



NATIONAL CANCER INSTITUTE Cancer Research Data Commons

datacommons.cancer.gov

Table of Contents

I. INT	RODUCT	ION		1
II. PR	EREQUIS	ITES		1
III. CR	RDC DATA	A MOD	ELS	2
IV. DO			N	2
V. КЕ VI ст	QUEST S			3 A
VII. SI	ONTINU			
VIII. A	ADDING /	AND M	ANAGING COLLABORATORS	
1.	Obtaini	ing Sub	mission Templates	8
2.	Downlo	bading	Data Dictionary and Submission Templates	11
	2.1	Subm	ssion Templates and Properties	12
		2.1.1	Special Columns	12
		2.1.2	Type Column	13
		2.1.3	Relationship Columns (Parent Mapping Columns)	13
3.	Upload	ing Dat	a Files and Metadata Manifests	14
	3.1	Uploa	der CLI Tool	14
		3.1.1	Introduction	14
	3.2	Down	loading the Uploader CLI Tool	14
		3.2.1	Download the Uploader CLI Tool from the CRDC Submission Portal	14
		3.2.2	Cloning the Uploader CLI Tool from GitHub	15
	3.3	Settin	g Up the Python Environment	16
	3.4	Using	the Uploader CLI Tool	16
		3.4.1	Uploader CLI Tool Configuration File	16
		3.4.2	File Manifest	19
	3.5	Startin	ng the Upload Process	20
	3.6	Using	the CRDC Submission Portal's Graphical Interface to Upload Metadata Submission	
		Temp	ates	20
4.	Runnin	g Valida	ations	23
	4.1	Review	wing Validation Results	23
	4.1.1	1Viewii	ng and Filtering Validation Results	24
	4.21	Revisin	g Released Metadata – Options to Keep or Overwrite Existing Content	25
	4.3	Correc	cting Errors	27
	4.4	Remo	ve Specific Files	27
5.	Submit	ting Yo	ur Final Dataset	28
6.	What to	o Expec	t After Submission	28

I. INTRODUCTION

This tutorial walks you through the process of submitting data to CRDC through the CRDC Submission Portal. If you have questions that are not answered here, contact the Data Concierge assigned to your submission once you successfully create a submission or email the CRDC Help Desk (NCICRDC@mail.nih.gov).

II. PREREQUISITES

Before starting your data submission, complete the following prerequisites:

- Secure approval from the CRDC Submission Review Committee to submit data. Approval/Rejection notification will be found on the portal under Submission Request and in an email sent to the requestor. If rejected, consider other repositories for sharing data at NIH.
- Create a Login.gov account. It is strongly recommended that the Login.gov identity be associated with the submitter's/user's organization or institution; however, it is not a requirement. Using an institutional email as your user identity helps us quickly determine your organization, but you may choose a personal email instead. NIH staff can log in using their PIV card.

Note: If you do not log in to the CRDC Submission Portal account within 60 days, your access to the portal will be deactivated. This applies even while you are working on an active, ongoing submission. To reactivate your access, contact the CRDC Help Desk (NCICRDC@mail.nih.gov).

- **Request Access to the Submitter role:** Users must request the *Submitter* role to submit data to the CRDC Submission Portal or oversee the relevant submissions. See "V." on page 3 for more details.
- If the submission contains controlled access data, the study must be registered at the database of Genotypes and Phenotypes(dbGaP). dbGaP will provide a dbGaP ID/Accession number (phs000####) upon registration. The CRDC Submission Portal will not allow users to submit controlled access data without a dbGaP ID/Accession number.
 - If the dbGaP ID/Accession number was not already shared on the approved CRDC Submission Request form, the user must email the dbGaP ID associated with their study to the CRDC Help Desk (NCICRDC@mail.nih.gov). The submitter can initiate the data submission process only when the dbGaP ID has been provided to CRDC.
 - Note: Data will be released on the respective CRDC Data Commons portal after the study is released to dbGaP. It is therefore recommended to register the study and work on dbGaP submissions concurrently with the submission to CRDC.
- The CRDC Submission Portal uses CRDC standard Common Data Elements (CDEs), and all submissions
 are expected to use these CDEs and comply with their permissible values. A comprehensive list of
 CRDC standard CDEs can be found at caDSR. It is recommended that submitters familiarize
 themselves with the CRDC standards before starting a submission. On the caDSR website, click the
 CRDC Standard Data Elements link in the Links to Favorites section or download them from the
 getCRDCList endpoint of the caDSR API.

III. CRDC DATA MODELS

CRDC's Data Commons use data models to organize data in a consistent and structured manner, ensuring accuracy and facilitating reusability. These data models are graph-based, and data are organized as nodes and relationships. Nodes contain properties and can have relationships with other nodes. Nodes are equivalent to tables in a relational model, and properties are equivalent to columns in a relational model. Relationships serve a similar purpose as foreign keys in a relational model. For example, in the GC Data Model, the Participant node is a child of the Study node and a parent of the Diagnosis, Treatment, Sample, and File nodes. See Figure 1.



Figure 1. Participant Node Relationships in the GC Data Model

IV. DOCUMENTATION

Submitters can find this document, as well as instructions on using APIs to submit data under the **Documentation** tab.



Figure 2. Documentation Menu Showing Available Documents

V. REQUEST SUBMITTER ROLE

- 1. Log in to the CRDC Submission Portal and go to your User Profile by clicking your name in the upperright corner of the page.
- 2. By default, the account has the User role. To request the Submitter role, click **Request Access**. See Figure 3.

		@nih.gov	
	Account Type	NIH	
	Email		
	First name*	Demo	
	Last name*	Account	
	Role	Submitter REQUEST ACCESS	
	Institution	NCI	
	Studies	CPTAC and other 4	
	Account Status	Active	
scions (8)			

Figure 3. Request Access

A **Request Access Form** appears (see Figure 4) where you provide the required information.

Role*	
Submitter	\$
Institution	•
100 characters allowed	Ç
Studies*	
Select one or more studies from t	the list 🗘
Additional Info	

Figure 4. Request Access Form

3. From the **Role** dropdown menu, select **Submitter** if you want to work on the data submission or oversee the submissions for specific studies.

- 4. From the **Institution** dropdown menu, select your Institution's name. If not listed, enter it manually. In the **Studies** dropdown menu, select the relevant studies from the list.
- 5. Optionally, Users can provide additional details about their role in the **Additional Info** box before submitting the form.

The System Administrator grants the Submitter Role, and the user is notified via their login identity. The user granted with the **Submitter** role can now start the data submission process.

Note: If you do not log in with the CRDC Submission Portal account within 60 days, your access to the portal will be deactivated. If that happens, contact the CRDC Help Desk (NCICRDC@mail.nih.gov) to reactivate your account.

VI. STARTING A NEW SUBMISSION

Once logged in with a Submitter role, navigate to the Data Submissions tab on the CRDC Data Submission portal. The submitter is taken to the Data Submission List page. If this is the Submitter's first data submission, the table showing the list of Data Submissions will be empty. If the Submitter has multiple submissions, use the filters at the top of the table to narrow down the list. See Figure 5.

	TIONAL CA	NCER INS	TITUTE ata Comm	ions								
Back to CRDC	Submis	sion Requ	ests Dat	a Subm	issions	Documenta	tion ~	Model Navigo	itor ~			DEMO ~
Data S	ubmi	ccior	List									•
Below is a list of d Please click on an	ata submissi y of the data	ons that are submission	associated w s to review of	vith your r continue	account. e work.					Create	a Data Submis	sion
Program	All		٥)	Status	6 statuses se	elected	\$	Data Commons	All	\$	C
Submission Name	Minimum 3 characters requi		3 characters required		dbGaP ID	Minimum 3 c	haracters	required	Submitter	All	٥	-
Submission Name	Submitter	Data Commons	Туре	DM Version	Program	Study	dbGaP II	D Status	Primary Contact	Record Count	Created Date	Last Updated ↓
Test1 CDS	Demo Accou	CDS	New/Update	<u>v6.0.2</u>	NCI	2354423	676	New		0	2/24/2025	2/24/2025
test-stats		CDS	New/Update	<u>v6.0.2</u>	NA	001TS		In Progress	mattuak+4	3	2/24/2025	2/24/2025

Figure 5. Create a Data Submission

To start a new data submission, click the **Create a Data Submission** button and a dialog box as shown in Figure 6 appears. Fill out all the required information as described below.

Submission Type* 🔘 New	/Update O Delete	
Data Type*	idata and Data Files 🔵 Metadata Only	
Data Commons*		
CDS	\$	
Study *		
CONTROLLED1	Please contact NCICRDC@mail.	nih.gov
dbGaP ID*	your dbGaP ID once you have re study on dbGap.	egistere
<not provided=""></not>	4	
Submission Name*		

Figure 6. Data Submission Dialog Box

- 1. Choose the Submission Type.
 - Select New/Update to create a new data submission or update an existing one.
 - Select the **Delete** option to remove files from a previous submission already released publicly by the Data Commons. Selecting the Delete option keeps only one Data Type option enabled, which would be **Metadata Only.** Submit the metadata associated with the data that needs to be deleted. The deletion request goes through a validation process by the CRDC Team and the CRDC Data Commons. Once approved by the CRDC Data Commons, the deletion is processed.
- 2. For Data Type, indicate whether you are submitting both metadata and data files or only metadata files by selecting one of the options, **Metadata and Data Files** or **Metadata Only**, respectively.
- 3. Select the **Data Commons** that you were approved to submit to by the Submission Review Committee (SRC) for your data, if it is not preselected already. If you do not see the CRDC Data Commons listed, email the CRDC Help Desk (NCICRDC@mail.nih.gov).
- 4. The **Study** dropdown menu displays the Study Title (or the Study Abbreviation) you previously shared through the Submission Request form. If you notice an error in this list, email the CRDC Help Desk (NCICRDC@mail.nih.gov).
- 5. If your study includes controlled access data, the **dbGaP ID** /accession number will be pre-populated as provided on the Submission Request Form. If you did not provide the dbGaP ID on the Submission Request Form, please email it to the CRDC Help Desk. The system will not allow the submitter to initiate a data submission if the submitter has not shared the dbGaP ID.
- 6. Submitters can give the submission a **Name** in the provided free-text field to label their submissions. This name appears in the Submissions List table on the Data Submissions List page, once the submission is created.
- 7. Click the **Create** button to create the new submission. This button will only be enabled if the dbGaP ID is provided.

Once a submission is created, CRDC assigns a Data Concierge to your submission. You can find the email address of the Data Concierge assigned to your submission on the dashboard of your data submissions page. From this step onwards, all questions related to the data submission should be directed to the assigned Data Concierge.

VII. CONTINUING AN EXISTING SUBMISSION

To access and update an existing submission, go to the **Data Submissions** tab. A table, listing all the existing submissions, appears on that page. Under the **Submission Name** column, select the submission you want to continue with. To see the full version of the submission name, hover over the name (see Figure 7). The submitters can customize the columns displayed in the submission table. Desired columns can be selected and added by clicking the **Table** icon in the top-right corner and applying the changes. The filters at the top of the table provide useful ways to refine searches, especially when the list of submissions is extensive.

	TIONAL CA	ANCER INS search Da	TITUTE Ita Comm	ons									
Back to CRDC	Submis	sion Reque	ests Date	a Submi	ssions D	ocumento	ation ~	Model Navig	ator ~				DEMO ~
													0
Data S	ubmi	ssior	ı List										
Below is a list of d Please click on an	ata submissi y of the data	ons that are submissions	associated w s to review or	vith your a continue	account. work.					[Create a	Data Submis	sion
										l			
Program	All		\$]	Status	6 statuses se	elected	\$	Data Commons	All		3	0
Submission Name	Minimum	3 characters i	required]	dbGaP ID	Minimum 3 c	characters r	required	Submitter	All		\$	
Submission Name	Submitter	Data Commons	Туре	Model Version	Program	Study	dbGaP ID	Status	Data Concierge	Record Count	Data File Size	Created Date	Last Updated ↓
Test 2	Demo Accou	GC	New/Update	<u>v8.0.0</u>	CPTAC	CPTAC	phs001287	/ New		0	0	5/29/2025	5/29/2025
test-2363		GC	New/Update	<u>v8.0.0</u>	DCCPS CIDF	CIDR: Disc	phs002250) New		0	0	5/28/2025	5/29/2025
Test for d	Demo Accou	GC	New/Update	<u>v8.0.0</u>	CPTAC	CPTAC	phs001287	V New		0	0	5/29/2025	5/29/2025

Figure 7. Active Data Submissions List

VIII. ADDING AND MANAGING COLLABORATORS

The submitter can add collaborators to a given submission and manage their access so that the collaborators can also upload and validate data through the CRDC Submission Portal. By default, the collaborators count is set to zero. As collaborators are added, the count is updated on the dashboard. See Figure 8.

NIH Cance	NAL CANCER INSTITUTE er Research Data Comr	nons		
Back to CRDC	Submission Requests	Data Submissions Documentation	Model Navigator \sim	DEMO ~
	232f456e-6f63-4bcf-8c0	5-15fcecd6a24d		
Revie	New ww Comments	SUBMISSION NAME Test 2 SUBMITTER Demo Account STUDY CPTAC PROGRAM NA	SUBMISSION TYPE New/Update COLLABORATORS Q DATA COMMONS GC DATA CONCIERGE	

Figure 8. Count of Collaborators on the Data Submission Dashboard

To add and manage collaborators, click the hyperlink next to **Collaborators** on the Submission dashboard, which opens a new Data Submission Collaborators window. See Figure 9.

Ensure that the collaborator has an authorized account on the CRDC Submission Portal with the Submitter role, is affiliated with the same study, and has permissions from the data owner or Study PI(s) to submit data. (see Section V. Request a Submitter Role). In the Data Submission Collaborators window, the submitter can select collaborator/s by selecting names from the dropdown list.

The submitter can **remove** a collaborator by clicking the remove icon [X] and add multiple collaborators by clicking the **Add Collaborator** button and repeating the process for each additional collaborator. Be sure to click **Save** before closing the window to retain any changes.

Below is a list of collab	orators who have be	en granted access to	this data s	ubmission.
Ince added, each colla unning validations, and	borator can contribu d submitting.	te to the submission	by uploadir	ig data,
Collaborator				Remove
Select Name			\$	×
Add Collaborator				
+ Add Collaborator				

Figure 9. Add and Remove Collaborators for a Given Submission

1. Obtaining Submission Templates

Submitting data to CRDC requires you to submit your metadata using the submission templates. The metadata is then used to validate the submitted actual raw data. For instance, the file name and the file size of each uploaded raw data file will be compared with the file name and file size specified in the metadata manifest.

To get to the submission templates, click **Model Navigator** in the menu bar (see Figure 10), which lists the Data Models of the various Data Commons integrated under the CRDC Data Submission Portal. . The submission templates for various Data Commons, including the General Commons (GC) Model, Clinical and Translational Data Commons (CTDC) Model, and Integrated Canine Data Commons (ICDC) Model, are provided. Models of other Data Commons will be added when they are integrated with CRDC Submission Portal. Select the Data Model respective to the Data Commons (DC) to which the Submission Review Committee (SRC) has approved you to submit data.



Figure 10. Use the Menu Bar to Navigate to the Data Model Viewer

Once you select the Data Model, you are taken to the Data Model Viewer page, as seen in Figure 11. On this page, you can view the data model in detail and download the submission templates from the dropdown list of Available Downloads. **Note:** Figure 11 shows the GC Data Model as an example.

	IATIONAL CANCER INSTITUTE Cancer Research Data Commons	
Back to CRD	DC Submission Requests Data Submissions Documentation	DEMO ~
GC Data Model		(README ?) (+ Available Downloads)
Filter & Search	Graph View Table View Version History	
Search in Dictionary Q	E program	Node Category 🗢
Filter By Nodes Category		ing study ing cove ing data file
 Assignment Class 	diagnosis	aa)
Filter By Property		
✓ Inclusion		
 UI Display 	genomic_inite image ima	

Figure 11. Data Model Viewer Graph View

Use the Data Model Viewer to explore the data elements that the Data Commons require or can accept. On the **Graph View** tab, the data model is represented in graphical nodes and relationships. Clicking a node in the graph shows its summary. At the bottom of the node summary, click the **View Properties** to open a **Table View** of the selected node. For instance, as shown in Figure 12 clicking on the **Diagnosis** node opens its summary, with View Properties option at the bottom.



Figure 12. Click a Node to View a Summary and Open the Table View

The **Table View** lists all the data elements/properties of the data model and includes the description of each of these properties (such as strings, integers, etc.) The Table View also shows which of these properties are required for a submission. See Figure 13. Each of these properties is mapped to the Common Data Elements (CDEs) of the caDSR standards where applicable. Please note that CRDC data validations will only accept **Permissible Values** (PVs) for elements mapped to the **Common Data Elements** (**CDEs**). Details about the CDEs can be accessed by clicking on the Public ID for that specific property. If you do not find the correct match in the provided values, you can request new CDEs and/or new Permissible Values by emailing the CRDC Help Desk (NCICRDC@mail.nih.gov). These requests will need to go through a CRDC approval process before you can use the new CDEs or PVs in a submission.

Additionally, in the **Table View**, the **Submission Template** and the associated **Data Dictionary** for the specific node can be downloaded as shown in Figure 13 and described in the next section.

Graph V	iew Table View	Version History			
Diagnosis	Text term used to de Diseases for Oncology	scribe the patient's histologic c (ICD-O).	liagnosis, as Assignment: Co	described by the World Health Organization's (WHO) International Classificatio	on of ary U
Property	Туре	CDE Info	Required	Description	Sour
study_diagnosis_id 🗪	"string"		Required	The property study_diagnosis_id is a compound property, combining the property diagnosis_id, string character "_" in the middle as the connector, and the parent property participant.study_participant_id. It is the ID property for the node diagnosis.	
diagnosis_id	"string"		Required	Internal identifier	
disease type	Acceptable Values: • Epithelial Neoplasms, NOS • Transitional Cell Papillomas and Carcinomas • Soft Tissue Tumors and Sarcomas, NOS • Paragangliomas and Glomus Tumors • Osseous and Chondromatous Neoplasms • Nevi and Melanomas • Nevi and Melanomas • Nevi and Melanomas • Nevi and Melanomas • Nevi, Mucinous and Serous Neoplasms • Granular Cell Tumors and Alveolar Soft Part Sarcomas …show more	CDE Full Name Diagnosis Disease or Disorder Morphology Disease Version 1.00 Public ID 13471160 Origin caDSR	Optional	Type of disease [?]	

Figure 13. Table View of an Excerpt of Diagnosis Node

Additionally, any updates to the data model, such as property modifications, additions, or updates to permissible values, are reflected in the **Version History** tab, ensuring users have access to the changes. See Figure 14. Currently Version History is applicable and available for the GC data model only.

Graph View Table View Version History
6.0.2 (Released 2/18/2025)
 Updated the version-history.md. Updated the relationship between the node: "<i>proteomic</i>" and the node: "<i>file</i>" to be many-to-one.
6.0.1 (Released 2/13/2025)
• Renaming the prop: " <i>gender</i> " to " <i>sex</i> ".
6.0.0 (Released 1/2/2025)
 Added the prop: "is_supplementary_file" to the node: "file".
 Added the prop: "release_datetime" to the node: "file".
 Removed the prop: "library_source" from the node: "file".
 Updated the description for the prop: "<i>file_id</i>".
 Updated the prop: "study_data_types"'s type to list.

Figure 14. Excerpt of Version History of the GC Data Model

2. Downloading Data Dictionary and Submission Templates

Each Data Commons has its own unique data model with corresponding data dictionary, set of nodes, properties, and submission templates. To download the Data Dictionary and the submission templates, select **Available Downloads**, as shown in Figure 15. Then click one of the options from the dropdown list to download the file in the selected format. You can also download CRDC Vocabularies for the selected Data Commons Model and Example Templates.

- Data Dictionary
 - All Properties (PDF, TSV, JSON)
 - Required Properties (PDF, TSV, JSON)
- Submission Templates (TSV)
- All Vocabularies (TSV, JSON)
- Example Templates

Submitters can use the **Submission Templates** to format and upload metadata to the CRDC Submission Portal. These templates <u>must</u> be in TSV format. Templates in any other format, such as Microsoft Excel (.xls, .xlsx) etc., will fail. You can use software such as <u>ModernCSV</u> to work with these submission templates, as it handles CSV and TSV as tables without automatically modifying the data.

The **Data Dictionary** provides detailed information about the metadata structure, content, and the required, preferred, and optional data elements for all nodes within the selected data model. Submitters can choose to download the Data Dictionary for either all properties or only the required properties. The **All Vocabularies** document contains the permissible values for the data elements. The **Example Templates** are examples of completed submission templates with mock data, designed to guide users in preparing the metadata manifest for their data. These can be useful to understand what each of the columns in the template is supposed to contain.



Figure 15. Using the Available Downloads Menu

The downloaded files are provided as a ZIP archive. The tab-separated text files can be viewed in any text editor or spreadsheet application like Microsoft Excel or OpenOffice Calc. Multiple metadata template files in TSV format are included within the ZIP archive called Submission Templates. The downloaded and unzipped Submission Templates folder for the GC model is shown in Figure 16.

Note: The exact content of the submission templates differs depending on the selected data model and the associated submission process requirements.

✓ GC_Data_Loading_Templates_v8.0.0
GC_Data_Loading_Template_pdx_v8.0.0.tsv
GC_Data_Loading_Template_proteomic_v8.0.0.tsv
GC_Data_Loading_Template_NonDICOMradiologyAllModalities_v8.0.0.tsv
GC_Data_Loading_Template_NonDICOMPETimages_v8.0.0.tsv
GC_Data_Loading_Template_NonDICOMpathologyImages_v8.0.0.tsv
GC_Data_Loading_Template_NonDICOMMRimages_v8.0.0.tsv
GC_Data_Loading_Template_NonDICOMCTimages_v8.0.0.tsv
GC_Data_Loading_Template_MultiplexMicroscopy_v8.0.0.tsv
GC_Data_Loading_Template_image_v8.0.0.tsv
GC_Data_Loading_Template_genomic_info_v8.0.0.tsv
GC_Data_Loading_Template_file_v8.0.0.tsv
GC_Data_Loading_Template_sample_v8.0.0.tsv
GC_Data_Loading_Template_treatment_v8.0.0.tsv
GC_Data_Loading_Template_diagnosis_v8.0.0.tsv
GC_Data_Loading_Template_participant_v8.0.0.tsv
GC_Data_Loading_Template_study_v8.0.0.tsv
GC_Data_Loading_Template_program_v8.0.0.tsv

Figure 16. Submission Templates Downloaded from the CRDC Submission Portal Model Viewer

2.1 Submission Templates and Properties

Each of the submission templates covers information relevant to a specific node in the model; for example, the template, 'GC_Data_Loading_Template_image_v8.0.0.tsv,' collects imaging data related information. Not all templates are required; rather, only those templates relevant to the data being submitted are required. For instance, if the submission does not include imaging data, the submitter does not need to fill or submit the 'GC_Data_Loading_Template_image_v8.0.0.tsv.'

For every template that will be submitted, review the Data Dictionary (accessible through the **Available Downloads** menu) to understand which properties are required, as each individual template has required properties, as well as preferred and optional properties in the node. Note that you should not edit the first row of each template as it contains the property names and other special columns, explained below.

2.1.1 Special Columns

Every template has two types of special columns, also called parent mapping columns: "type" and "relationship."

2.1.2 Type Column

A "type" column contains the name of the node type (such as study or genomic_info) and is required by all template files. In the downloaded template, the second row in the first column is prefilled with the correct node name for that specific template (e.g., in the 'GC_Data_Loading_Template_study_v8.0.0.tsv' template, the second row in the first column is filled in as 'Study.'). All rows should contain the same node name in the "type" column. Mixing multiple node types in one file is not supported. For example, the node type 'Sample' should not be mixed with the node type 'File.'

2.1.3 Relationship Columns (Parent Mapping Columns)

A "relationship" column is used to specify relationships between the current node and its related nodes. A relationship column has a header in the form of "<parent node name>.<parent ID property name>." Values in the relationship columns are IDs of the related nodes (like a foreign key in a relational model).

For example, for the study node, the "program.program_acronym" column indicates that the study node has "program" node as its parent node, and the property used to identify the program node is "program_acronym." Each value in the "program.program_acronym" column is an acronym used for a program, such as HTAN.

3. Uploading Data Files and Metadata Manifests

You can move files from your local environment to the CRDC through the Submission Portal in the following two ways:

- **Uploader CLI Tool** This command-line interface is used to transfer primary data files like genomic sequence files or imaging data files to CRDC.
- **Graphical interface** The graphical interface can be used to upload metadata files such as the Submission Templates.

Note: Submit primary data files using the Uploader CLI Tool only. Do not attempt to upload data files using the CRDC Submission Portal's graphical interface.

3.1 Uploader CLI Tool

3.1.1 Introduction

The CRDC Submission Portal provides a command-line interface (CLI), the Uploader CLI Tool, for uploading data to its temporary CRDC storage. You can install and use the Uploader CLI Tool on any system capable of running Python 3.6 or higher. Binary versions of the Uploader CLI Tool are also available, which don't require any installation or Python.

Notes:

- There are detailed instructions on downloading, installing, and running the Uploader CLI Tool in the README file of the <u>GitHub repository</u>.
- The Uploader CLI Tool does not have to be downloaded for each submission; this is a Python script that can be used for any upload to the CRDC Submission Portal. The only aspect that must be tailored to each submission is the configuration file, which is discussed below. However, submitters should ensure that they are using the **latest version** of the Uploader CLI Tool and the configuration file.

3.2 Downloading the Uploader CLI Tool

You can download the Uploader CLI Tool either directly from the CRDC Submission Portal or by cloning the GitHub repository. However, downloading from the CRDC Submission Portal is recommended as it ensures you are using the latest version. Starting from version 4.0, CLI tool will automatically check if your CLI tool is compatible with CRDC Submission Portal backend, and prompt you to download a new version when it is no longer compatible.

3.2.1 Download the Uploader CLI Tool from the CRDC Submission Portal

Open the menu by clicking your user profile name, found in the upper-right corner of the Data Submission page. See Figure 17. Select **Uploader CLI Tool** from the menu to open a pop-up window with the latest version of the CLI tool.

		AL CANCER INSTITUTE Research Data Comn	ons					
Вс	ack to CRDC S	Submission Requests	Data Submissions Docume	entation ~	Model Navigator ~			DEMO ^
	User Profile	Uploader CLI Tool	API Token Logout					
						STANS -		094 <i>248</i> 8

Figure 17. Menu with the Uploader CLI Tool Download Option

Click the **Download** icon next to the available Package Type options to download the Uploader CLI Tool. The download comes with accompanying instructions (see Figure 18). A ZIP archive will be saved to your local machine.

ne Uploader CLI Ibmission files fr	is a command-line om your workstatio	interface tool designed for directly uploading dat n to the CRDC Submission Portal cloud storage.
o download the t	ool and access the	accompanying instructions, please choose from
e available dowr	nload options belov	ν.
Package Type	Platform	Download
Source	Any	<u>crdc-datahub-cli-uploader-src.zip</u> ↓
Binary	Windows x64	crdc-datahub-cli-uploader-windows.zip
Binary	MacOS x64	crdc-datahub-cli-uploader-mac-x64.zip
Diridiry	MacOS ARM	<u>crdc-datahub-cli-uploader-mac-arm.zip</u> ↓
Binary	MacOS ARM	

Figure 18. Download the Uploader CLI Tool

3.2.2 Cloning the Uploader CLI Tool from GitHub

The latest version of the Uploader CLI Tool can also be cloned from the Data Hub <u>GitHub repository</u> (see Figure 19). To clone the repository to your local machine, use the following command:

```
git clone --recurse-submodules
https://github.com/CBIIT/crdc-datahub-cli-uploader.git
```

🐉 master 👻 🐉 49 Branches 🛇 44 Tags	Q Go to file t Add file -	<> Code -	About
S vshand11 Merge pull request #53 from CBI	T/3.1.0 🚥 0f02167 · last month	🕚 149 Commits	CRDC datahub data uploader is a command line interface for end use
igithub/workflows	Build the windows binary on windows	4 months ago	upload cancer research files and metadata to Datahub.
Configs	Updated example config files base on QA feedback.	7 months ago	🖾 Readme
src	Fixed	last month	-∿ Activity
🗋 .gitignore	update	6 months ago	E Custom properties
] .gitmodules	Init the repo with all necessary files	2 years ago	③ 15 watching
README-technical.md	complete	8 months ago	父 2 forks
C README.md	Updated README to include correct binary commands an	4 months ago	
requirements.txt	update requirenebts.txt	2 years ago	Releases 15

Figure 19. Uploader CLI Tool as it Appears in GitHub

3.3 Setting Up the Python Environment

Binary versions of the Uploader CLI Tool are self-contained and don't require Python or installation of any dependencies. The Source version of Uploader CLI Tool has Python library dependencies that you must install before running the CLI. These dependencies can be installed by running the command pip3 install -r requirements.txt. The requirements.txt file contains the list of dependencies, described below. If you want to install the dependencies individually, install the following libraries:

- pyyaml
- boto3
- requests
- requests_aws4auth
- Rich
- pandas

3.4 Using the Uploader CLI Tool

3.4.1 Uploader CLI Tool Configuration File

The behavior of the **Uploader CLI Tool** is controlled by the configuration file. You can directly download the configuration file from the CRDC Submission Portal by clicking the Download Configuration File button, which is shown in Figure 20. You can also access the **Data Submission Instructions** document, the version of the **Data Model** your submission utilizes, and the **Uploader CLI Tool** on this same page.

No files selected		Upload	
1011 v14 0			
ION: <u>V4.0</u>			
files to CRDC Submission Portal and requires a co d however the configuration file needs to be custo yample configuration files found in the CLI Tool	onfiguration file to work. omized for each	Download Configuration File	ו
load a configuration file customized for this submi	ission.		J -
	~		
e n	files to CRDC Submission Portal and requires a co ad however the configuration file needs to be cust example configuration files found in the <u>CLI Tool o</u> nload a configuration file customized for this subm	files to CRDC Submission Portal and requires a configuration file to work. ad however the configuration file needs to be customized for each example configuration files found in the <u>CLI Tool download</u> [2], or you can nload a configuration file customized for this submission.	In files to CRDC Submission Portal and requires a configuration file to work. In ad however the configuration file needs to be customized for each example configuration files found in the <u>CLI Tool download</u> [2], or you can nload a configuration file customized for this submission.

Figure 20. Download Configuration File

After you click the **Download Configuration File** button, the **Download Configuration File** pop-up window appears, as shown in Figure 21.

ull Path to Data Files Folder * 🛡
Users/me/my-data-files-folder
ull Path to Manifest File * ♥ /Users/me/my-metadata-folder/my-file-manifest.tsv

Figure 21. Dependencies to Download Configuration File

Enter the path to the local or S3 folder containing your data files and metadata manifest file (explained in Section 3.4.2) in the designated text boxes.

Note: If you enter an S3 URL in the **full path to Data Files folde**r, the Uploader CLI Tool initiates an S3-to-S3 transfer.

Clicking the **Download** button will download the configuration file in YML format to your computer with pre-populated fields. Please note that if your data files are organized in a nested folder structure, the CLI tool will automatically rename the data files to avoid filename collision. To correctly upload data files from nested folder, provide the relative path to each data file in the file_name property of the file manifest.

For example, also shown in Figure 22, if data files are stored in a folder named "data files" and the file is located at data files/folder1/folder2/abc.bam, then you should enter "folder1/folder2/abc.bam" in the file_name column – not just abc.bam.

As a result, the uploaded file will be renamed to folder1_folder2_abc.bam in the CRDC storage. However, the file_name property in the manifest will still reflect the original file name, abc.bam.



Figure 22. CLI Tool Requires Relative path to Data File in a Nested Directory

Additionally, if you choose to populate the configuration file manually, you can find examples in the *configs* directory of either the extracted ZIP file or the cloned GitHub repository. The examples provided are the same configuration file modified for the two different upload types.

- Uploader-metadata-config.example.yml This file is an example of the configuration file needed by the Uploader CLI Tool to upload metadata submission templates rather than submitting them via the CRDC Submission Portal graphical interface.
- **Uploader-file-config.example.yml** This is an example of the configuration file needed by the Uploader CLI Tool for uploading large primary data files such as BAM files. Files uploaded this way go through the file validation system rather than the metadata validation system.

These configuration files are in YAML format and the Uploader CLI Tool will fail if the file is not a valid YAML. YAML-aware text editors such as Microsoft Visual Studio Code, Sublime Text, or Notepad++ can be extremely helpful in preserving YAML formatting. The fields in this file follow.

- **api-url** This field provides the Uploader CLI Tool with the URL/location of the temporary CRDC storage used for API communications and upload.
- token This is the API access token that is obtained from the CRDC Submission Portal's graphical interface. To obtain an API token, log into the CRDC Submission Portal graphical interface to bring up the user menu, then select API Token. This opens a dialog box that allows you to create and copy an API token to your clipboard.
- submission This is the submission ID that identifies which study the uploaded files will be associated with. To find the correct submission ID, log into the system and select the study from the Data Submissions List by clicking the submission name. You can copy the Submission ID from the upper-left corner of the interface by clicking the icon to the right of the Submission ID number.
 Note: A study consists of one or more submissions (often many more), with each Submission ID linked to the parent study. A single user working on multiple studies must carefully track which Submission IDs they are uploading to ensure the data is associated with the correct study.
- **type** This tells the system if this is a metadata upload or a data file upload. Enter the term *metadata* if the upload contains submission templates and *file* if the upload contains data files.
- **data** This is the local path to the directory that contains the files to be uploaded.
- manifest (Data file upload only) This is the local path to the manifest file.
- **retries** This is the number of retries the Uploader CLI Tool will perform after a failed upload.

- **overwrite** If this is set to *true*, the Uploader CLI Tool overwrites the file with the same name that already exists in the CRDC Submission Portal target storage. If set to *false*, the Uploader CLI tool does not upload if a file with the same name and size exists in the CRDC Submission Portal target storage.
- **dryrun** If this is set to *true*, the Uploader CLI Tool does not upload any files to the CRDC Submission Portal target storage. If set to *false*, CLI uploads files to the CRDC Submission Portal target storage.

While users are expected to provide paths to their data folder and manifest file, they may choose to customize the values of the three parameters—**retries**, **overwrite**, and **dryrun**—to suit their needs.

3.4.2 File Manifest

The Uploader CLI Tool uses a document called a file manifest to upload the data files to the temporary CRDC storage. The file manifest is a simple table (a TSV file) with all the required properties as defined by the data model except the file IDs, which are generated by the Uploader CLI Tool. Submitters can use the file.tsv template downloaded from the Data Model viewer page, to create this file manifest, saving the effort of creating a duplicate file.

Instructions for using file.tsv template with CLI tool to upload data files:

- 1. Download the file.tsv template from the Data Model Viewer page. The file.tsv does not include a column for file IDs/Keys because the IDs are generated by the Uploader CLI tool while uploading the data.
- 2. Also download templates for child nodes which have empty file ID columns.
- 3. Populate the templates to prepare the file.tsv manifest (template populated with metadata) and the manifests for child node that include file ID columns
 - Ensure these child node manifests **do** include columns for file IDs.
 - Ensure correct file names are filled in columns for file IDs. The Uploader CLI Tool will generate final manifests for child nodes and replace the file names with corresponding file IDs.
- 4. Organize files.
 - Place the file.tsv and all child node metadata manifests into the same folder.
- 5. Run the Uploader CLI Tool.
 - Use the CLI tool to upload the data files.
 - During the upload the tool processes the manifests and generates file IDs.
- 6. Generate the final version of manifests.
 - Once the data upload is complete, the CLI tool automatically creates the "final" versions of the manifests (e.g., file-final.tsv).
 - These final manifests are saved in the same folder as the original manifests.
- 7. The final manifests are uploaded to the Submission Portal by the Uploader CLI Tool, so users do not need to make changes in the final manifests.
 - If you need to make changes to the final manifests after the Uploader CLI Tool has uploaded them to the portal, ensure you edit the final version (e.g., file-final.tsv saved locally on your computer) and then re-upload it through the user interface and not the Uploader CLI Tool.

• If you need to re-upload your data files for any reason, you can reuse the final manifest generated by the Uploader CLI Tool. This will retain the existing file IDs/keys instead of generating new ones.

3.5 Starting the Upload Process

Once the configuration file has been downloaded or edited, the upload script can be started. The only required parameter is --config, which should provide the full path and file name for the completed configuration file. The command should look something like the following, though the exact details may be customized depending on how the tool (and Python) were installed. Also, the following commands assume that your current directory is in the unzipped CLI directory.

\$ python3 scr/uploader.py --config path/to/metadata-upload.yml

When running the Windows version, the command should look like the following:

\$ uploader.exe --config path/to/metadata-upload.yml

When running Mac version, the command should be:

\$./uploader --config path/to/metadata-upload.yml

3.6 Using the CRDC Submission Portal's Graphical Interface to Upload Metadata Submission Templates

The Upload Metadata feature in the CRDC Submission Portal's graphical interface is intended for submitting completed metadata templates. Users can access tooltips and a link to the data submission user guide directly from the interface. The Upload Activities table tracks the uploading process of data and metadata files and provides details on any errors related to failed uploads. After the data upload is complete, the system automatically runs the basic validation process, and results are shown in the Validation Results table. You can delete the specific data or metadata files for a submission from the Data View table.

If you intend to upload the metadata in batches, you should keep your associated metadata separated by participants. For example, if a study has 100 participants, the submitted template for the first batch could either contain all 100 or a subset of that 100, with the remainder submitted in later batch uploads. If there is no overlap in participants between the different uploads, the system will not flag an error. However, mixing new data from previously uploaded participants with new participants will result in an error, as the system knows about the previously uploaded participants. To make corrections, select the file you want to delete in the Data View table and click the **Delete** button.

As depicted in Figure 23, to start the upload process, click the **Choose Files** button and then select the metadata submission manifests you want to submit. The total number of files that you have selected appears. If that number is correct, click the **Upload** button to start the upload. The Status column in the Upload Activities table displays *Uploading* until the upload and primary validation are completed. Once you have selected and uploaded the files, the CRDC Submission Portal automatically performs basic validations on the files and reports the results on the Validation Results table. Successful files show *Uploaded* in the Status column. If a file fails the primary validation checks, the data is not uploaded and the status displays as *Failed*.

1 UPLOAD Metadat	ta Files Choose Fil	No files selected				Upload
2 UPLOAD The CLI The CLI submiss click the	DATA FILES UPLOADER CL Tool is used to upload of Tools is a one-time dow ion. You can either edit button on the right to	IVERSION: <u>v3.2</u> data files to CRDC Subm wnload however the conf the example configuratio download a configuratior	ission Portal and requir figuration file needs to b n files found in the <u>CLI</u> n file customized for this	es a configuration file to work e customized for each Tool download [2], or you ca submission.		Download Configuration File
3 VALIDAT Validati Validati	E DATA on Type: () Val on Target: () Ne	lidate Metadata w Uploaded Data	Validate Data File	s O Both		Validate
		Upload Ac	tivities Validati	on Results Data View		
Batch ID	Batch Type	File Count	Status	Uploaded Date 🥠	Uploaded By	Upload Errors
5	Metadata	<u>16</u>	Uploaded	2/25/2025	Demo Account	
5	Metadata Metadata	<u>16</u> 2	Uploaded Failed	2/25/2025 2/25/2025	Demo Account Demo Account	2 Errors

Figure 23. Upload and Validate Data and Metadata

Clicking the number under the File Count column displays a list of all files uploaded in that batch. If there are errors in the files being uploaded, the Status column displays *Failed* and the Upload Errors column displays a link to the errors. Clicking that link opens a dialog box that explains what errors have been encountered. Correct all identified errors and reupload the file(s). If multiple files are uploaded in a batch, a failure in one of the files fails the entire batch. All files in a failed batch must be reuploaded. An example of a batch upload error message is presented in Figure 24.

Validation	Type:	Batch 63 Upload Errors	Validate
valuation	rarget.	33 ERRORS 1. "non_targeted_therapy.tsv.2": duplicated data detected:	
		"non_targeted_therapy_id": CTDC-NON-TARGETED-THERAPY-425649.	
Batch ID	Batch Ty	 "non_targeted_therapy.tsv:42": duplicated data detected: "non_targeted_therapy_id": CTDC-NON-TARGETED-THERAPY-425649. 	Upload Errors
67	Metadat	3. "non_targeted_therapy.tsv:65": duplicated data detected:	
66	Metadat	"non_targeted_therapy_id": CTDC-NON-TARGETED-THERAPY-425649.	
65	Metadat		
64	Metadat	Close	
63	Metadat		33 Errors

Figure 24. Batch Upload Errors

You can remove specific metadata files previously uploaded in the system in the Data View table, as depicted in Figure 25. For instance, to remove metadata for a participant, select the file and click the **Delete** icon to remove it from your submission. Note, this action will also delete the associated metadata from the child nodes.

de Type	e data file Status	All	\$	Submitted ID	
			Rows pe	er page: 20 👻 1-4 of 4	
•	File Name 🤟	Status	Orphaned	File Size	Uploaded Date/Time
	TCL01_Study_Protocol.docx	Error	Y	219492	2024-08-21T19:31:35.000Z
	TCL01_CellLine_CDS_Variants.csv	Passed	N	4173475	2024-08-01T19:33:26.000Z
	TCL01_CellLineSourcing_STRs_DrugSensitivity.xlsx	Error	N	41862	2024-08-21T19:31:34.000Z
		Error	Y		2024-08-01T19:32:57.000Z

Figure 25. Delete and Download Metadata File(s)

You can download the selected metadata files in the Data View table by selecting the file(s) and clicking the cloud-shaped **download** icon. See Figure 25.

To view the associated metadata from the child nodes, click the File Name in the first column of the Data View table. This action opens a window displaying metadata linked to the selected item. For example, clicking on the file name *TCLO1_Study_Protocol.docx* opens a window like the one in Figure 26. The window displays the related metadata from the child nodes such as Study, Participant, Diagnosis and Specimen. Note that the specific child nodes may vary depending on the data model used by each CRDC Data Commons.

Study (1)	Participant (1)	Diagnosis (1)	Specimen (1)	
		Rov	vs per page: 20 ▼ 1-1 of 1	< 1 >
study_short_name ↑	study_name	study_id	study_description	study_type
CTDC-STUDY-424369	scientia,_aere_perennius	damnant_quod_non_intellegunt	barba_crescit_caput_nescit	Proteomic Resear
		Rov	vs per page: 20 🔻 1-1 of 1	< 1 >

Figure 26. Visualize Associated Metadata from Child Nodes

4. Running Validations

Validations can be run at any point in the submission process, without restriction on when or how often. To run a validation, select options in the Validate Data panel and then click the **Validate** button. Refer to Figure 27.

3 VALIDATE DATA	VALIDATION COMPLETED			
Validation Type:	🔘 Validate Metadata	🔵 Validate Data Files	O Both	Validate
Validation Target:	🔘 New Uploaded Data	All Uploaded Data		



The first step is selecting which files to validate. The **Validate Metadata** option runs validations only on the submitted Metadata Manifest, not on any of the uploaded data files. The **Validate Data Files** option does the reverse and checks all the uploaded data files. The Both option validates both.

By default, only newly uploaded files are validated. This can be a significant time saver for large submissions as some validations take considerable time and the system keeps a record of previously submitted files that have already passed validation. However, if there is a need to check the entire submission, regardless of previous validations, the All Uploaded Data option validates everything that has been uploaded so far.

4.1 Reviewing Validation Results

After validations are run, the graphics on the page are updated to give a summary of the results as depicted in Figure 28.



Figure 28. Validation Summary

The left graph in Figure 28 displays a list of the nodes and indicates the status of the uploaded and validated data, with green representing data that passed the validations and red representing data that failed the validations. Hovering over each bar generates a more detailed summary for that node. For

instance, in Figure 28, the demographic node displays 150 uploaded files, 72 successfully passed validation, and 78 have errors. The blue color indicates data that has been uploaded but has not yet been validated. The graphs on the right are a node-by-node description of the results with the left and right arrows moving between the nodes that have been submitted to date.

4.1.1 Viewing and Filtering Validation Results

Submitters can view validation results in Aggregated (default) or Expanded format by clicking the toggle in the top left. Aggregated format groups validation results by Issue Type as shown in Figure 29.

Severity All			
Aggregated Expanded		Rows per page: 20 * 1-4 of 4 <	1 >
Issue Type	Severity	Count 🗸	Expand
Value not permitted	Error	1,659	Expand
Invalid Property	Warning	300	Expand
Invalid integer value	Error	300	Expand
Updating existing data	Warning	16	Expand

Figure 29. Aggregated Validation Results Table

To explore errors shown in the Expanded format (see Figure 30) in more detail, follow these steps:

- a. Click the **Expand** link under the Expand column to view detailed validation errors by Issue Type.
- b. Additional filters (Issue Type, Batch ID, Node Type, and Severity) can be helpful to narrow down the results. Selecting a specific **Node Type** refines the search to relevant validation errors. Choosing **All** from this menu displays a list of all files containing errors or warnings. The Node Type *file* pertains to metadata , while the *data file* refers to raw data, such as sequencing data.
- c. View details for a specific Issue under the Issue column by clicking the **See details** link to access more specific information regarding each validation error.
- d. Users can download all the validation errors in a table by clicking on the cloud icon.

		Upload Activities	Validation Results	Data View	
Issue Type	Value not permitted	Satch ID All	Node Type	All 🗘	Severity All
Aggregated	Expanded		Rows per page: 20 👻	1-20 of 1001 < 1 2	
Batch ID ↓	Node Type	Submitted Identifier	Severity	Validated Date	Issues
4	Demographic	CTDC-DEMOGRAPHI	Error	2/26/2025	Value not permi See details.
4	Demographic	CTDC-DEMOGRAPHI	Error	2/26/2025	Value not permi See details.
4	Demographic	CTDC-DEMOGRAPHI	Error	2/26/2025	Value not permi See details.
					the second second

Figure 30. Expanded Validation Results Table

This table shows all the errors that were found after the validations were run. The information in the columns can be interpreted as follows:

- **Batch ID** This correlates with the Batch ID shown on the Data Activity tab and indicates the specific upload with which the error is associated. This helps to identify which files may be involved.
- Node Type These correspond to the various metadata submission templates. In the example above, selecting ALL in the Node Type displays all nodes with errors. The figure demonstrates that the error is located within the Demographic metadata template.
- **Submitted Identifier** This is the identifier supplied by the user and is not a CRDC Submission Portal identifier. Again, this should specifically identify what object is causing the error.
- **Severity** Severity will either be *Error* (which must be corrected before the submission can be finalized) or *Warning* (which should be fixed but is not required to be fixed).
- Validated Date This is the date that the validation was run.
- **Issues** This gives a brief description of the error and a link to open a dialog box with more details about the error. Figure 31 presents an example of validation issue details.



Figure 31. Validation Error Details

4.2 Revising Released Metadata – Options to Keep or Overwrite Existing Content

After uploading the metadata to the CRDC Submission Portal and performing the validation process, if your metadata files include properties with values that already exist in the system from a previously released data submission, the system will generate a warning message. You can view these warnings in the Validation Results table by selecting the issue type - **update existing data**. Figure 32 shows how the warning messages appear for this Issue Type and Node Type *Genomic_info*, and Figure 33 illustrates the detailed information displayed when you click the "see details" link in the **Issues** column.

		Upload Activities	Validation Results	Data View		
Issue Type	Ipdating existing d 💲	Batch ID All	Node Type All		Severity All	\$
Aggregated	Expanded			Rows per page: 20 👻	- 1-17 of 17 < 1 >	۵
Batch ID ↓	Node Type	Submitted Identifier	Severity	Validated Date	Issues	
3	Genomic_info	dg.4DFC/840c7e9	Warning	6/5/2025	Updating existi <u>See det</u> a	<u>ails.</u>
3	Genomic_info	dg.4DFC/840c7e9	Warning	6/5/2025	Updating existi <u>See deta</u>	<u>ails.</u>
3	Genomic_info	dg.4DFC/840c7e9	Warning	6/5/2025	Updating existi See deta	aits.
3	Genomic_info	dg.4DFC/840c7e9	Warning	6/5/2025	Updating existi <u>See deta</u>	ails.

Figure 32. Visualize Validation Issue Type – Update Existing Data

/alidatio	on Target:	Valio For Genor	dation Issues nic_info Node IDdg.4DFC/840c7e9c-8900-4d1c	- 1-96667c	e99c1e_1	×	Validate
ype A	Expand	1. (War "genom exists a	rning) [cds_mock_genomic_infoSofia-final.tsv: lin nic_info_id: dg.4DFC/840c7e9c-8900-4d1d-966i nd will be updated.	e 2] "gen 67ce99c1	omic_info": e_1" : already		y All 7 < 1 > ≪
ID ↓	25						
	Genomic_info		genomic_info_id	bases	number_of_reads		ating existi <u>See details.</u>
	Genomic_info	Existing	dg.4DFC/840c7e9c-8900-4d1d-96667ce99c1e_1	648	116		ating existi See details.
		New	dg.4DFC/840c7e9c-8900-4d1d-96667ce99c1e_1	648	116		
	Genomic_info Genomic_info Genomic_info Genomic_info Genomic_info						ating existi See details. ating existi See details.
	Image			ating existi <u>See details.</u>			
	Image		Close				ating existi <u>See details.</u>
	Image						ating existi <u>See details.</u>

Figure 33. Illustration of the Details of Warning Message – Update Existing Data

It is important to follow the guidelines below to ensure the existing information, from previous submissions, is not unintentionally removed or overwritten. Note that validations will not be complete until the necessary actions are taken.

What Happens Automatically

- If a property or column in your metadata file(tsv) is left blank or not included, the system will keep the existing value for that field. See Table 1.
 - o This applies to both required and optional fields
 - For example, if the existing data for a sample has "sample_type = Tumor", and you leave sample_type blank or omit this property in your metadata file, the system will keep "Tumor."

What Submitter Should Not Do

- Do not leave the fields blank if your intention is to delete or clear a value. Blanks will not remove data.
- Do not try to remove a required field using special commands (see below); required fields cannot be deleted

How to Delete a Value (if needed)

- If you need to **delete** a value from the previously released submission (e.g., an incorrect entry), use command word <**delete>** as the value for that specific property.
- It is not **case-sensitive**: <delete>, <Delete>, and <DELETE> all work the same. This word must appear exactly as shown, with angle brackets (< >).
- Use this only for **optional fields**. If you try to delete a **required field**, your data submission cannot be submitted

Property Type	Property Name	Existing Value in System	Value in Your Metadata file	What Happens?
Required	sample_type	Tumor	(blank)	Keeps "Tumor"
Required	sample_type	Tumor	<delete></delete>	Error – can't delete required field
Optional	sample_type_category	DNA	<delete></delete>	Deletes DNA
Optional	sample_type_category	DNA	(blank)	Keeps DNA
Optional	sample_type_category	(empty)	(not included)	Still empty, nothing changes

Table 1. Summary of how existing metadata is updated during submission

4.3 Correcting Errors

Errors should be corrected by addressing the issues in local files, re-uploading the corrected file, and running the validation again. This process should be repeated until all errors have been addressed, and the validation returns no errors.

Anything marked as an *Error* in the Severity table must be fixed before the dataset can be formally submitted. Anything marked as a *Warning* will not block the final submission; however, users are <u>strongly</u> <u>encouraged</u> to fix warnings as well.

4.4 Remove Specific Files

After validating the uploaded data and metadata files, users can view the high-level details of these files in the **Data View** table. As shown in Figure 34, select either **data file** or **file** from the Node Type dropdown menu. To remove a specific data file or file ID along with its associated metadata, select the file and then click the delete icon. Users can also download the contents displayed in the Data View table as a TSV file by choosing the Node Type and clicking the download icon.

Node Type data file	\$ Sta	tus All	\$	Submitted I	D
data file file			Row	s per page: 20 👻 1-4 of	f4 < 1 > 🙆 🗟
🗕 🛛 File Name 🗸		Status	Orphaned	File Size	Uploaded Date/Time
TCL01_Study_Prote	ocoldocx	Error	Y	219492	2024-08-21T19:31:35.000Z
TCL01_CellLine_CE	OS_Variants.csv	Passed	Ν	4173475	2024-08-01T19:33:26.000Z
TCL01_CellLineSou	rcing_STRs_DrugSensitivity.xlsx	Error	Ν	41862	2024-08-21T19:31:34.000Z
		Error	Y		2024-08-01T19:32:57.000Z

Figure 34. Removing Specific Files

The contents of the Data View table can be interpreted as

- File Name This column displays the name of the uploaded data file.
- **Status** This indicates the validation status, which can be *Error*, *Warning*, or *Passed*. Details of any errors are available under the Validation Results. Warnings do not halt the submission process.
- **Orphaned** This column shows whether the data file has associated metadata uploaded in the system. *Y* means the file is orphaned and lacks associated metadata, while *N* means the file has associated metadata uploaded in the system.
- File Size This column lists the file size in bytes.
- Uploaded Data/Time This column records the date and time when the file was uploaded.

5. Submitting Your Final Dataset

When a dataset has passed all validations with no outstanding errors, the Submit button at the bottom of the page is activated. Clicking the Submit button locks the submission and passes control to the Data Concierge for a final check. No further changes will be allowed. Should the Submit button be clicked in error, contact the assigned Data Concierge and they can reject the submission and return the control to you.

6. What to Expect After Submission

Once the final dataset has been submitted, the CRDC Submission Team will perform some final checks to make sure everything is as required by the destination Data Commons (for example, GC, ICDC, and CTDC). If those checks pass, the submission will be released to the appropriate CRDC Data Commons, and you will receive notification that the Data Commons is now responsible for the next steps. The respective Data Commons will be responsible for indexing and releasing the files for secondary sharing and will be made available and accessible on their portal.

If the final checks reveal some unexpected issues or the Data Commons team has any questions or concerns, the Data Concierge for your submission will reach out with additional questions and may reopen the submission to allow additional corrections.