UCD IMPROVE Technical Information #351A

Data Ingest

Interagency Monitoring of Protected Visual Environments Air Quality Research Center University of California, Davis

> October 3, 2022 Version 1.0

Prepared By:	Docusigned by: Indu Huckemeppilly Sivakumar 19F6B63B1B17443	Date:	10/7/2022
Reviewed By:	DocuSigned by: Dominique Young BB55DBA34BAB407	Date:	10/7/2022
Approved By:	DocuSigned by: Marcus Langston 0A10CECE79B0452	Date:	10/7/2022



DOCUMENT HISTORY

Date Modified	Initials	Section/s Modified	Brief Description of Modifications
03/14/22	SRS	All	Previously anthologized version separated into individual TIs.

TABLE OF CONTENTS

1.	Purpose and Applicability				
2.	Sum	mary of the Method			
3.	Definitions				
4.	Health and Safety Warnings				
5.	Caut	ions			
6.	Inter	ferences			
7.	Perso	onnel Qualifications			
8.	Equi	pment and Supplies			
9.	Proce	edural Steps5			
9.	.1 (Carbon Results			
9.	.2 I	Ion Results7			
9.	.3 I	Element and Optical Absorption Results			
9.	.4 F	Re-ingesting			
9.	.5 I	Issue Tracking			
10.	Qu	ality Assurance and Quality Control			
10	0.1	Code Development			
10	0.2	Bug Reporting			
10	0.3	Data Validation10			
11.	Re	ferences			

LIST OF FIGURES

Figure 1. Carbon analysis results upload page	6
Figure 2. Ions analysis results upload page	8

LIST OF TABLES

Table 1. Automated validity checks performed during carbon data upload	7
Table 2. Automated validity checks performed during the ions data upload	8

1. PURPOSE AND APPLICABILITY

The purpose of this technical information (TI) is to provide information on handling electronic laboratory records from samples collected in the Interagency Monitoring of Protected Visual Environments (IMPROVE) network. This document is intended to guide users on the receiving and validating of IMPROVE laboratory records and ingestion to the University of California, Davis (UCD), IMPROVE database. These include ion analysis results from RTI International (RTI), carbon analysis results from Desert Research Institute (DRI), and gravimetric mass, elemental, and filter absorption analysis results from UCD.

2. SUMMARY OF THE METHOD

Ion analysis results from RTI and carbon analysis results from DRI are received in data files, typically delivered as .csv files and XML files, respectively. The files are ingested into the UCD IMPROVE database using the UCD IMPROVE Data Management website. Gravimetric mass, elemental, and filter absorption analysis results from UCD are automatically ingested.

3. **DEFINITIONS**

- AQRC: Air Quality Research Center.
- CSV: a comma-separated value file that is the common format for delivery files.
- **DRI:** Desert Research Institute.
- Energy Dispersive X-Ray Fluorescence (EDXRF): An analytical technique used to determine the concentration of elements.
- Hybrid Integrating Plate/Sphere (HIPS): An analytical technique for optical absorption.
- Ion Chromatography (IC): An analytical technique used to determine the concentration of ions.
- Interagency Monitoring of Protected Visual Environments (IMPROVE): Federal PM_{2.5} and PM₁₀ sampling network directed by the National Park Service, with sites located principally in remote rural areas.
- **IMPROVE database:** A SQL Server database that is the central warehouse of IMPROVE preliminary and final data at UCD.
- **PM:** Particulate Matter. $PM_{2.5}$ is particulate matter with diameters 2.5 micrometers (μ m) and smaller. PM_{10} is particulate matter with diameters 10 μ m or smaller.
- **RTI:** Research Triangle Institute, International.
- **SQL:** database management system used by AQRC.
- Thermal Optical Analysis (TOA): An analytical technique used to determine the concentration of carbon. Also referred to as TOR (Thermal Optical Reflectance) and TOT (Thermal Optical Transmittance).
- UCD: University of CA—Davis.

• Extensible Markup Language (XML): a markup language defining a set of rules for encoding documents in a particular format; used for IMPROVE carbon files.

4. HEALTH AND SAFETY WARNINGS

Not applicable.

5. CAUTIONS

Not applicable.

6. INTERFERENCES

Not applicable.

7. PERSONNEL QUALIFICATIONS

The UCD Air Quality Research Center (AQRC) Data & Reporting Group staff assigned to tasks described in this document have advanced training in database programming and database management.

8. EQUIPMENT AND SUPPLIES

The hardware and software used for IMPROVE data ingest are described in the associated UCD IMPROVE SOP #351: Data Processing & Validation.

9. PROCEDURAL STEPS

Prior to data processing and validation, data are ingested for each of the analysis pathways: (1) carbon results from DRI, (2) ions results from RTI, and (3) elemental and optical absorption results from UCD.

9.1 Carbon Results

Carbon analysis results are sent from DRI to UCD via email in .xml format, including three files:

- 1. CarbonData.xml
- 2. CarbonInformation.xml
- 3. CarbonLaser.xml

All three files are included in a zip folder which will be saved to U:\IMPROVE\RawDataReceived\Carbon DRI\To be Imported. The files will be extracted by right clicking on zipped folder, selecting 7-Zip and then 'Extract Here'. A new folder with the same name as the zipped folder will be created in the same location on the U Drive and contain all three files as named above.

All three files are ingested using the UCD IMPROVE Management website. Figure 1 shows a screenshot of the carbon data upload page, which is accessed via the Analysis Data Section as described in section 8, selecting the **Carbons** tab, and clicking the **Ingest Data** button. To ingest the files from the data upload page, select the relevant files, create a name for the import batch under *Batch Label*, and click **Submit**. The suggested batch label is the filename from DRI (First sample date - Last Sample date).

CarbonInformation, *CarbonLaser*, and *CarbonData* are ingested simultaneously, and an automated validity check is performed (Table 1). Results from the validity check will indicate upload failures. The Quality Assurance Officer will review the upload results and notify the Lead Quality Assurance Officer if there are upload failures from validation errors. After ingest, the source files are stored on the file server at

U:\IMPROVE\RawDataReceived\Carbon DRI\Imported, within a folder which is named the same as the *Batch Label* created for ingest. After successfully ingesting the results, the Quality Assurance Officer will save a copy of the Carbon Data Ingestion summary page by printing the page to a PDF and saving to

U:\IMPROVE\RawDataReceived\Carbon DRI\IngestRecord\. Further details of the ingest file are recorded in a log file located at U:\IMPROVE\RawDataReceived\Carbon DRI\Carbon_Ingest_log.xlsx.

Figure 1. Carbon analysis results upload page.

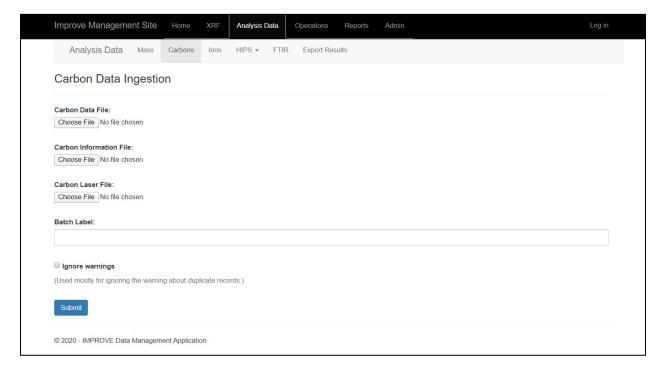


Table 1. Automated validity checks performed during carbon data upload.

Check	Action
Basic schema validation on xml files	Error
No filter found for record	Warning
Filter.Module doesn't match record Site field	Warning
Record is marked as re-analysis	Warning
Carbon Laser file has records missing wavelength	Warning
Found more parameter records than expected for an analysis	Warning
Parameter missing for an analysis	Warning
Comment from DRI on analysis	Note
Parameter/record already recorded in database	Warning
Incomplete analysis record (missing entries in either Carbon/Carbon	Warning
Laser/Carbon Info file)	

9.2 Ion Results

Ions analysis results are sent as one file from RTI to UCD via email in .csv format. The naming convention of the ions data includes the year followed by the ions data set number (e.g. '2020 1 2 3 data export to UCD'). The file is saved to the file server at U:\IMPROVE\RawDataReceived\Ions RTI\To Be Ingested.

The ion analysis records are ingested using the UCD IMPROVE Management website. Figure 2 shows a screenshot of the ions data upload page, accessed via the Analysis Data Section as described in section 8, selecting the Ions tab, and click the **Upload Data** button. To ingest a file from the data upload page, select the relevant file and click **Continue**. An automated validity check is performed, and the validity check results will indicate if there are upload failures (Table 2). The Quality Assurance Officer will review the upload results and notify the Lead Quality Assurance Officer if there are upload failures from validation errors. An ingest page after passing validation is shown in Figure 3.If the file is ready for ingestion, click **Submit**. After ingest, the source files are stored on the file server at U:\IMPROVE\RawDataReceived\Ions RTI\Ingested. After successfully ingesting the results, the Quality Assurance Officer will save a copy of the Ions Data Ingestion Status summary page by printing the page to a PDF and saving to U:\IMPROVE\RawDataReceived\Ions RTI\Ingest_record\. Further details of the ingest file are recorded in a log file located at U:\IMPROVE\RawDataReceived\Ions RTI\Ions_DataIngest_Log.xlsx.

Data Ingest UCD TI #351A, Version 1.0 October 3, 2022 Page 8 of 10

Figure 2. Ions analysis results upload page.

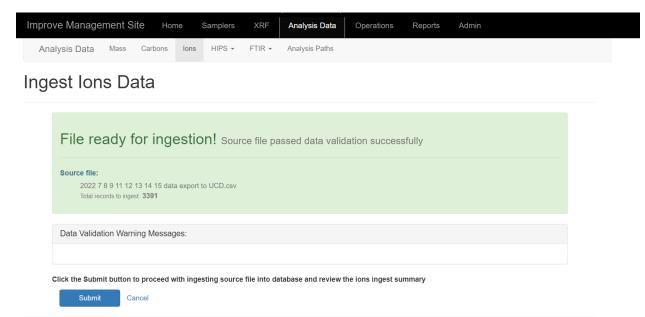
Improve Management Site	Home	XRF	Analysis Data	Operations Reports	Admin			Log off
Analysis Data Mass	Carbons	lons	HIPS + FTIR	Export Results				
Ions Analysis Upload Data								
Select lons analysis source file (. Choose File 2019 45 202ucd_N		required):					
Continue >>								

© 2020 - IMPROVE Data Management Application

Table 2. Automated validity checks performed during the ions data upload.

Check	Action
Basic schema validation on csv files	Error message
No filter is found for record	Error message
Data already exists for filter record	Warning message
Parameter missing for a filter e.g., parameter has null value or	The ingest process will
parameter column isn't present at all	abort with an error message
Parameter already recorded in database	Warning

Figure 3. Ions ingest page after passing validation



9.3 Element and Optical Absorption Results

Elemental analysis is performed at UCD. The PANalytical XRF software generates results files and automatically transmits them to a directory on an AQMT file server. A service on the server (internally named *XRF Data Transfer*) monitors the transmission directory, checking every five minutes for new files. The XRF results files are standard text files with the extension .*qan*. The file name includes XRF analysis dates and times in the format *YYYYMMDDHHMMSS.qan*. The results files and contents are automatically parsed and ingested into tables in the UCD IMPROVE database.

Optical absorption analysis is performed at UCD. The HIPS instrument generates results which are then verified by the operator to be complete and then written to the database. The data are then available on the UCD IMPROVE database.

9.4 Re-ingesting

If errors are identified in the source files from DRI or RTI that cause the import to fail, or if results are updated as part of the validation and reanalysis process, new files must be requested and provided for ingestion. Upload the new files using the process described in sections 9.1 and 9.2.

For carbon, whether the files contain new batches of data or reanalysis results, take care to ingest with the *ignore warnings* box unchecked. Scrutinize the messages and warnings to check for errors and take note of further actions that may be required after the data is ingested (e.g., changing analysis QC codes). The import process indicates if there are matching existing records, if existing records are not updated, or if only new records are added (including cases with different analysis results from the sample filter). Once the messages have been reviewed and addressed, re-run the ingest process with the *ignore warnings* box checked. For carbon, if the reanalysis results are used, the analysis QC code can be updated using the tool available at

https://improve.aqrc.ucdavis.edu/AnalysisData/Carbons/CarbonsQcReview.

For ions, the data are ingested without any changes to the original process; the QC code is updated using the tool available at https://improve.aqrc.ucdavis.edu/AnalysisData/Ions/IonsQcReview.

9.5 Issue Tracking

Software bugs and data management issues are tracked through JIRA tracking software. All users who have access to the internal UCD JIRA website can submit, track, and comment on issues. Users requesting new tools, modifications to existing tools, or to report bugs specific to the IMPROVE data should add JIRA tickets to the IMPROVE Data Management Software project at

https://improve.atlassian.net/jira/software/c/projects/IMPSW/issues/

10. QUALITY ASSURANCE AND QUALITY CONTROL

10.1 Code Development

Software for data management, processing, and validation is developed in-house by professional software engineers. Source code is managed through a code repository. Development of code changes and new applications is conducted on a development environment that parallels the production environment. Prior to deployment in production, all code changes undergo testing within a separate test environment. The testing, which is conducted by developers, managers, and users, is targeted both at the identification of software bugs and the confirmation of valid data equivalent to the production system.

10.2 Bug Reporting

Software bugs and data management issues are tracked through JIRA tracking software. All UCD users have access to an internal JIRA website and can submit, track, and comment on bug reports.

10.3 Data Validation

Data integrity is enforced within the UCD IMPROVE database via unique primary keys and non-nullable records. Data completeness and data quality are thoroughly checked through the data validation process, as described elsewhere in this SOP.

11. REFERENCES

Not applicable.