

# FTP Server Import Integration

Use this data connector to directly import data from your FTP server to Treasure Data.

For sample workflows on importing data from your FTP server, view [Treasure Boxes](#).

- [Prerequisites](#)
- [Requirements](#)
- [About Incremental Data Loading](#)
  - [Limitations, Supported, Suggestions](#)
  - [About Incremental Loading for Integrations](#)
  - [Incremental Loading for Connectors](#)
- [Use the TD Console to Create Your Connection](#)
  - [Create a New Connection](#)
  - [Transfer Your Data to Treasure Data](#)
  - [Filters](#)
  - [Data Preview](#)
  - [Data Placement](#)
- [Validate Connection](#)
- [Optionally Configure Export Results in Workflow](#)

## Prerequisites

- Basic knowledge of Treasure Data
- Basic knowledge of FTP

## Requirements

- Make sure you have a valid protocol. If you intend to *FTP* or *FTPS*, you can use this Data Connector for FTP. If *SFTP*, use the [SFTP Integration](#).
- If you're using a firewall, check your accepted IP range/port. Server administrators sometimes change the default port number for security reasons.
- Be sure that FTP uses *TCP/21* as the default control port but also uses any TCP ports as a data transfer port when you're using passive mode. This port range will depend on your server's settings.
- Check that you're connecting with *passive* mode. *activeP* mode generally doesn't work because it establishes the connection from the FTP server-side.
- If you're using FTPS, there are 2 modes *Explicit* and *Implicit*. Explicit mode is typically used.

## About Incremental Data Loading

Incremental loading is the activity of loading only new or updated records from a source into Treasure Data. Incremental loads are useful because they run efficiently when compared to full loads, and particularly for large data sets.

Incremental loading is available for many of the Treasure Data integrations. In some cases, it is a simple checkbox choice and in others, after you select incremental loading you are provided with other fields that must be specified.

## Limitations, Supported, Suggestions

- For some integrations, if you choose incremental loading, you might need to make sure that there is an index on the columns to avoid a full table scan.
- Only Timestamp, Datetime, and numerical columns are supported as `incremental_columns`.
- For the raw query, the `incremental_columns` is required because it won't be able to detect the Primary keys for a complex query.

## About Incremental Loading for Integrations

Treasure Data Incremental loading has 4 patterns (3 types of data connector + 1 workflow `td_load` operator.), then the 3 data connector loading examples are as follows:

- Cloud storage service (e.g. AWS S3, GCS and etc.)
  - Lexicographic order of file name
- Query (e.g. MySQL, BigQuery and etc.)
  - Date time
- Variable period (Google Analytics, etc)

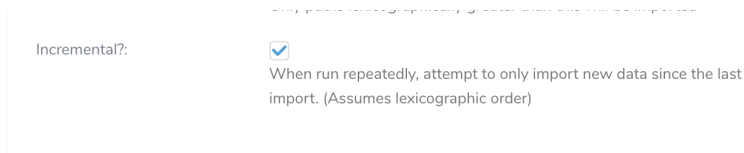
- Use `start_date` for loading

## Incremental Loading for Connectors

If incremental loading is selected, data for the connector is loaded incrementally.

This mode is useful when you want to fetch just the object targets that have changed since the previously scheduled run.

For example, in the UI:



The screenshot shows a UI element with the label 'Incremental?:' followed by a checked checkbox. Below the checkbox is the text: 'When run repeatedly, attempt to only import new data since the last import. (Assumes lexicographic order)'. The checkbox is a small square with a blue checkmark inside.

Database integrations, such as MySQL, BigQuery, and SQL server, require column or field names to load incremental data. For example:



The screenshot shows two UI elements. The first is 'Incremental Loading:' with a checked checkbox and the text 'When enabled, attempt to ingest only new data from the last import'. The second is 'Incremental Field:' with an empty text input box and the text 'The timestamp column to compare with last ingestion' below it.

Learn more [About Database-based Integrations](#).

## Use the TD Console to Create Your Connection

### Create a New Connection

In Treasure Data, you must create and configure the data connection prior to running your query. As part of the data connection, you provide authentication to access the integration.

1. Open **TD Console**.
2. Navigate to **Integrations Hub > Catalog**.
3. Search for and select FTP.



4. Select **Create Authentication**.

## New Authentication

FTP



1 Credentials > 2 Details

Host:

Port:

User:

Password:

Passive mode?:

ASCII mode?:

Use FTPS/FTPES:

[Learn more](#) [Continue](#)

5. Enter the required credentials for your remote FTP instance. Depending on your selections, the fields you see might vary:

Field	Description
Host	The host information of the remote FTP instance, for example, an IP address.
Port	The connection port on the remote FTP instance the default is 21.
User	The user name used to connect to the remote FTP instance.
Password	The password used to connect to the remote FTP instance.
Passive mode	Use passive mode (default: checked)
ASCII mode	Use ASCII mode instead of binary mode (boolean, default: unchecked)
Use FTPS /FTPES	Use FTPS (SSL encryption). (boolean, default: unchecked)
Verify cert	Verify the certification provided by the server. By default, the connection fails if the server certificate is not signed by one of the CAs in JVM's default trusted CA list.
Verify hostname	Verify server's hostname matches the provided certificate.
Enable FTPES	FTPES is a security extension to FTPS
SSL CA Cert Content	Paste the contents of the certificate file

6. Select **Continue**.

7. Enter a name for your connection.

8. Select **Continue**.

## Transfer Your Data to Treasure Data

After creating the authenticated connection, you are automatically taken to Authentications.

1. Search for the connection you created.
2. Select **New Source**.
3. Type a name for your **Source** in the Data Transfer field.
4. Select **Next**.

The Source Table dialog opens.

1 Connection	Path prefix:	<input type="text"/>
2 Source Table	Path regex:	<input type="text"/>
3 Data Settings		Only files matching this regex will be included
4 Filters	Incremental?:	<input checked="" type="checkbox"/>
5 Data Preview		When run repeatedly, attempt to only import new data since the last import
6 Data Placement	Start after path:	<input type="text"/>
		Only paths lexicographically greater than this will be imported

5. Edit the following parameters:

Parameters	Description
Path prefix	The prefix of target files (string, required). For example, resultoutputtest.
Path regex	Type a regular expression to query file paths. If a file path doesn't match the specified pattern, the file is skipped. For example, if you specify the pattern <code>.csv\$ #</code> , then a file is skipped if its path doesn't match the pattern.
Incremental	Enables incremental loading (boolean, optional. default: true. If incremental loading is enabled, the config diff for the next execution will include <code>last_path</code> parameter so that the next execution skips files before the path. Otherwise, <code>last_path</code> is not included.
Start after path	Only paths lexicographically greater than this will be imported.

6. Select **Next**.

The Data Settings page can be modified for your needs or you can skip the page.

Optionally, you can modify data settings and then see your changes in Data Preview. [Skip This Step](#)

▼ DECODERS

[Add](#)

Type: Bzip2 ⊗

- Bzip2
- Gzip

▼ PARSER

**Create Source**  
Using meg\_ftpbetter

Optionally, you can modify data settings and then see your changes in Data Preview. [Skip This Step](#)

- 1 Connection
- 2 Source Table
- 3 Data Settings**
- 4 Filters
- 5 Data Preview
- 6 Data Placement

▼ DECODERS

[Add](#)

▼ PARSER

Type:

Delimiter:   
Delimiter character such as , for CSV, "\t" for TSV, "|" or any single-byte character

Quote character:   
The character surrounding a quoted value. Setting null disables quoting.

Escape character:   
Escape character to escape a special character. Set to null to disable escaping.

Cancel [Back](#) [Next](#)

Skip header lines:   
Skip this number of lines first. Set to 1 if the file has a header line.

Null string:   
If a value is this string, converts it to NULL. For example, set \N for CSV files created by mysqldump

Trim if not quoted?:   
Remove spaces of a value if the value is not surrounded by the quote character

Comment line prefix:   
Skip a line if the line begins with this string

Allow optional columns?:   
If true, sets omitted columns to null. If false, skips rows with insufficient columns.

Allow extra columns?:   
If true, ignores extra columns. If false, skips rows with too many columns.

Max quoted size limit:   
Maximum number of bytes of a quoted value. If a value exceeds the

7. Optionally, edit the parameters.

8. Select **Next**.

## Filters

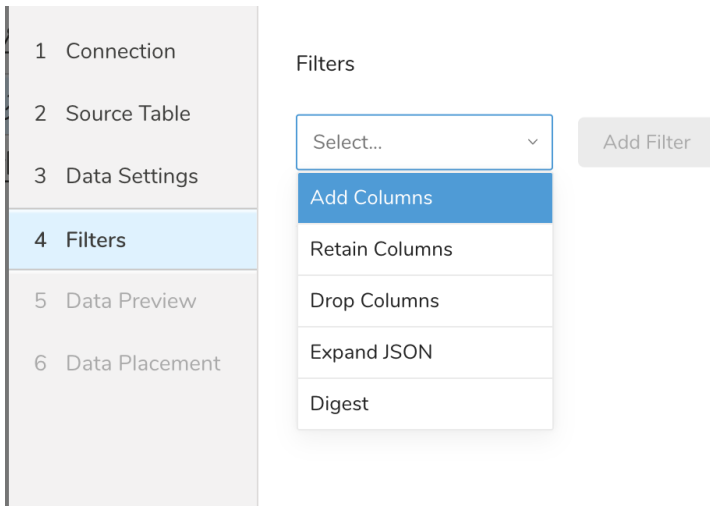
Import Integration Filters enable you to modify your imported data after you have completed [Editing Data Settings](#) for your import.

To apply import integration filters:

1. Select **Next** in Data Settings.

The Filters dialog opens.

2. Select the filter option you want to add.



### 3. Select **Add Filter**.

The parameter dialog for that filter opens.

### 4. Edit the parameters.

For information on each filter type, see one of the following:

- [Retaining Columns Filter](#)
- [Adding Columns Filter](#)
- [Dropping Columns Filter](#)
- [Expanding JSON Filter](#)
- [Digesting Filter](#)

5. Optionally, to add another filter of the same type, select **Add** within the specific column filter dialog.

6. Optionally, to add another filter of a different type, select the filter option from the list and repeat the same steps.

7. After you have added the filters you want, select **Next**.  
The Data Preview dialog opens.

## Data Preview

You can see a [preview](#) of your data before running the import by selecting Generate Preview.

Data shown in the data preview is approximated from your source. It is not the actual data that is imported.

1. Click **Next**.  
Data preview is optional and you can safely skip to the next page of the dialog if you want.
2. To preview your data, select **Generate Preview**. Optionally, click **Next**.
3. Verify that the data looks approximately like you expect it to.

**Create Source**  
Using onetrust\_demo

The preview shows a subset of data from the source based on the data settings. Refer to [help document](#) to learn more about preview data.

8 columns

	Ab_id	Ab_language	Ab_identifier	last_updated_date	Ab_link_token	
1	f7abf910-b5da-47c2-bbee-37f4c86...	NULL	Quan3	2020-09-25 22:42:59...	NULL	0
2	9022117f-cf3c-418c-b527-a8bd9a9...	NULL	Quan2	2020-08-05 03:48:19...	NULL	0
3	a432b52f-3d93-483b-b65f-3c7530...	NULL	Quan4	2020-08-05 03:48:19...	NULL	0
4	233ec0c2-70ab-4de4-ac48-a4a048f...	NULL	Quan5	2020-08-05 03:48:19...	NULL	0
5	f78be70b-8b5d-404e-b663-b606a2...	NULL	Quan1	2020-08-05 03:48:19...	NULL	0
6	db5d9f89-c264-4d82-a246-5939e5...	NULL	example@otprivacy.com	2020-08-06 17:51:12...	NULL	0
7	5ef9542c-315d-4b56-ad1c-c63ad0...	NULL	Michael.White@gmail.com	2020-09-09 20:01:45...	NULL	0
8	3f1dfcb9-1904-4517-9087-0cc45f0...	NULL	Robert.Brown@gmail.com	2020-09-09 20:01:45...	NULL	0
9	4a3a88dd-11a3-4c8b-a1d9-d7043f...	NULL	Mary.Anderson@gmail.com	2020-09-09 20:01:46...	NULL	0
10	4f69893a-9e49-46dc-9519-1cf9dea...	NULL	Elizabeth.Scott@gmail.com	2020-09-09 20:01:47...	NULL	0
11	33342e5d-4c95-4cfe-a622-4e91dc5...	NULL	David.Miller@aol.com	2020-09-09 20:01:47...	NULL	0
12	f54b0d7c-df75-4bf3-934a-dc19a96...	NULL	Robert.Anderson@att.com	2020-09-10 04:57:16...	NULL	0
13	43bfe156-dfba-43b8-964d-1b2a4ae...	NULL	Elizabeth.Miller@google.com	2020-09-10 04:57:16...	NULL	0

Cancel Back Next

4. Select **Next**.

## Data Placement

For data placement, select the target database and table where you want your data placed and indicate how often the import should run.

1. Select **Next**. Under Storage you will create a new or select an existing database and create a new or select an existing table for where you want to place the imported data.

1 Connection

2 Source Table

3 Data Settings

4 Data Preview

5 Data Placement

STORAGE

Database: chung\_default\_db

Table: sftp\_v2\_devproxy

Method:

- Append: Add records into existing table.
- Always Replace: Always clear the destination table before adding records.
- Replace on new data: When there is new data, delete existing data, and insert new data.

Timestamp-based Partition Key: time

Data Storage Timezone: UTC (default)

SCHEDULE

Repeat:  Off  On

Scheduling Timezone: Asia/Saigon

2. Select a **Database** > **Select an existing** or **Create New Database**.
3. Optionally, type a database name.
4. Select a **Table** > **Select an existing** or **Create New Table**.
5. Optionally, type a table name.
6. Choose the method for importing the data.
  - **Append** (default)-Data import results are appended to the table. If the table does not exist, it will be created.
  - **Always Replace**-Replaces the entire content of an existing table with the result output of the query. If the table does not exist, a new table is created.
  - **Replace on New Data**-Only replace the entire content of an existing table with the result output when there is new data.
7. Select the **Timestamp-based Partition Key** column.  
If you want to set a different partition key seed than the default key, you can specify the long or timestamp column as the partitioning time. As a default time column, it uses upload\_time with the add\_time filter.

8. Select the **Timezone** for your data storage.
9. Under **Schedule**, you can choose when and how often you want to run this query.
  - Run once:
    - a. Select **Off**.
    - b. Select **Scheduling Timezone**.
    - c. Select **Create & Run Now**.
  - Repeat the query:
    - a. Select **On**.
    - b. Select the **Schedule**. The UI provides these four options: *@hourly*, *@daily* and *@monthly* or custom *cron*.
    - c. You can also select **Delay Transfer** and add a delay of execution time.
    - d. Select **Scheduling Timezone**.
    - e. Select **Create & Run Now**.

After your transfer has run, you can see the results of your transfer in **Data Workbench > Databases**.

## Validate Connection

Review the job log. Warning and errors provide information about the success of your import. For example, you can [identify the source file names associated with import errors](#).

## Optionally Configure Export Results in Workflow

Within Treasure Workflow, you can specify the use of this data connector to export data.

Learn more at [Using Workflows to Export Data with the TD Toolbelt](#).

### Example Workflow for FTP

```
timezone: UTC

schedule:
  daily>: 02:00:00

sla:
  time: 08:00
  +notice:
    mail>: {data: Treasure Workflow Notification}
    subject: This workflow is taking long time to finish
    to: [meg@example.com]

_export:
  td:
    dest_db: dest_db
    dest_table: dest_table
  ftp:
    ssl: true
    ssl_verify: false

+prepare_table:
  td_ddl>:
    database: ${td.dest_db}
    create_tables: ["${td.dest_table}"]

+load_step:
  td_load>: config/daily_load.yml
  database: ${td.dest_db}
  table: ${td.dest_table}
```