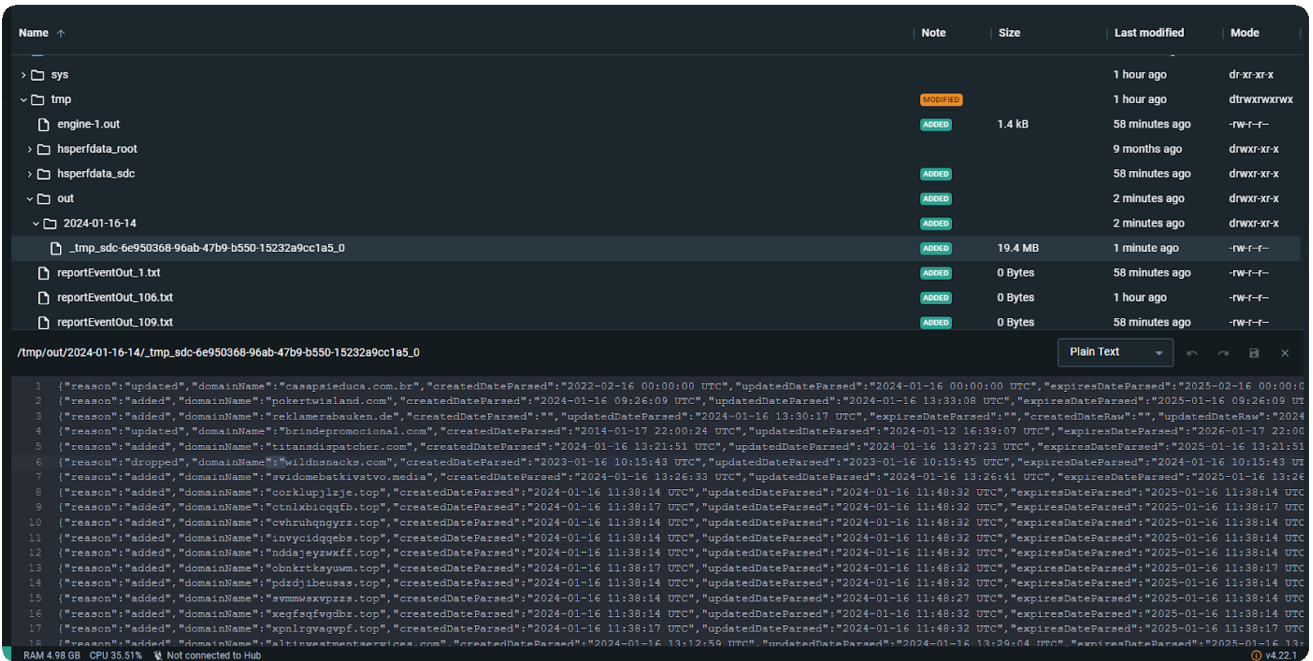UIで「Draft Run」、「Start Pipeline」を選択してパイプラインを実行します。



上記の手順が完了した後、パイプラインを起動すると、下図のようなUIが表示されます。

Docker コンテナ内では、WebSocket から取得したデータを含むファイルが指定したディレクトリに作成されていることを確認できます。

# まとめ

本ガイドでは、Dockerを使用してStreamSets UIでパイプラインをセットアップおよび実行する手順を概説しました。WebSocket Clientの設定からLocal FSステージの定義に至る各ステップを踏むことで、最終的にデータ処理のパイプラインを構築できます。検証ステップはパイプラインの整合性を保証するもので、実行に成功すると、UIがDockerコンテナ内の指定されたディレクトリ内に出力ファイルを表示します。本ガイドで示した手順により、合理化されたデータ処理パイプラインの作成、設定、実行を正常に行い、効果的なデータ統合および管理を実現できます。