

# GDC Data Transfer Tool User's Guide

NCI Genomic Data Commons (GDC)

# Contents

<b>1</b>	<b>Getting Started</b>	<b>3</b>
	Getting Started . . . . .	3
	The GDC Data Transfer Tool: An Overview . . . . .	3
	Downloading the GDC Data Transfer Tool . . . . .	3
	System Recommendations . . . . .	3
	Binary Distributions . . . . .	3
	Release Notes . . . . .	3
<b>2</b>	<b>Preparing for Data Download and Upload</b>	<b>4</b>
	Overview . . . . .	4
	Obtaining UUIDs for Data Download . . . . .	4
	Obtaining an Authentication Token for Data Downloads . . . . .	4
	Obtaining UUIDs for Data Uploads . . . . .	7
	Obtaining an Authentication Token for Data Uploads . . . . .	9
<b>3</b>	<b>Data Transfer Tool Command Line Documentation</b>	<b>10</b>
	Downloads . . . . .	10
	Downloading Data Using a Manifest File . . . . .	10
	Downloading Data Using GDC File UUIDs . . . . .	10
	Resuming a Failed Download . . . . .	10
	Download Latest Version of a File . . . . .	10
	Downloading Controlled-Access Data . . . . .	11
	Directory structure of downloaded files . . . . .	11
	Uploads . . . . .	11
	Uploading Data Using a Manifest File . . . . .	11
	Uploading Data Using a GDC File UUID . . . . .	11
	Resuming a Failed Upload . . . . .	12
	Deleting Previously Uploaded Data . . . . .	12
	Troubleshooting . . . . .	12
	Invalid Token . . . . .	12
	dbGaP Permissions Error . . . . .	12
	File Availability Error . . . . .	12
	GDC Upload Privileges Error . . . . .	12
	File in Uploaded State Error . . . . .	13
	Microsoft Windows Executable Error . . . . .	13
	Help Menus . . . . .	13
	Root menu . . . . .	13
	Download help menu . . . . .	13
	Upload help menu . . . . .	15
<b>4</b>	<b>Release Notes - Command Line</b>	<b>17</b>
	Data Transfer Tool Release Notes . . . . .	17
	V1.6.1 . . . . .	17
	New Features and Changes . . . . .	17
	Bugs Fixed Since Last Release . . . . .	17
	Known Issues and Workarounds . . . . .	17

V1.6.0 . . . . .	18
New Features and Changes . . . . .	18
Bugs Fixed Since Last Release . . . . .	18
Known Issues and Workarounds . . . . .	18
V1.5.0 . . . . .	18
New Features and Changes . . . . .	18
Bugs Fixed Since Last Release . . . . .	18
Known Issues and Workarounds . . . . .	18
V1.4.0 . . . . .	18
New Features and Changes . . . . .	19
Bugs Fixed Since Last Release . . . . .	19
Known Issues and Workarounds . . . . .	19
v1.3.0 . . . . .	19
New Features and Changes . . . . .	19
Bugs Fixed Since Last Release . . . . .	19
Known Issues and Workarounds . . . . .	19
v1.2.0 . . . . .	20
New Features and Changes . . . . .	20
Bugs Fixed Since Last Release . . . . .	20
Known Issues and Workarounds . . . . .	20
v1.1.0 . . . . .	20
New Features and Changes . . . . .	20
Bugs Fixed Since Last Release . . . . .	20
Known Issues and Workarounds . . . . .	20
v1.0.1 . . . . .	21
New Features and Changes . . . . .	21
Bugs Fixed Since Last Release . . . . .	21
Known Issues and Workarounds . . . . .	21
v1.0.0 . . . . .	21
New Features and Changes . . . . .	21
Bugs Fixed Since Last Release . . . . .	21
Known Issues and Workarounds . . . . .	21
<b>5 Data Transfer Tool UI Documentation</b>	<b>22</b>
<b>6 Release Notes - UI</b>	<b>29</b>
Data Transfer Tool UI Release Notes . . . . .	29
v1.0.0 . . . . .	29
New Features and Changes . . . . .	29
Bugs Fixed Since Last Release . . . . .	29
Known Issues and Workarounds . . . . .	29
v0.6.0 . . . . .	29
New Features and Changes . . . . .	29
Bugs Fixed Since Last Release . . . . .	30
Known Issues and Workarounds . . . . .	30
v0.5.4 . . . . .	30
New Features and Changes . . . . .	30
Bugs Fixed Since Last Release . . . . .	30
Known Issues and Workarounds . . . . .	30
v0.5.3 . . . . .	30
New Features and Changes . . . . .	30
Bugs Fixed Since Last Release . . . . .	30
Known Issues and Workarounds . . . . .	30

# Chapter 1

## Getting Started

### Getting Started

#### The GDC Data Transfer Tool: An Overview

Raw sequence data, stored as BAM files, make up the bulk of data stored at the NCI Genomic Data Commons (GDC). The size of a single file can vary greatly. Most BAM files stored in the GDC are in the 50 MB - 40 GB size range, with some of the whole genome BAM files reaching sizes of 200-300 GB.

The GDC Data Transfer Tool, a command-line driven application, provides an optimized method of transferring data to and from the GDC and enables resumption of interrupted transfers.

#### Downloading the GDC Data Transfer Tool

##### System Recommendations

The system recommendations for using the GDC Data Transfer Tool are as follows:

- OS: Linux (Ubuntu 16.x or later), OS X (10.9 Mavericks or later), or Windows (8 or later)
- CPU: At least two 64-bit cores, Intel or AMD
- RAM: At least 8 GiB
- Storage: Enterprise-class storage system capable of at least 1 Gb/s (gigabit per second) write throughput and sufficient free space for BAM files.

##### Binary Distributions

Binary distributions are available on the GDC Transfer Tool page. To install the GDC Data Transfer Tool, download the respective binary distribution and unzip the distribution's archive to a location on the target system. It is recommended that the binary be copied to a location that is in the user's path so that it is accessible from any location within the terminal or command prompt.

##### Release Notes

Release Notes are available on the GDC Data Transfer Tool Release Notes Page.

## Chapter 2

# Preparing for Data Download and Upload

###Preparing for Data Downloads and Uploads

### Overview

The GDC Data Transfer Tool is intended to be used in conjunction with the GDC Data Portal and the GDC Data Submission Portal to transfer data to or from the GDC. First, the GDC Data Portal's interface is used to generate a manifest file or obtain UUID(s) and (for Controlled-Access Data) an authentication token. The GDC Data Transfer Tool is then used to transfer the data files listed in the manifest file or identified by UUID(s).

####Downloads ### Obtaining a Manifest File for Data Download

The GDC Data Transfer Tool supports downloading multiple files listed in a GDC manifest file. Manifest files can be generated and downloaded directly from the GDC Data Portal:

First, select the data files of interest. Click the *Cart* button in the row corresponding to the file desired. The button will turn green to indicate that the file has been selected.

Once all files of interest have been selected, click on the *Cart* button in the upper right-hand corner. This will bring up the cart page, which provides an overview of all currently selected files. This list of files can be downloaded as a manifest file by clicking on the green *Download* button and selecting *Manifest* from the drop down.

### Obtaining UUIDs for Data Download

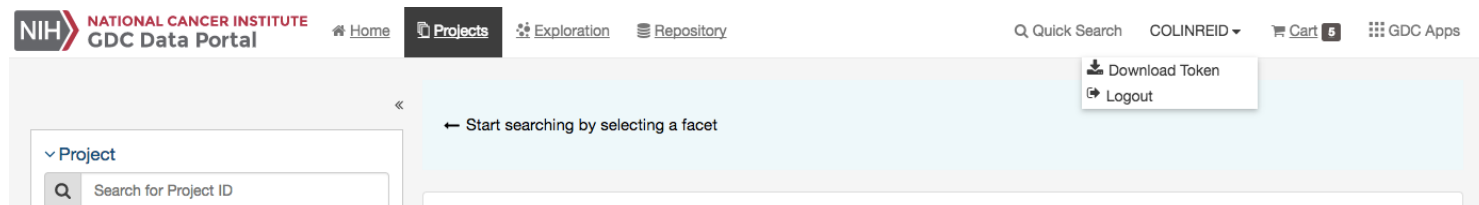
A manifest file is not required to download files from GDC. The GDC Data Transfer Tool will accept file UUID(s) instead of a manifest file for downloading individual data files. To obtain a data file's UUID from the GDC Data Portal, click the file name to find its detail page including its GDC UUID.

### Obtaining an Authentication Token for Data Downloads

The GDC Data Transfer Tool requires an authentication token to download from GDC data portal to download Controlled-Access Data. Tokens can be generated and downloaded directly from the GDC Data Portal.

To generate a token, first log in to the GDC Data Portal by clicking the *Login* button in the top right corner of the page. This will redirect to the eRA Commons login page. After successful authentication, the GDC Data Portal will display the username in place of the *Login* button. Here, the user Ian Miller is logged in to the GDC Data Portal, indicated by the username IANMILLER.

Clicking the username will open a drop-down menu. Select *Download Token* from the menu to generate an authentication token.



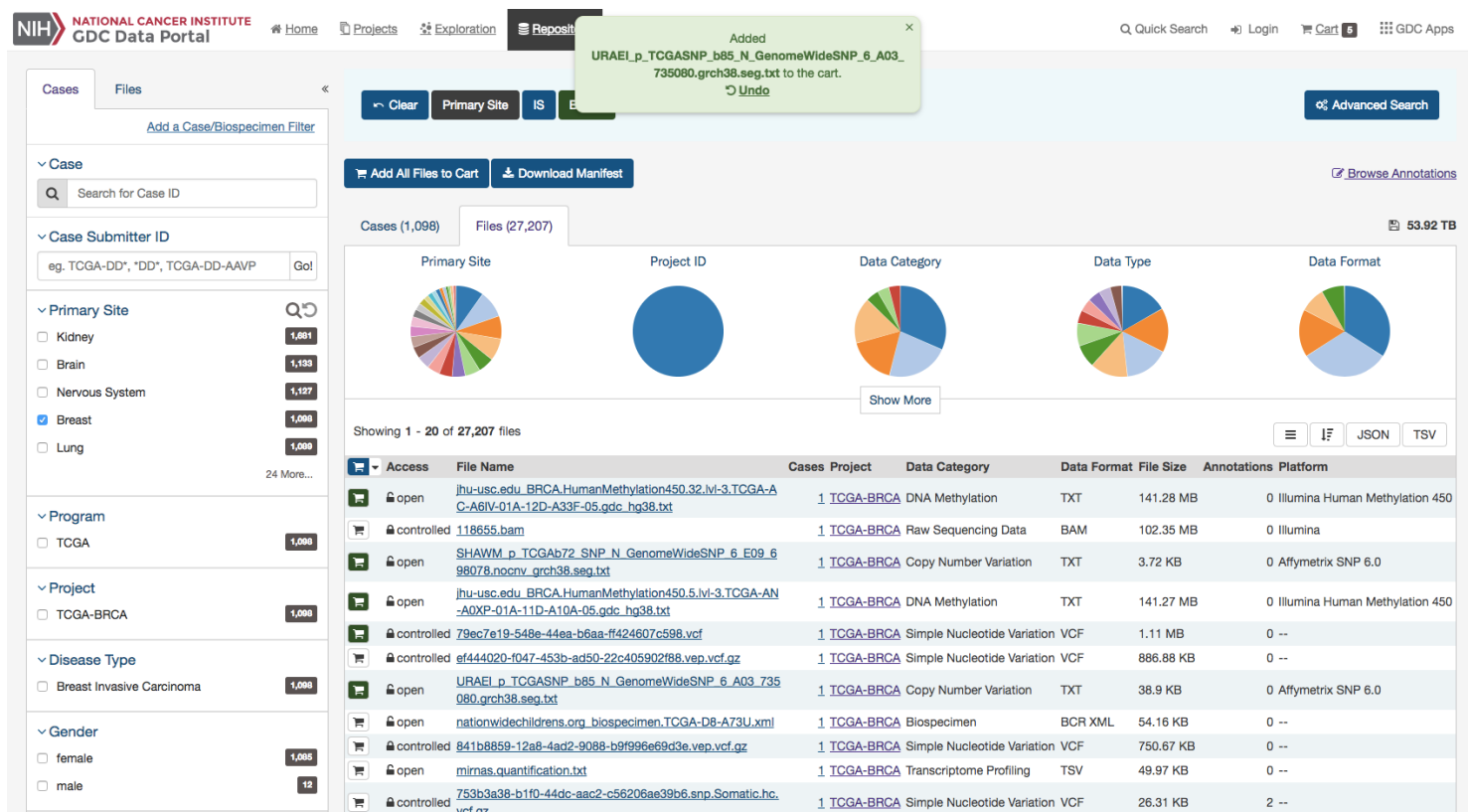


Figure 2.1: GDC Data Portal: Selecting Files of Interest

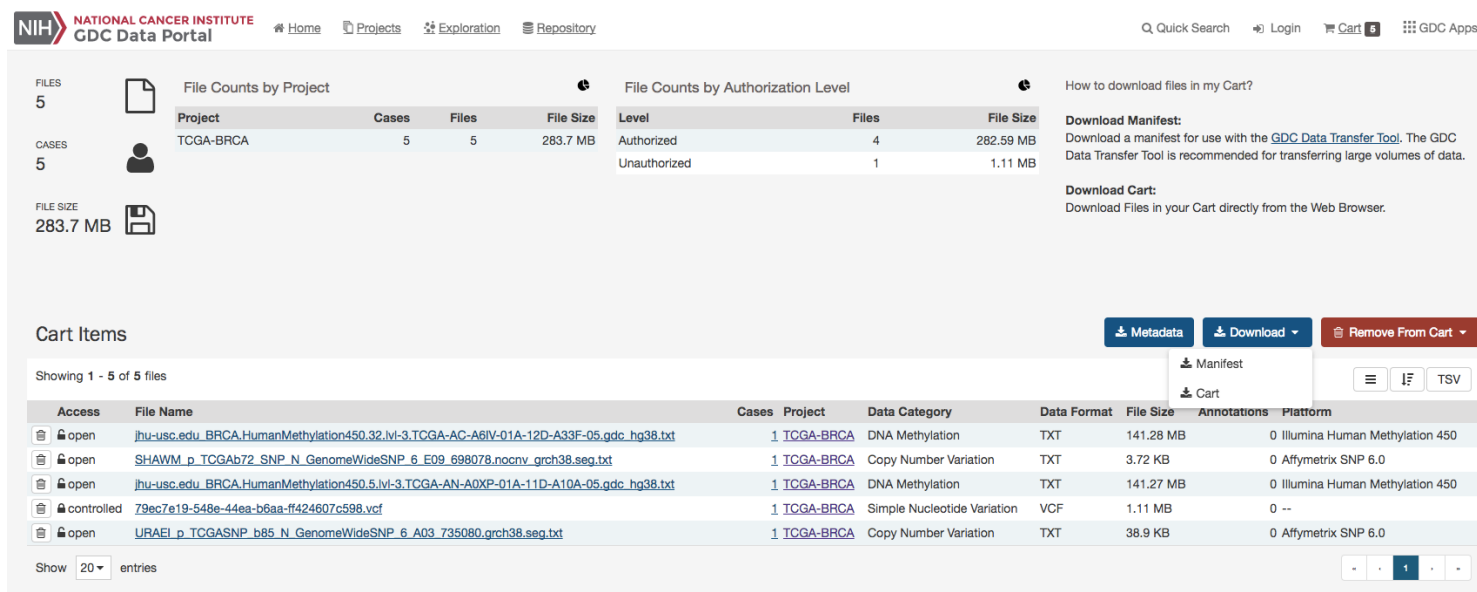


Figure 2.2: GDC Data Portal: Cart Page

**FL** 512994a9-2bf7-4fe3-994c-f2a80c57b0f1

[Add to Cart](#)
[BAM Slicing](#)
[Download](#)

### File Properties

<b>Name</b>	09001cae-e831-4af8-bd86-e22cf21c2525_gdc_realn_rehead.bam
<b>Access</b>	controlled
<b>UUID</b>	512994a9-2bf7-4fe3-994c-f2a80c57b0f1
<b>Submitter ID</b>	09001cae-e831-4af8-bd86-e22cf21c2525
<b>Data Format</b>	BAM
<b>Size</b>	2.8 GB
<b>MD5 Checksum</b>	3bc97d784f95e3f37d945a0583380a10
<b>Archive</b>	--
<b>Project ID</b>	<a href="#">TCGA-COAD</a>

### Data Information

<b>Data Category</b>	Raw Sequencing Data
<b>Data Type</b>	Aligned Reads
<b>Experimental Strategy</b>	RNA-Seq
<b>Platform</b>	Illumina

### Associated Cases/Biospecimen

Type to filter cases.

Entity Id	Entity Type	Case UUID	Annotations
<a href="#">c30ce88d-5dff-4503-b090-01b4b6aa0b80</a>	aliquot	<a href="#">c0b8c55c-b993-481d-aeaa-9ebfa64ee20e</a>	0

### Analysis

<b>Analysis ID</b>	dbb0f6f6-ca45-487a-9b9a-7711b7d40c8b
<b>Workflow Type</b>	STAR 2-Pass
<b>Workflow Completion Date</b>	2017-03-04
<b>Source Files</b>	0

### Reference Genome

<b>Genome Build</b>	GRCh38.p0
<b>Genome Name</b>	GRCh38.d1.vd1

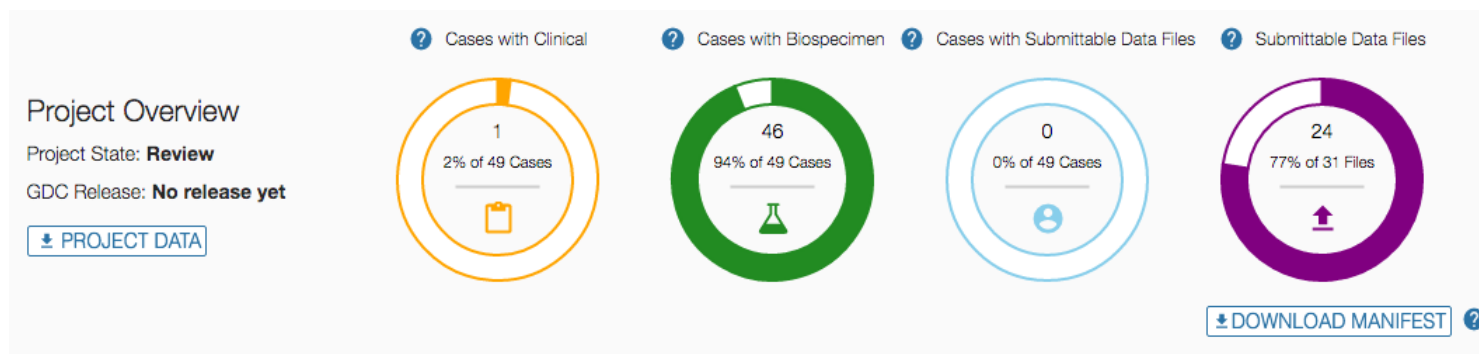
### Read Groups

Read Group ID	Is Paired End	Read Length	Library Name	Sequencing Center	Sequencing Date
803b829d-0e7a-4aa4-8b77-963d0d01b2ce	true	76	unknown	UNC	--

Figure 2.3: GDC Data Portal: Detailed File Page

**NOTE:** The authentication token should be kept in a secure location, as it allows access to all data accessible by the associated user.

####Uploads #### Obtaining a Manifest File for Data Uploads Multiple data file uploads are supported by the GDC Data Transfer Tool via a manifest file. Manifest files can be generated and downloaded directly from the GDC Submission Portal. A project's manifest file can be downloaded from the projects's dashboard.



**NOTE:** To download a project's manifest file click on the *Download Manifest* button located on the home page of the project, just below the four status charts. A manifest will be generated for the entire project or if previous files have already been upload only the files that remain to be uploaded.

A manifest for individual files can also be downloaded from the transaction tab and browse tab pages of the submission portal's project. More information on the process can be found under the Submission Portal's documentation section entitled Uploading the Submittable Data File to the GDC.

## Obtaining UUIDs for Data Uploads

A UUID can be used for data submission with the Data Transfer Tool. The UUID for submittable data uploads can be obtained from the Submission Portal or from the API GraphQL endpoint. In the Submission Portal the UUID for a data file can be found in the Manifest YAML file located in the *id:* row located under the file size entry.

```
files:
- data_category: Raw Sequencing Data
  data_format: FASTQ
  data_type: Unaligned Reads
  experimental_strategy: WGS
  file_name: GDC-INTERNAL-000084-S1-Q1-RG1.fastq.zip
  file_size: 138112000
  id: c414a205-376e-4993-af48-2a4689eb433e
  local_file_path: GDC-INTERNAL-000084-S1-Q1-RG1.fastq.zip
  md5sum: e0bb0367ffbc287dcf10ed4212a740a2
  project_id: GDC-INTERNAL
  read_groups:
  - id: 4231ef42-4f24-48f1-88da-aa98b492e57e
    submitter_id: GDC-INTERNAL-000084-S1-Q1-RG1
  state_comment: null
  submitter_id: GDC-INTERNAL-000084-S1-Q1-RG1.fastq.zip
  type: submitted_unaligned_reads
```

Figure 2.4: Submission Manifest yaml file

A second location to obtain a UUID in the Submission Portal is on the Browse Tab page. Under the Submittable Data Files section a UUID can be found by opening up the file's detail page. By clicking on the Submitter ID of the upload file a new



window will display a Summary of the file's details, which contains the UUID.

The screenshot shows the NIH National Cancer Institute GDC Data Submission Portal. The main content area displays a table of 'Submitted Unaligned Reads'. The table has columns for Submitter ID, Type, Case ID, Status, File Status, and Last Updated. A specific record is highlighted, showing a UUID: 745ce3f6-adc3-43d3-8621-8a70c7b9cbdf. The right sidebar shows the 'BLOOD-00001-aliquot\_lane1\_barcode...' details, including a 'SUMMARY' section with fields like Type, UUID, Project Id, Submitter Id, Created Datetime, File State, State, and Updated Datetime. The 'DETAILS' section shows Data Category, Data Format, Data Type, Experimental Strategy, File Name, File Size, and Md5sum.

Figure 2.5: Submission Portal Browse Page Details

## #####GraphQL

A UUID can be obtained from the API GraphQL endpoint. An overview of what GraphQL and its uses is located on the API documentation page section Querying Submitted Data Using GraphQL

The following example will query the endpoint to produce a UUID along with submitter\_id, file\_name, and project\_id.

```
1 {
2   submitted_unaligned_reads (project_id: "GDC-INTERNAL", submitter_id:
3     "Blood-00001-aliquot_lane1_barcode23.fastq") {
4     id
5     submitter_id
6     file_name
7     project_id
8   }
```

```
1 {\n  submitted_unaligned_reads (project_id: \"GDC-INTERNAL\", submitter_id:\n    \"Blood-00001-aliquot_lane1_barcode23.fastq\") {\n    id\n    submitter_id\n    file_name\n    project_id\n  }\n}
```

```
1 {
2   "query": "{\n  \n  submitted_unaligned_reads (project_id: \"GDC-INTERNAL\", submitter_id:\n    \"Blood-00001-aliquot_lane1_barcode23.fastq\") {\n    id\n    submitter_id\n    file_name\n    project_id\n  }\n}",
3   "variables": null
4 }
```

```

1 export
   token=ALPHANUMERICTOKEN-01234567890+AlPhAnUmErIcToKeN=0123456789-ALPHANUMERICTOKEN-01234567890+AlPhAnUmErIcT
2 $ curl --request POST --header "X-Auth-Token: $token" 'https://api.gdc.cancer.gov/v0/submission/graphql'
   -d@data.json

1 {
2   "data": {
3     "submitted_unaligned_reads": [
4       {
5         "file_name": "dummy.fastq",
6         "id": "616eab2f-791a-4641-8cd6-ee195a10a201",
7         "project_id": "GDC-INTERNAL",
8         "submitter_id": "Blood-00001-aliquot_lane1_barcode23.fastq"
9       }
10    ]
11  }


```

## Obtaining an Authentication Token for Data Uploads

While biospecimen and clinical metadata may be uploaded via the GDC Data Submission Portal, file upload must be done using the Data Transfer Tool or API. An authentication token is required for data upload and can be generated on the GDC Data Submission Portal.



To generate a token, first log in to the GDC Data Submission Portal by clicking the *Login* button in the top right corner of the page. This will create a popup window that will redirect to the eRA Commons login page. After successful authentication, the GDC Submission Portal will display the username in place of the *Login* button. Here, the user Ian Miller is logged in to the GDC Submission Portal, indicated by the username IANMILLER.

Clicking the username will open a drop-down menu. Select *Download Token* from the menu to generate an authentication token.



**NATIONAL CANCER INSTITUTE**  
**GDC Data Submission Portal**

IANMILLER

 Download Token  
 Logout

GDC APPS

WELCOME TO THE GDC DATA SUBMISSION PORTAL

The GDC Data Submission Portal allows researchers to submit and release clinical, biospecimen, and experimental data for studies registered in dbGaP into GDC. Select a project from the project list to submit and release data as well as view previously submitted data and transactions.

DOCUMENTATION

[User's Guide](#)

Tutorial (Coming soon)

[Submission Workflow](#)

PROJECTS (15)

FILTER PROJECTS

ID	Name	Primary Site	Submission State	Release ?	Last Updated
<a href="#">GDC-INTERNAL</a>	Internal	Lung	Open	<div>RELEASE</div>	2016-10-26 11:59
<a href="#">TRIO-CRU</a>	Ukrainian National Research Center for Radiation Medicine Trio Study		Open	<div>RELEASE</div>	2016-10-21 08:37

## Chapter 3

# Data Transfer Tool Command Line Documentation

###Data Transfer Tool Command Line Documentation

## Downloads

### Downloading Data Using a Manifest File

A convenient way to download multiple files from the GDC is to use a manifest file generated by the GDC Data Portal. After generating a manifest file (see [Preparing for Data Download and Upload](#) for instructions), initiate the download using the GDC Data Transfer Tool by supplying the **-m** or **-manifest** option, followed by the location and name of the manifest file. OS X users can drag and drop the manifest file into Terminal to provide its location.

The following is an example of a command for downloading files from GDC using a manifest file:

```
1 gdc-client download -m /Users/JohnDoe/Downloads/gdc_manifest_6746fe840d924cf623b4634b5ec6c630bd4c06b5.txt
```

### Downloading Data Using GDC File UUIDs

The GDC Data Transfer Tool also supports downloading of one or more individual files using UUID(s) instead of a manifest file. To do this, enter the UUID(s) after the download command:

```
1 gdc-client download 22a29915-6712-4f7a-8dba-985ae9a1f005
```

Multiple UUIDs can be specified, separated by a space:

```
1 gdc-client download e5976406-473a-4fbd-8c97-e95187cdc1bd fb3e261b-92ac-4027-b4d9-eb971a92a4c3
```

### Resuming a Failed Download

The GDC Data Transfer Tool supports resumption of interrupted downloads. To resume an incomplete download, repeat the download of the manifest or UUID(s) in the same folder as the initial download. Failed downloads will appear in the destination folder with a `.partial` extension. This feature allows users the ability to identify quickly where the download stopped. For large downloads this feature can let the user identify where the download was interrupted and edit the manifest accordingly.

```
1 gdc-client download f80ec672-d00f-42d5-b5ae-c7e06bc39da1
```

### Download Latest Version of a File

The GDC Data Transfer Tool supports file versioning. Our backend data storage supports multiple file versions so older and current versions can be accessible to our users. For information about accessing file versioning information with our API and finding older UUID information from current UUIDs please check out the [API User Guide](#) section in our [API documentation](#). When working with older manifests or older lists of UUIDs the latest version of a file can always be download with the `-latest` flag.

```
1 gdc-client download 426de656-7e34-4a49-b87e-6e2563fa3cdd --latest -t gdc-user-token.2018.txt
```

```
1 Downloading LATEST versions of files
```

```
2 Latest version for 426de656-7e34-4a49-b87e-6e2563fa3cdd ==> 6633bfbd-87f1-4d3a-a475-7ad1e8c2017a
```

```
3 100%
```

```
[#####]
```

```
Time: 0:01:16 14.10 MB/s
```

```
4 Successfully downloaded: 1
```

## Downloading Controlled-Access Data

A user authentication token is required for downloading Controlled-Access Data from GDC. Tokens can be obtained from the GDC Data Portal (see instructions in Obtaining an Authentication Token). Once downloaded, the token *file* can be passed to the GDC Data Transfer Tool using the **-t** or **-token-file** option:

```
1 gdc-client download -m gdc_manifest_e24fac38d3b19f67facb74d3efa746e08b0c82c2.txt -t
gdc-user-token.2015-06-17T09-10-02-04-00.txt
```

## Directory structure of downloaded files

The directory in which the files are downloaded will include folders named by the file UUID. Inside these folders, along with the the data and zipped metadata or index files, will exist a logs folder. The logs folder contains state files that insure that downloads are accurate and allow for resumption of failed or prematurely stopped downloads. While a download is in progress a file will have a .partial extension. This will also remain if a download failed. Once a file is finished downloading the extension will be removed. If an identical manifest is retried another attempt will be made to download files containing a .partial extension.

```
1 C501.TCGA-BI-A0VR-10A-01D-A10S-08.5_gdc_realn.bam.partial logs
```

## Uploads

### Uploading Data Using a Manifest File

GDC Data Transfer Tool supports uploading molecular data using a manifest file to the Data Submission Portal. The manifest file for submittable data files can be retrieved from the GDC Data Submission Portal, or directly from the GDC Submission API given a submittable data file UUID. The user authentication token file needs to be specified using the **-t** or **-token-file** option.

First, generate an upload manifest, either using the GDC Data Submission Portal, or using a call to the GDC Submission API **manifest** endpoint (as in the following example):

```
1 export
token=ALPHANUMERICTOKEN-01234567890+AlPhAnUmErIcToKeN=0123456789-ALPHANUMERICTOKEN-01234567890+AlPhAnUmErIcT
2
3 curl --header "X-Auth-Token: $token"
'https://api.gdc.cancer.gov/submission/CGCI/BLGSP/manifest?ids=460ad2fe-5a7f-4797-9e18-336d33e21444'
>manifest.yml
```

```
1 gdc-client upload --manifest manifest.yml --token-file token.txt
```

### Uploading Data Using a GDC File UUID

The GDC Data Transfer Tool also supports uploading molecular data using a file UUID. The tool will first make a request to get the filename and project id from GDC API, and then upload the corresponding file from the current directory.

```
1 gdc-client upload cd939bdd-b607-4dd4-87a6-fad12893932d -t token.txt
```

## Resuming a Failed Upload

By default, GDC Data Transfer Tool uses multipart transfer to upload files. If an upload failed but some parts were transmitted successfully, a resume file will be saved with the filename *resume\_*[manifest\_filename]**. Running the upload command again will resume the transfer of only those parts of the file that failed to upload in the previous attempt.

```
1 gdc-client upload -m manifest.yml -t token
```

## Deleting Previously Uploaded Data

Previously uploaded data can be replaced with new data by deleting it first using the **-delete** switch:

```
1 gdc-client upload -m manifest.yml -t token --delete
```

## Troubleshooting

### Invalid Token

An error message about an ‘invalid token’ means that a new authentication token needs to be obtained from the GDC Data Portal or the GDC Data Submission Portal as described in [Preparing for Data Download and Upload](#).

```
1 403 Client Error: FORBIDDEN: {
2   "message": "Your token is invalid or expired, please get a new token from GDC Data Portal"
3 }
```

### dbGaP Permissions Error

Users may see the following error message when attempting to download a file from GDC:

```
1 403 Client Error: FORBIDDEN: {
2   "message": "You don't have access to the data: Please specify a X-Auth-Token"
3 }
```

This error message indicates that the user does not have dbGaP access to the project to which the file belongs. Instructions for requesting access from dbGaP can be found [here](#).

### File Availability Error

Users may also see the following error message when attempting to download a file from GDC:

```
1 403 Client Error: FORBIDDEN: {
2   "message": "You don't have access to the data: Requested file abd28349-92cd-48a3-863a-007a218de80f
3   does not allow read access"
3 }
```

This error message means that the file is not available for download. This may be because the file has not been uploaded or released yet or that it is not a file entity.

### GDC Upload Privileges Error

Users may see the following error message when attempting to upload a file:

```
1 Can't upload: {
2   "message": "You don't have access to the data: You don't have create role to do 'upload'"
3 }
```

This means that the user has dbGaP read access to the data, but does not have GDC upload privileges. Users can contact The database of Genotypes and Phenotypes (dbGaP) to request upload privileges.

## File in Uploaded State Error

Re-uploading a file may return the following error:

```
1 Can't upload: {
2     "message": "File in uploaded state, upload not allowed"
3 }
```

To resolve this issue, delete the file using the **-delete** switch before re-uploading.

## Microsoft Windows Executable Error

Attempting to run gdc-client.exe by double-clicking it in the Windows Explorer will produce a window that blinks once and disappears.

This is normal, the executable must be run using the command prompt. Click ‘Start’, followed by ‘Run’ and type ‘cmd’ into the text bar. Then navigate to the path containing the executable using the ‘cd’ command.

## Help Menus

The GDC Data Transfer Tool comes with built-in help menus. These menus are displayed when the GDC Data Transfer Tool is run with flags -h or -help for any of the main arguments to the tool. Running the GDC Data Transfer Tool without argument or flag will present a list of available command options.

```
1 gdc-client --help

1 usage: gdc-client [-h] [--version] {download,upload,settings} ...
2
3 The Genomic Data Commons Command Line Client
4
5 optional arguments:
6   -h, --help            show this help message and exit
7   --version              show program's version number and exit
8
9 commands:
10  {download,upload,settings}
11      for more information, specify -h after a command
12  download              download data from the GDC
13  upload                upload data to the GDC
14  settings              display default settings
```

The available menus are provided below.

### Root menu

The GDC Data Transfer Tool displays the following output when executed without any arguments.

```
1 gdc-client

1 usage: gdc-client [-h] [--version] {download,upload,settings} ...
2 gdc-client: error: too few arguments
```

### Download help menu

The GDC Data Transfer Tool displays the following help menu for its download functionality.

```
1 gdc-client download --help
```

```

1 usage: gdc-client download [-h] [--debug]
2                             [--log-file LOG_FILE]
3                             [--color_off] [-t TOKEN_FILE]
4                             [-d DIR] [-s server]
5                             [--no-segment-md5sums]
6                             [--no-file-md5sum]
7                             [-n N_PROCESSES]
8                             [--http-chunk-size HTTP_CHUNK_SIZE]
9                             [--save-interval SAVE_INTERVAL]
10                            [--no-verify]
11                            [--no-related-files]
12                            [--no-annotations]
13                            [--no-auto-retry]
14                            [--retry-amount RETRY_AMOUNT]
15                            [--wait-time WAIT_TIME]
16                            [--latest] [--config FILE] [-u]
17                            [-m MANIFEST]
18                            [file_id [file_id ...]]
19
20 positional arguments:
21 file_id                The GDC UUID of the file(s) to download
22
23 optional arguments:
24 -h, --help              show this help message and exit
25 --debug                Enable debug logging. If a failure occurs, the program
26                        will stop.
27 --log-file LOG_FILE    Save logs to file. Amount logged affected by --debug
28 --color_off            Disable colored output
29 -t TOKEN_FILE, --token-file TOKEN_FILE
30                        GDC API auth token file
31 -d DIR, --dir DIR      Directory to download files to. Defaults to current
32                        dir
33 -s server, --server server
34                        The TCP server address server[:port]
35 --no-segment-md5sums    Do not calculate inbound segment md5sums and/or do not
36                        verify md5sums on restart
37 --no-file-md5sum       Do not verify file md5sum after download
38 -n N_PROCESSES, --n-processes N_PROCESSES
39                        Number of client connections.
40 --http-chunk-size HTTP_CHUNK_SIZE, -c HTTP_CHUNK_SIZE
41                        Size in bytes of standard HTTP block size.
42 --save-interval SAVE_INTERVAL
43                        The number of chunks after which to flush state file.
44                        A lower save interval will result in more frequent
45                        printout but lower performance.
46 --no-verify            Perform insecure SSL connection and transfer
47 --no-related-files     Do not download related files.
48 --no-annotations       Do not download annotations.
49 --no-auto-retry        Ask before retrying to download a file
50 --retry-amount RETRY_AMOUNT
51                        Number of times to retry a download
52 --wait-time WAIT_TIME  Amount of seconds to wait before retrying
53
54 --latest               Download latest version of a file if it exists
55 --config FILE          Path to INI-type config file
56 -u, --udt              Use the UDT protocol.
57 -m MANIFEST, --manifest MANIFEST
58                        GDC download manifest file

```

## Upload help menu

The GDC Data Transfer Tool displays the following help menu for its upload functionality.

```
1 gdc-client upload --help
```

```
1 usage: gdc-client upload [-h] [--debug]
2                               [--log-file LOG_FILE]
3                               [--color_off] [-t TOKEN_FILE]
4                               [--project-id PROJECT_ID]
5                               [--path path]
6                               [--upload-id UPLOAD_ID]
7                               [--insecure] [--server SERVER]
8                               [--part-size PART_SIZE]
9                               [--upload-part-size UPLOAD_PART_SIZE]
10                              [-n N_PROCESSES]
11                              [--disable-multipart] [--abort]
12                              [--resume] [--delete]
13                              [--manifest MANIFEST]
14                              [--config FILE]
15                              [file_id [file_id ...]]
16 positional arguments:
17   file_id                The GDC UUID of the file(s) to upload
18
19 optional arguments:
20   -h, --help              show this help message and exit
21   --debug                 Enable debug logging. If a failure occurs, the program
22                           will stop.
23   --log-file LOG_FILE     Save logs to file. Amount logged affected by --debug
24   --color_off             Disable colored output
25   -t TOKEN_FILE, --token-file TOKEN_FILE
26                           GDC API auth token file
27   --project-id PROJECT_ID, -p PROJECT_ID
28                           The project ID that owns the file
29   --path path, -f path    directory path to find file
30   --upload-id UPLOAD_ID, -u UPLOAD_ID
31                           Multipart upload id
32   --insecure, -k          Allow connections to server without certs
33   --server SERVER, -s SERVER
34                           GDC API server address
35   --part-size PART_SIZE
36                           DEPRECATED in favor of [--upload-part-size]
37   --upload-part-size UPLOAD_PART_SIZE, -c UPLOAD_PART_SIZE
38                           Part size for multipart upload
39   -n N_PROCESSES, --n-processes N_PROCESSES
40                           Number of client connections
41   --disable-multipart     Disable multipart upload
42   --abort                 Abort previous multipart upload
43   --resume, -r            Resume previous multipart upload
44   --delete                Delete an uploaded file
45   --manifest MANIFEST, -m MANIFEST
46                           Manifest which describes files to be uploaded
47   --config FILE           Path to INI-type config file
```

####Data Transfer Tool Configuration File The DTT has the ability to save and reuse configuration parameters in the format of a flat text file via a command line argument. A simple text file needs to be created first with an extension of either txt or dtt. The supported section headers are upload and download which can be used independently of each other or used in the same configuration file. Each section header corresponds to the main functions of the application which are to either download data from the GDC portals or to upload data to the submission system of the GDC. The configurable parameters are those listed in



the help menus under either download or upload displayed under the output tabs.

Example usage:

```
1 gdc-client download d45ec02b-13c3-4afa-822d-443ccd3795ca --config my-ddt-config.dtt
```

Example of configuration file:

```
1 [upload]
2 path = /some/upload/path
3 upload_part_size = 1073741824
4
5
6 [download]
7 dir = /some/download/path
8 http_chunk_size = 2048
9 retry_amount = 6
```

#####Display Config Parameters This command line flag can be used with either the download or upload application feature to display what settings are active within a custom data transfer tool configuration file.

```
1 gdc-client settings download --config my-ddt-config.dtt
2 [download]
3 no_auto_retry = False
4 no_file_md5sum = False
5 save_interval = 1073741824
6 http_chunk_size = 2048
7 server = http://exmple-site.com
8 n_processes = 8
9 no_annotations = False
10 no_related_files = False
11 retry_amount = 6
12 no_segment_md5sum = False
13 manifest = []
14 wait_time = 5.0
15 no_verify = True
16 dir = /some/download/path
```

## Chapter 4

# Release Notes - Command Line

### Data Transfer Tool Release Notes

Version	Date
v1.6.1	May 17, 2021
v1.6.0	July 8, 2020
v1.5.0	January 30, 2020
v1.4.0	December 18, 2018
v1.3.0	August 22, 2017
v1.2.0	Oct 31, 2016
v1.1.0	September 7, 2016
v1.0.1	June 2, 2016
v1.0.0	May 26, 2016

### V1.6.1

- **GDC Product:** Data Transfer Tool
- **Release Date:** May 17, 2021

### New Features and Changes

- None

### Bugs Fixed Since Last Release

- Fixed issue with resuming large file downloads.
- Fixed issue with error reporting.
- Improved multi-part file upload.

### Known Issues and Workarounds

- Use of non-ASCII characters in token passed to Data Transfer Tool will produce incorrect error message “Internal server error: Auth service temporarily unavailable”.
- On some terminals, dragging and dropping a file into the interactive client will add single quotes ( ' ') around the file path. This causes the interactive client to misinterpret the file path and generate an error when attempting to load a manifest file or token.
  - *Workaround:* Manually type out the file name or remove the single quotes from around the file path.
- When any files mentioned in the upload manifest are not present in the upload directory the submission will hang at the missing file.
  - *Workaround:* Edit the manifest to specify only the files that are present in the upload directory for submission or copy the missing files into the upload directory.

## V1.6.0

- **GDC Product:** Data Transfer Tool
- **Release Date:** July 8, 2020

### New Features and Changes

- None

### Bugs Fixed Since Last Release

- Fixed issue with file upload requiring a manifest
- Restored multi-part upload feature
- Fixed error reporting issue

### Known Issues and Workarounds

- Use of non-ASCII characters in token passed to Data Transfer Tool will produce incorrect error message “Internal server error: Auth service temporarily unavailable”.
- On some terminals, dragging and dropping a file into the interactive client will add single quotes ( ' ') around the file path. This causes the interactive client to misinterpret the file path and generate an error when attempting to load a manifest file or token.
  - *Workaround:* Manually type out the file name or remove the single quotes from around the file path.
- When any files mentioned in the upload manifest are not present in the upload directory the submission will hang at the missing file.
  - *Workaround:* Edit the manifest to specify only the files that are present in the upload directory for submission or copy the missing files into the upload directory.

## V1.5.0

- **GDC Product:** Data Transfer Tool
- **Release Date:** January 30, 2020

### New Features and Changes

- Data transfer tool code now uses Python 3.

### Bugs Fixed Since Last Release

- Problems with downloading associated annotations is fixed.

### Known Issues and Workarounds

- Use of non-ASCII characters in token passed to Data Transfer Tool will produce incorrect error message “Internal server error: Auth service temporarily unavailable”.
- On some terminals, dragging and dropping a file into the interactive client will add single quotes ( ' ') around the file path. This causes the interactive client to misinterpret the file path and generate an error when attempting to load a manifest file or token.
  - *Workaround:* Manually type out the file name or remove the single quotes from around the file path.
- When any files mentioned in the upload manifest are not present in the upload directory the submission will hang at the missing file.
  - *Workaround:* Edit the manifest to specify only the files that are present in the upload directory for submission or copy the missing files into the upload directory.

## V1.4.0

- **GDC Product:** Data Transfer Tool
- **Release Date:** December 18, 2018

## New Features and Changes

- Enabled download latest file version feature
- Removal of Interactive mode
- Enabled display of all default settings
- Standardized upload and download help menus

## Bugs Fixed Since Last Release

- Download flag `-no-related-files` bug preventing file downloads fixed
- File name handling with forward slashes bug fixed
- Download flag `-no-segment-md5sums` bug fixed.

## Known Issues and Workarounds

- Use of non-ASCII characters in token passed to Data Transfer Tool will produce incorrect error message “Internal server error: Auth service temporarily unavailable”.
- On some terminals, dragging and dropping a file into the interactive client will add single quotes ( ' ') around the file path. This causes the interactive client to misinterpret the file path and generate an error when attempting to load a manifest file or token.
  - *Workaround:* Manually type out the file name or remove the single quotes from around the file path.
- When any files mentioned in the upload manifest are not present in the upload directory the submission will hang at the missing file.
  - *Workaround:* Edit the manifest to specify only the files that are present in the upload directory for submission or copy the missing files into the upload directory.

## v1.3.0

- **GDC Product:** Data Transfer Tool
- **Release Date:** August 22, 2017

## New Features and Changes

- Faster performance when downloading many small files
- Faster performance overall
- Better handling of time outs
- Uses new default API URL (<https://api.gdc.cancer.gov>)
- Better logging

## Bugs Fixed Since Last Release

- Submission manifest **local\_file\_path:** will now modify path as expected
- Upload flags `-path/-f` will modify the upload path as expected
- When deleting uploaded files you will no longer need a file in the current directory of the same name
- Can specify manifest path for upload

## Known Issues and Workarounds

- Use of non-ASCII characters in token passed to Data Transfer Tool will produce incorrect error message “Internal server error: Auth service temporarily unavailable”.
- On some terminals, dragging and dropping a file into the interactive client will add single quotes ( ' ') around the file path. This causes the interactive client to misinterpret the file path and generate an error when attempting to load a manifest file or token.
  - *Workaround:* Manually type out the file name or remove the single quotes from around the file path.
- When any files mentioned in the upload manifest are not present in the upload directory the submission will hang at the missing file.
  - *Workaround:* Edit the manifest to specify only the files that are present in the upload directory for submission or copy the missing files into the upload directory.

## v1.2.0

- **GDC Product:** Data Transfer Tool
- **Release Date:** Oct 19th 2016

### New Features and Changes

- Better handling of connectivity interruptions

### Bugs Fixed Since Last Release

- Uploads via manifest file has been fixed.
- Legacy `-i/-identifier` flag removed.
- Improved error messaging when uploading without a token.

### Known Issues and Workarounds

- Use of non-ASCII characters in token passed to Data Transfer Tool will produce incorrect error message “Internal server error: Auth service temporarily unavailable”.
- On some terminals, dragging and dropping a file into the interactive client will add single quotes ( ' ') around the file path. This causes the interactive client to misinterpret the file path and generate an error when attempting to load a manifest file or token.
  - *Workaround:* Manually type out the file name or remove the single quotes from around the file path.
- When any files mentioned in the upload manifest are not present in the upload directory the submission will hang at the missing file.
  - *Workaround:* Edit the manifest to specify only the files that are present in the upload directory for submission or copy the missing files into the upload directory.
- Upload flags `-path/-f` do not modify the upload path as expected.
  - *Workaround:* Copy the Data Transfer Tool into the root of the submittable data directory and run from there.
- Submission manifest field **local\_file\_path:** does not modify upload path expected.
  - *Workaround:* Run Data Transfer Tool from root of the submittable data directory so that data is in the current working directory of the Data Transfer Tool.

## v1.1.0

- **GDC Product:** Data Transfer Tool
- **Release Date:** September 7, 2016

### New Features and Changes

- Partial extension added to all download files created during download. Removed after successful download.
- Number of processes started by default changed to 8 (`-n` flag).

### Bugs Fixed Since Last Release

- None to report.

### Known Issues and Workarounds

- Use of non-ASCII characters in token passed to Data Transfer Tool will produce incorrect error message “Internal server error: Auth service temporarily unavailable”.
- On some terminals, dragging and dropping a file into the interactive client will add single quotes ( ' ') around the file path. This causes the interactive client to misinterpret the file path and generate an error when attempting to load a manifest file or token.
  - *Workaround:* Manually type out the file name or remove the single quotes from around the file path.
- Use of a manifest file for uploads to the Submission Portal will produce an error message “ERROR: global name ‘read\_manifest’ is not defined”.
  - *Workaround:* Upload files via UUID instead or use the API/Submission Portal.

## v1.0.1

- **GDC Product:** Data Transfer Tool
- **Release Date:** June 2, 2016

### New Features and Changes

- MD5 checksum verification of downloaded files.
- BAM index files (.bai) are now automatically downloaded with parent BAM.
- UDT mode included to help improve certain high-speed transfers between the GDC and distant locations.

### Bugs Fixed Since Last Release

- None to report.

### Known Issues and Workarounds

- Use of non-ASCII characters in token passed to Data Transfer Tool will produce incorrect error message “Internal server error: Auth service temporarily unavailable”.
- On some terminals, dragging and dropping a file into the interactive client will add single quotes ( ' ') around the file path. This causes the interactive client to misinterpret the file path and generate an error when attempting to load a manifest file or token.
  - *Workaround:* Manually type out the file name or remove the single quotes from around the file path.

## v1.0.0

- **GDC Product:** Data Transfer Tool
- **Release Date:** May 26, 2016

### New Features and Changes

- Single-thread and multi-threaded download capability
- User-friendly command line interface
- Progress bars provide visual representation of transfer status
- Optional interactive (REPL) mode
- Detailed help menus for upload and download functionality
- Support for authentication using a token file
- Support for authentication using a token string
- Resumption of incomplete uploads and downloads
- Initiation of transfers using manifests
- Initiation of transfers using file UUIDs
- Advanced configuration options
- Binary distributions available for Linux (Ubuntu), OS X, and Windows

### Bugs Fixed Since Last Release

- None to report.

### Known Issues and Workarounds

- Use of non-ASCII characters in token passed to Data Transfer Tool will produce incorrect error message “Internal server error: Auth service temporarily unavailable”.
- On some terminals, dragging and dropping a file into the interactive client will add single quotes ( ' ') around the file path. This causes the interactive client to misinterpret the file path and generate an error when attempting to load a manifest file or token.
  - *Workaround:* Manually type out the file name or remove the single quotes from around the file path.

## Chapter 5

# Data Transfer Tool UI Documentation

### ###Data Downloads with the Data Transfer Tool UI

####Data Transfer Tool UI: Overview The UI version of the Data Transfer Tool was created for users who prefer a graphical interface over the command line or have limited command line experience. The command line version is recommended for those users with more command line experience, require large data transfers of GDC data, or need to download a large numbers of data files.

### ####System Recommendations

The system recommendations for using the GDC Data Transfer Tool are as follows:

- OS: Linux (Ubuntu 14.x or later), OS X (10.9 Mavericks or later), or Windows (7 or later)
- CPU: At least four 64-bit cores, Intel or AMD
- RAM: At least 2 GiB
- Storage: Enterprise-class storage system capable of at least 1 Gb/s (gigabit per second) write throughput and sufficient free space for BAM files.

### ####Binary Distributions

Binary distributions are available on the GDC Transfer Tool page. To install the GDC Data Transfer UI download the respective binary distribution and unzip the distribution's archive to a location on the target system that is easily accessible.

####Binary Installation Once the binary has been positioned in an appropriate location on the client's file system the application will need to run through a one-time installation process. On first execution the binary install splash screen will appear showing the progress of the installation. A hidden directory is created within the user's home directory labeled dtt that holds configuration and executable files.

### ####Preparing for Data Download

The GDC Data Transfer Tool UI is a stand-alone client application intended to work with data file information stored on the GDC Data Portals. Data download information must first be gathered from the GDC Data Portal. From there a manifest file can be generated to supply the client. Alternatively, individual file UUIDs can be provided to the UUID entry window located on the Download tab in the client.

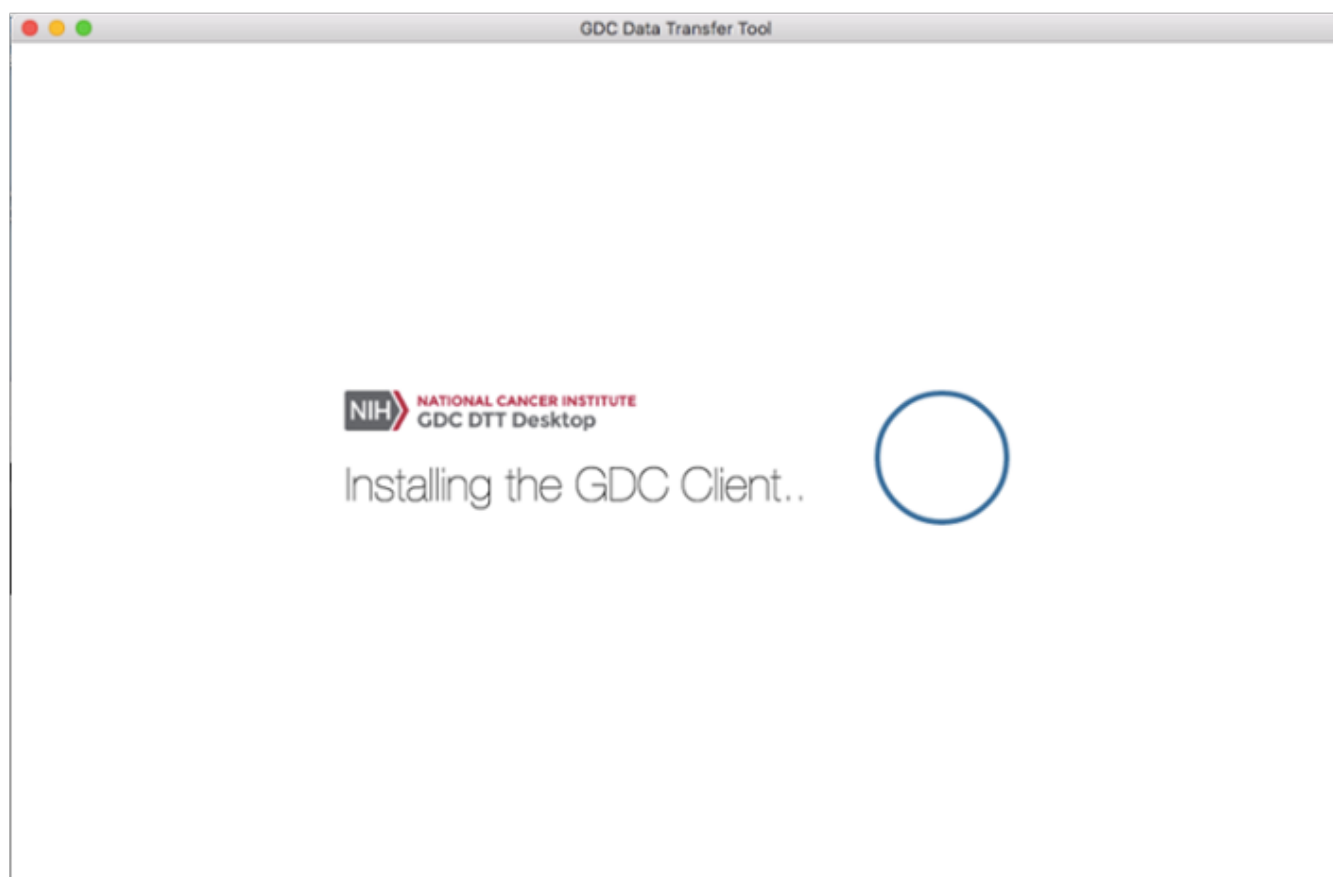
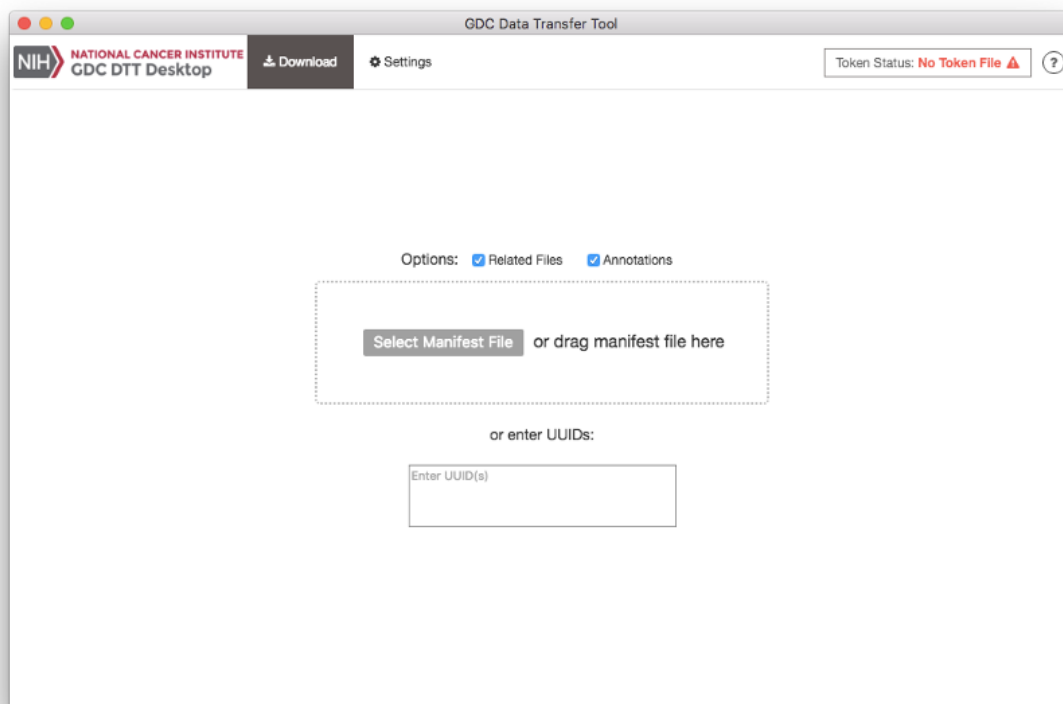


Figure 5.1: GDC DTT UI Installation





####Downloads with UUIDs The Data Transfer Tool UI can download files by individual UUID. UUIDs can be entered into the client while on the download tab. The single entry field labeled “Enter UUID(s)” allows the user to enter UUIDs individually. To obtain a data file’s UUID from the GDC Data Portal, click on the file name to display the file’s summary page which includes vital information such as its GDC UUID.

####Downloads with Manifest A portal-generated manifest file can be used with the Data Transfer Tool UI. From the Download tab home page click on the Select Manifest File button. A file system search window will popup allowing navigation to the manifest file.

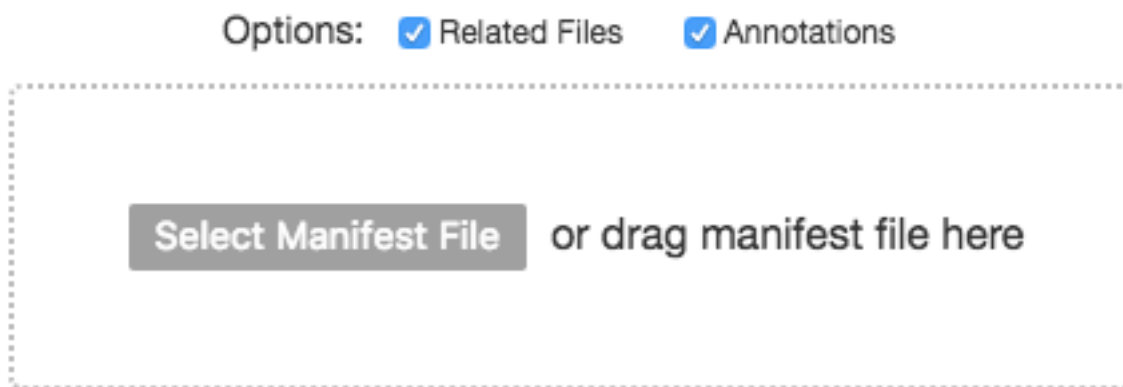


Figure 5.2: GDC DTT UI Manifest Button Example

####Download Progress Page The Download Progress Page is the command console for the Data Transfer Tool UI and allows users to monitor downloads. Progress of all downloads including the ability to start, stop, and restart a download are performed on the Download Progress Page. Once file UUIDs or a manifest has been added to the queue the download can be started by clicking on the download button located at the lower right hand side of the page.

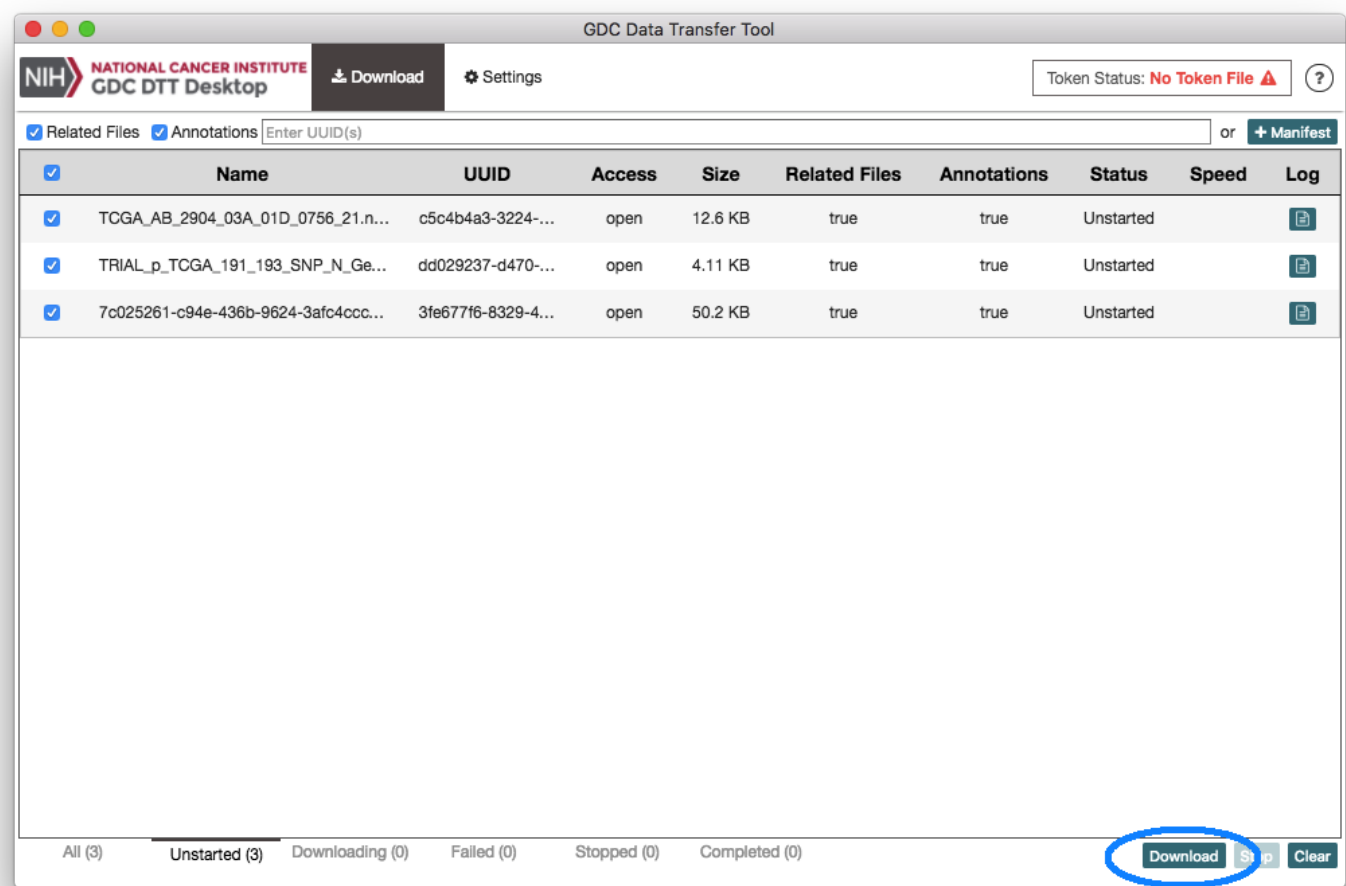


Figure 5.3: GDC DTT UI Download Progress Page\_Download

Once a download has completed, information about the downloads can be viewed from the Completed tab located at the bottom of the page. Any Stopped or Failed downloads can also be viewed from their respective labels located at the bottom of the Status page.



Figure 5.4: DTT\_UI\_Download\_Completed\_Tab

####Controlled Access File Downloads

Some files in the GDC are controlled access. If you require access to these files please review the process outlined in the documentation Obtaining Access to Controlled Data. After appropriate authorization has been granted an access token can be generated to allow the Data Transfer Tool UI application access to the requested data files. Documentation explaining the process of generating a token is located in the Obtaining an Authentication Token for Data Downloads. Once a token has been downloaded to a secure location on the client’s local filesystem the Data Transfer Tool UI can now access it.

The current status of client authorization is viewable in the upper right corner of the application. If the image and wording on



Figure 5.5: No\_Token\_icon

the token manager access button is in red then no valid client token file has been uploaded. To upload a valid token file click on the token status button. The token manager window should appear allowing either a drag and drop token file upload or a file navigation window can be opened to navigate to the file location.

The token manager will verify access and display the projects for which the user has access. To complete the token upload process click on the save button within the Token Manager window.

### ####Settings and Advanced Settings

While the default download options will work for the majority of use cases, there are a vareity of ways to customize or modify the download process within the DTT UI. Details of each of the settings are listed below.

Settings	Details
Server URL	Default: <a href="https://api.gdc.cancer.gov">https://api.gdc.cancer.gov</a>
Number of Client Connections: Default (3)	Number of concurrent client threads
Destination Folder: Default (User's Home Directory)	User selectable download file location
Calculate Inbound Segment and check Md5sum on Restart: Default (On)	Verify previous partial downloaded files via segment check sum
Calculates check sums on previous downloaded files Default (On)	Verify downloaded files via file level check sum
Save Logs: - Download Navigation windows for client downloads	Export download or token log files via drop down and export log button
Debug Logging: Default (Off)	Enable debug level logging for file downloads
Block Size (Bytes): Default (1048576)	HTTP chunk size transfers
Save Interval (Bytes): Default (1000000)	save interval in bytes
Auto Retry: Default (On)	Enables auto retries of failed downloads
Retry(s): Default (5)	Number of retry attempts to download a file after failure
Seconds between Retrys: Default (5)	Number of seconds between retires

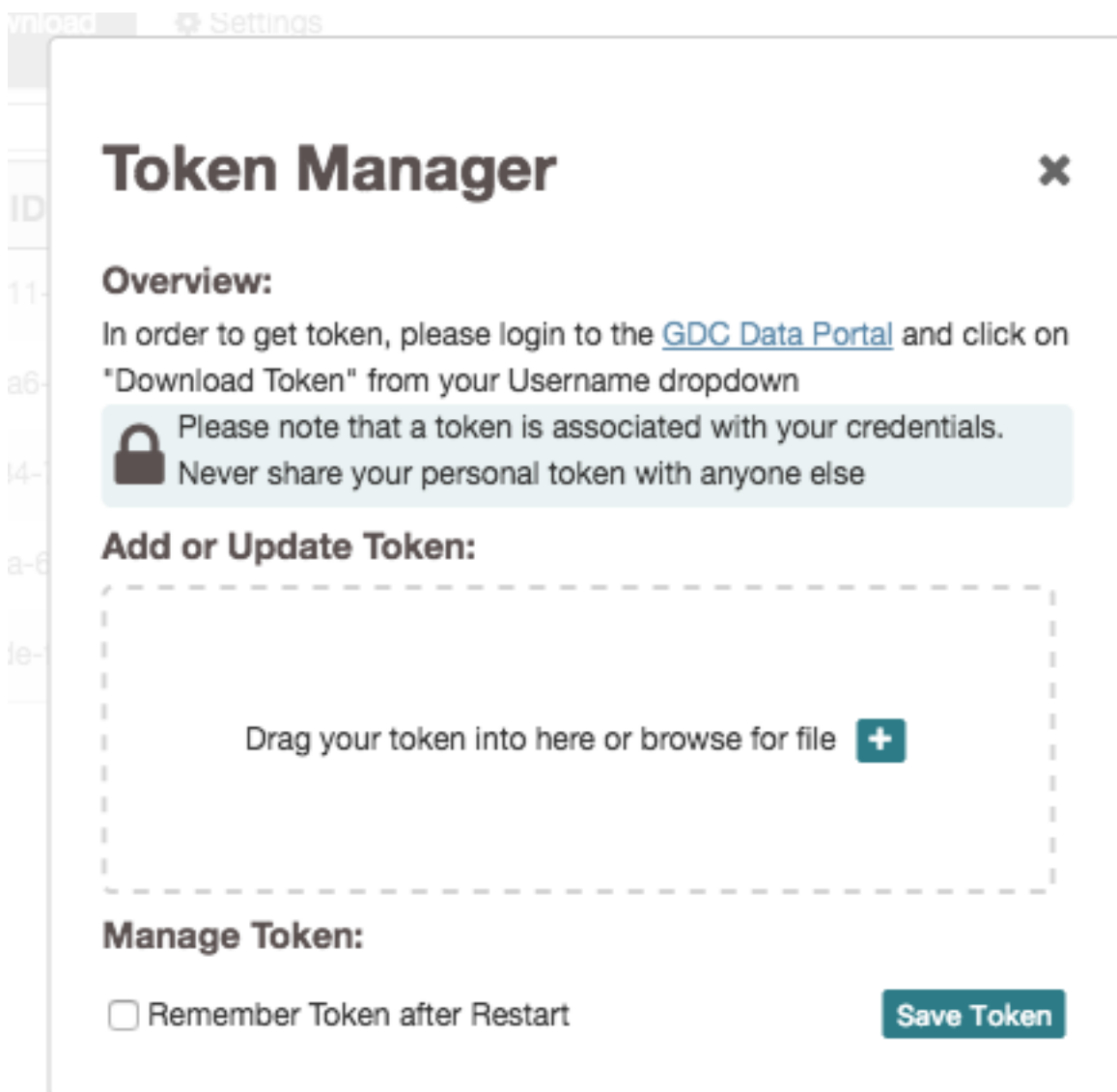


Figure 5.6: Token Manager Window

### Manage Token:

Selected gdc-user-token.2017v16.txt

The current token gives you download access to:

phs001179, CCLEV2, ims000000, phs001134, phs001444,  
phs001449, TEST4, VAL01, phs000892, phs001145, DEV1,  
phs001140, phs001163, phs001160, phs001374, phs000748,  
phs000235, phs001453, phs000178, phs001287, phs000218, TEST,  
gdc000000

☐ Remember Token after Restart

Save Token

Figure 5.7: Valid Token

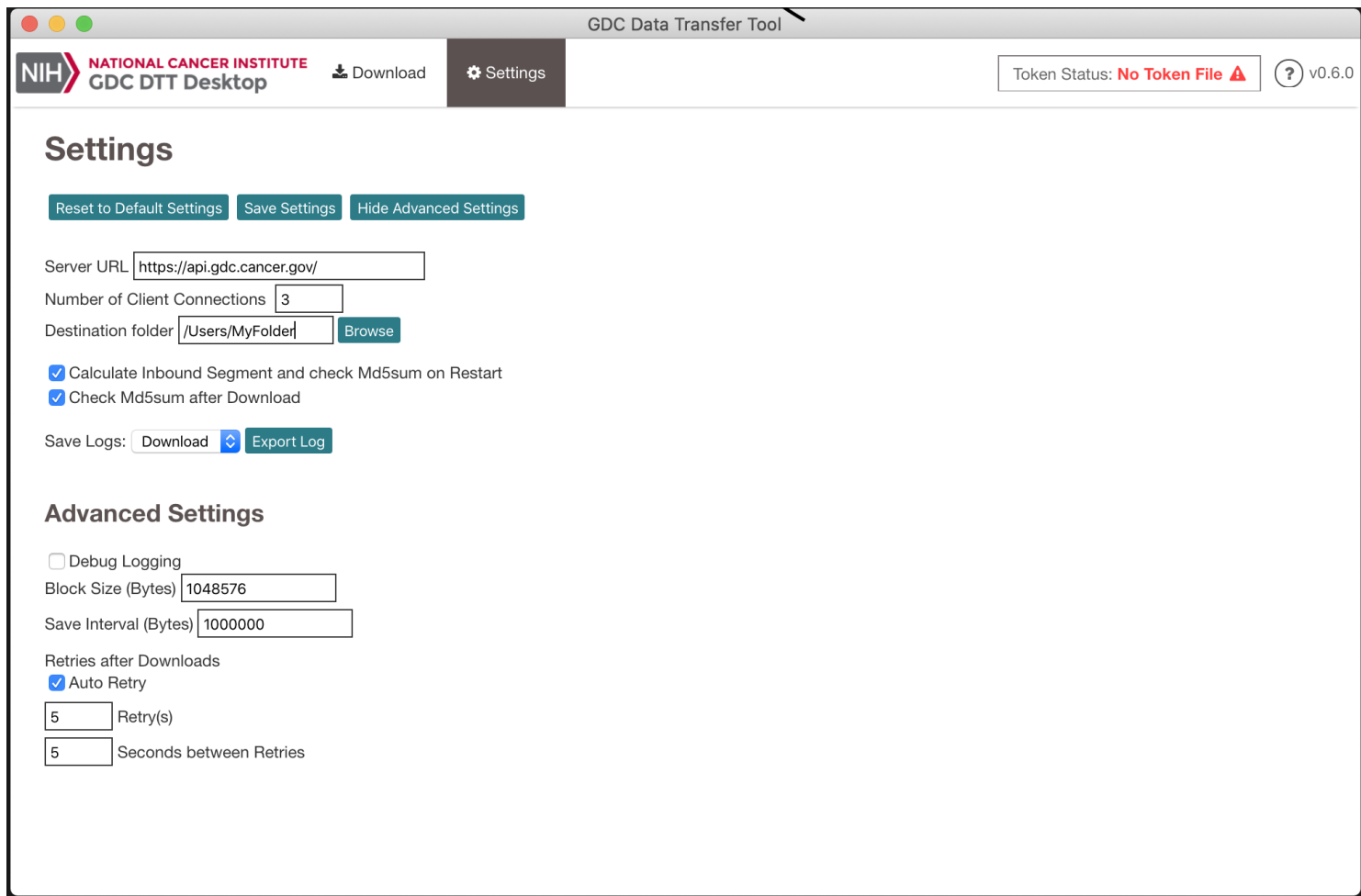


Figure 5.8: DTT Settings and Advanced Settings Page

# Chapter 6

## Release Notes - UI

### Data Transfer Tool UI Release Notes

Version	Date
v1.0.0	July 8, 2021
v0.6.0	August 12, 2020
v0.5.4	April 5, 2018
v0.5.3	December 14, 2017

#### v1.0.0

- **GDC Product:** Data Transfer Tool UI
- **Release Date:** July 8, 2021

#### New Features and Changes

- The DTT-UI now uses the latest version of the Data Transfer Tool (v1.6.1)

#### Bugs Fixed Since Last Release

- None

#### Known Issues and Workarounds

- Download speeds for large numbers of small files may be better handled with the Command Line version of the Data Transfer Tool
- Data Submission to the GDC is not supported in the Data Transfer Tool UI. Instead users must use the Command Line Data Transfer Tool

#### v0.6.0

- **GDC Product:** Data Transfer Tool UI
- **Release Date:** August 12, 2020

#### New Features and Changes

- The server option can now be indicated within the interface. Previously the DTT-UI server defaulted to `api.gdc.cancer.gov`
- The DTT-UI now uses Data Transfer Tool v1.6, which uses Python3.

## Bugs Fixed Since Last Release

- None

## Known Issues and Workarounds

- Download speeds for large numbers of small files may be better handled with the Command Line version of the Data Transfer Tool
- Data Submission to the GDC is not supported in the Data Transfer Tool UI. Instead users must use the Command Line Data Transfer Tool

## v0.5.4

- **GDC Product:** Data Transfer Tool UI
- **Release Date:** April 5, 2018

## New Features and Changes

- None

## Bugs Fixed Since Last Release

- Download is now enabled for GDC reference and publication files.

## Known Issues and Workarounds

- Download speeds for large numbers of small files may be better handled with the Command Line version of the Data Transfer Tool
- Data Submission to the GDC is not supported in the Data Transfer Tool UI. Instead users must use the Command Line Data Transfer Tool

## v0.5.3

- **GDC Product:** Data Transfer Tool UI
- **Release Date:** December 14, 2017

## New Features and Changes

- This is the first release for the Data Transfer Tool User Interface. It allows users to download controlled access data using a simplified point and click interface. This is a beta release and we welcome feedback on the user experience. Important updates compared to the Command Line version include:
  - Upload and store authentication token between sessions
  - Easily view progress on a download manifest as the files are completed
  - View download history

## Bugs Fixed Since Last Release

- None

## Known Issues and Workarounds

- Download speeds for large numbers of small files may be better handled with the Command Line version of the Data Transfer Tool
- Data Submission to the GDC is not supported in the Data Transfer Tool UI. Instead users must use the Command Line Data Transfer Tool