

Overview Guide for the GCE Data Toolbox for MATLAB

Richard Cary

Wade M. Sheldon, Jr.

John F. Chamblee

Introduction

This guide will cover various methods of importing data files into the GCE Data Toolbox for MATLAB software, working with the data, creating a metadata template for the file, joining multiple data sets, QA/QC processing and exporting documented data products. At this point, this guide is not designed to be a stand-alone document, but is instead written to accompany oral presentations and working demonstrations that elaborate upon the described concepts and help you work through the presented examples. You have been given a series of data files to use in completing these exercises. Specific file names are mentioned throughout the text, as they are referenced. We recommend keeping all the files in a single folder, labeled “workshop_products” on your desktop. We will refer to the workshop_products folder explicitly or implicitly assume its use throughout the guide.

The guide uses several text formatting conventions to alert you to different kinds of information.

Section Headers are bold and underlined. Subsection Headers are underlined, but not bold. Code snippets and MATLAB “**commands**” (i.e. for typing in the MATLAB command window) are in enclosed in quotes and are bolded. GCE Data Toolbox menu operations (e.g. ***File>Load Structure>Load Structure from File***) are enclosed in quotes, bold and italicized. Each level in a menu operation is separated by a right angular bracket (“>”). Figure captions are italicized.

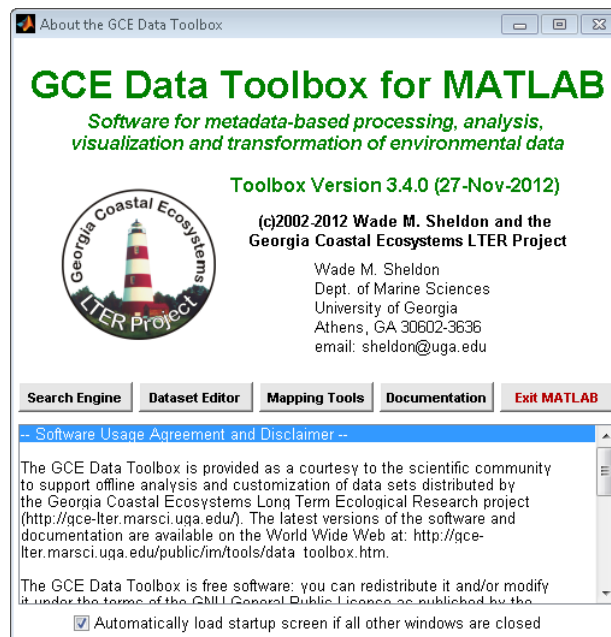
Starting the Toolbox and Importing Data

This section is designed to introduce you to the GCE Toolbox interface and environment. The goal is to become familiar with the overall interface by importing a raw, undocumented text file and then manipulating that file to explore basic toolbox functionality.

Getting Started

1. In order to use the GCE Data Toolbox for MATLAB, you must have a licensed and activated copy of MATLAB installed. You can find instructions for MATLAB installation on the Mathworks website at <http://www.mathworks.com/support/install.html>. We recommend installing MATLAB in a directory named “MATLAB” on the root of a local hard drive rather than in C:\Program Files\MATLAB (e.g. C:\MATLAB\R2012B) as spaces in pathnames can cause problems with scripts and third party applications that call MATLAB (e.g. Kepler).
2. You must also have a copy of the Toolbox code library on the local hard drive or a network-accessible directory. Download a complete distribution package as a Zip file from the GCE Toolbox Trac web site at https://gce-svn.marsci.uga.edu/trac/GCE_Toolbox/wiki/Downloads. You can also check out the code from the GCE Subversion repository at https://gce-svn.marsci.uga.edu/svn/GCE_Toolbox/trunk using an SVN client (login required – contact Wade Sheldon for details).
3. To install the toolbox, simply extract the downloaded files to any directory accessible to MATLAB. For beginning users, we recommend installing the files in a local folder called “GCE_Toolbox” within the MATLAB installation folder (e.g. C:\MATLAB\GCE_Toolbox). Note that write access to the toolbox root and \userdata directories is required, so avoid installing the toolbox in a write-protected server directory.

4. To start the Toolbox, you need to start up MATLAB, navigate to the folder where the toolbox was installed, and run the startup.m script. You can use the MATLAB path browser tool to change the working directory then double click on the startup.m file in the "Current Folder" file list. You can also use the "**cd**" command to change the working directory (as in Unix/DOS) and type "**startup**" in the command window to run the script. You can also create a MATLAB shortcut to change to the directory and run the startup script then just click on the shortcut to start the toolbox (e.g. "cd \path\to\toolbox; startup").
5. A GUI startup dialog is displayed by default when first starting the Toolbox and when all GUI dialogs are closed. Buttons on this dialog are used to launch primary Toolbox programs, display the documentation viewer, and exit the MATLAB environment.

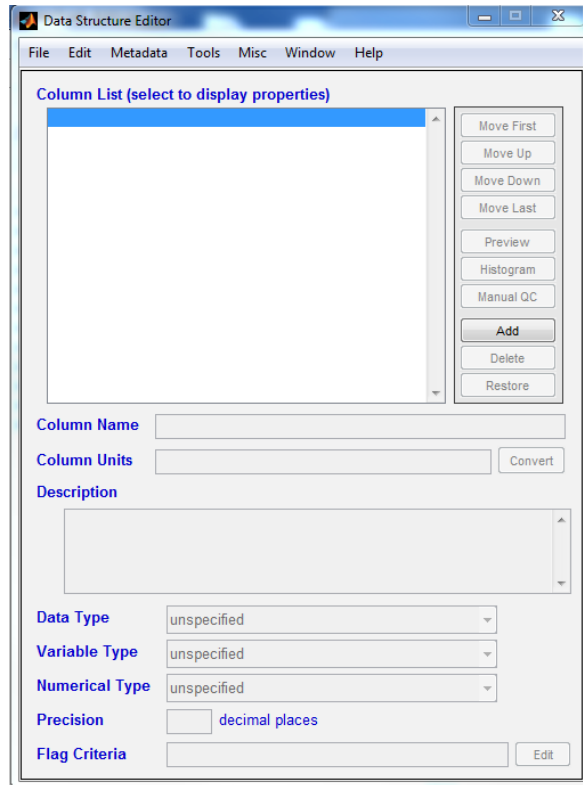


The GCE Data Toolbox startup screen

Loading a Raw, Undocumented ASCII File-

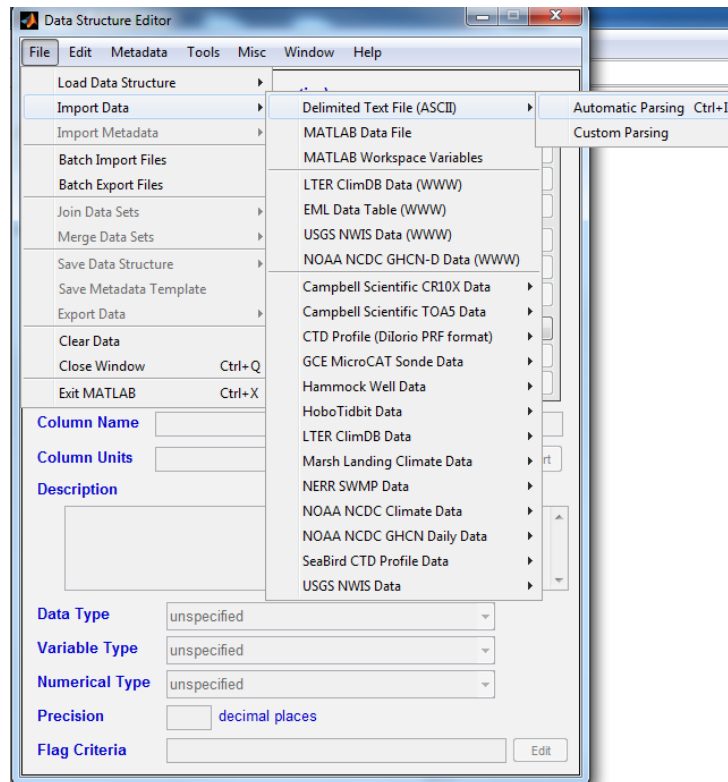
This section covers how to load a simple delimited ASCII file with 1-row header into the GCE Data Toolbox, and provides a brief overview of the features of the Data Structure Editor application.

1. From the initial GCE Data Toolbox startup screen, click the "Dataset Editor" button to bring up the Editor window. The window will initially be empty with most menus and buttons disabled until data are imported or loaded.



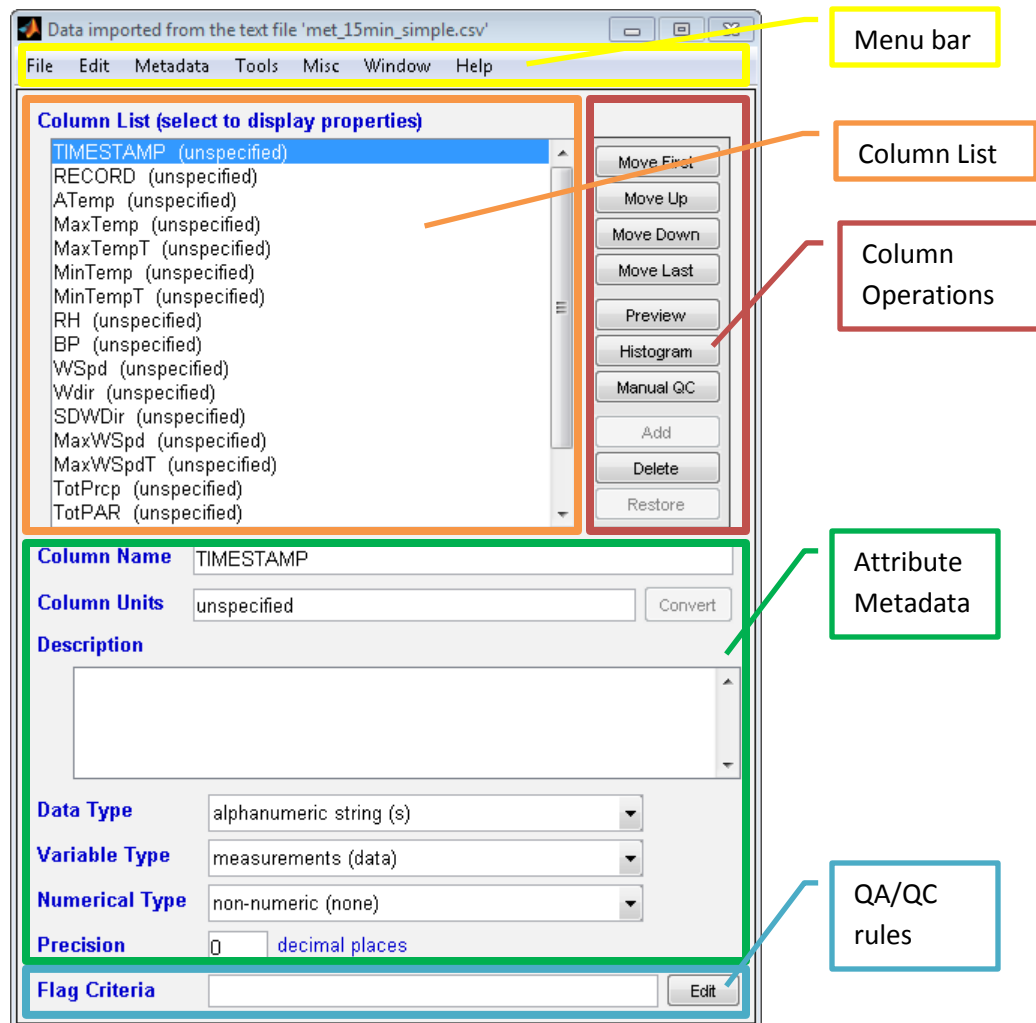
A blank Dataset Editor screen

2. Using the menu bar at the top, select **"File > Import Data > Delimited Text File (ASCII) > Automatic Parsing"**. When prompted, navigate to the **met_15min_simple.csv** sample file and click the Open button on the file loading dialog. The file is then parsed and the data are loaded into the Data Structure Editor. This is an example raw data file in comma-separated value (CSV) text format that only contains a 1-line header with column names followed by data rows – open the file in a text editor if desired to examine the native format.



Loading a raw ASCII file using the Data Structure Editor

3. Now that the .csv file has been loaded, we can take a look at the Data Structure Editor window itself. The screenshot below describes the general layout of the window.



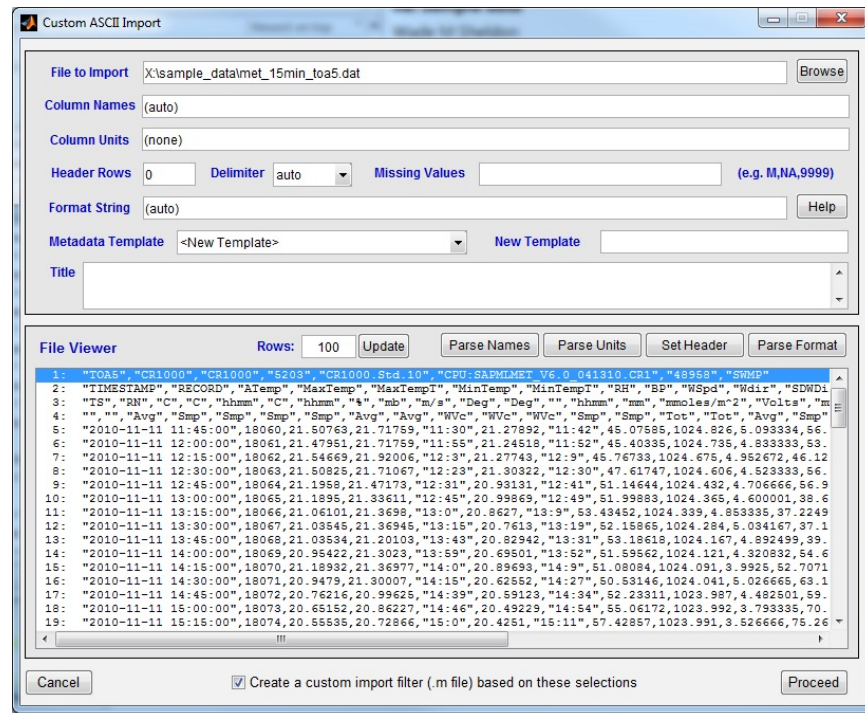
Items in the Data Structure Editor: 1) Menu Bar, 2) Data Column List, 3) Column order and basic data operation buttons, 4) Attribute Metadata, 5) Column QA/QC rules

- When met_15min_simple.csv is loaded into the Data Structure Editor, the GCE Toolbox assigns basic attribute metadata for the file by inspecting column formats and numeric scales. Click on each name in the column list to view the attribute metadata for the column. Note that because this is a simple ASCII file with only column names in the header the dataset will not initially contain any unit information, descriptions or QA/QC flag criteria.

Importing a Custom Delimited ASCII File

Now we will cover how to import a more complex delimited ASCII file using the "Custom Parsing" option. ASCII files that have multiline headers or missing value codes other than NaN, Inf or -Inf (the IEEE standards recognized by MATLAB) are incompatible with "Automatic Parsing", and will result in an error. More information is required in order to parse the data.

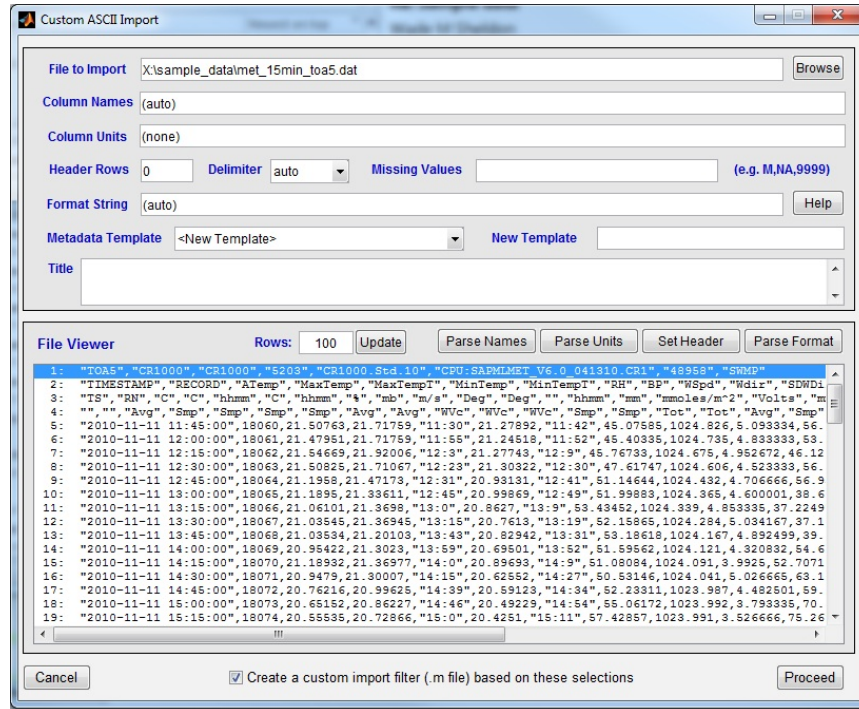
1. From the Data Structure Editor Window, select **"File > Import Data > Delimited Text File (ASCII) > Custom Parsing"**. This will bring up a new Custom ASCII Import window.
2. In the "File to Import" field, click on the **Browse** button to open a file loading dialog and navigate to the file **met_15min_toa5.dat**.
3. The first 100 rows of the file are displayed in the File Viewer listbox to aid in filling in the information required to parse the file. You can scroll through the data rows and load more rows if desired by editing the **Rows** field and hitting **Update**.



Custom ASCII Import window with data loaded

4. Now we need to identify the column names and the column units. In this data set, the column names are located in row 2. Select row 2 in the File Viewer window and then click the **"Parse Names"** button. The column names from the second row will be added to the Column Names field above. Repeat this for the column units (found in row 3) using the **"Parse Units"** button.
5. To complete the Format String field, select a representative data row that does not contain any missing values in the File Viewer listbox, and then click the **"Parse Format"** button. This will determine the data type of each field and fill in the appropriate field token automatically (click on **"Help"** next to the Format String field for more information)
6. Enter the missing value code for this data set in the Missing Values field. The code is **"NAN"** including the quotation marks. If the missing value code is not specified, MATLAB will display an error message and list the line number and text where the error occurred.
7. Fill out the Header Rows, either by entering the number of rows that comprise the data set header, or selecting the last header row in the file viewer and the hitting the **"Set Header"** button.

8. If you already had a metadata template created for this data set, you could select it from the Metadata Template drop-down menu and it would be applied to this dataset following completion of the import. Since we don't have one for this data set yet, we will leave this blank. Additionally, you can create a new empty metadata template based on the information entered in the import dialog. Enter **"Test Template"** in to the **New Template** field to create a new template for future use.
9. Select **comma** as the Delimiter type, and enter a Title. It is possible to save the import filter values for use with similarly formatted data, so leave the **"Create a custom import file"** box at the bottom checked and click the **"Proceed"** button.



Custom ASCII importer that has been filled out.

10. A file save dialog box will be displayed for saving the new import filter. Give the file a name, e.g. **"test_filter.m"** and save it in the \userdata Toolbox directory.
11. The data will now be imported into the Data Structure Editor. If you look at the attribute metadata, you can see that additional metadata such as the attribute units have been added to each attribute.
12. A new custom import filter has also been added to the Toolbox. To use the new filter to re-import the data file, use **"File > Import Data > User-Defined Text File Format"** and select the entry for your new filter.
13. Note: to edit the name of the filter or delete it from the GCE Toolbox menus, use **"Misc > Add/Edit Import Filters"**

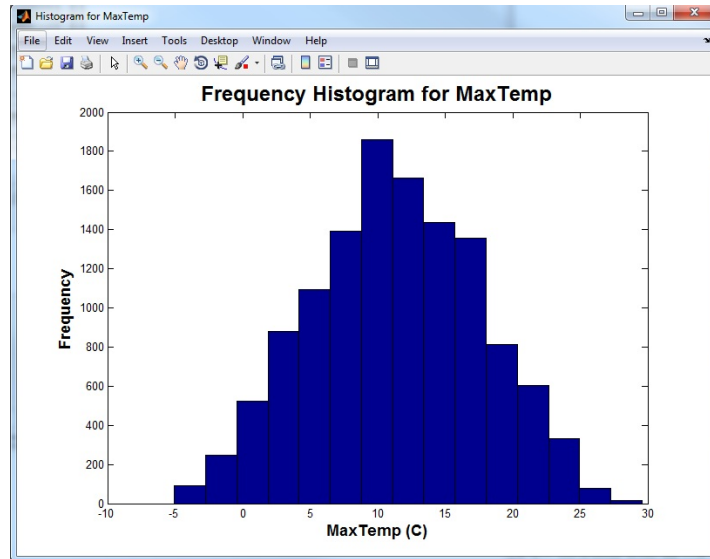
Exploring the Data Structure Editor

1. All of the parsed data columns are displayed in the Column List in order of occurrence in the file (i.e. first column listed is the left-most column and so on)
2. You can use the Data Structure Editor to manage the the order of columns and delete unneeded columns using the buttons on the right side of the Data Structure Editor. For example, select column MinTemp in the list, and press "**Move Up**" three times so that it is in the third position. You can delete one or more columns selecting them and hitting the "**Delete**" button. Pressing "**Restore**" will restore all deleted columns to the list.
3. There are also a number of useful tools on the right side of Data Structure Editor window that can be used to perform basic data inspection. The first of these is the "**Preview**" button. Hitting this button will bring up a new window containing just the data for the selected column in the column list to preview the formatting.



Preview window with maximum temperature data.

4. The next tool located here is "**Histogram**". Hitting this button will bring up a frequency histogram plot of the data in the selected column in the column list. This can be used to take a quick look at the range of column data values, or can be used for quick visual inspection for outliers or extreme values.



Example of Histogram function output.

- The last tool in this section is the “**Manual QC**” function. Clicking this button will open a Data Editor screen containing only the data for the selected column and any QA/QC flags assigned. Since we have not developed any flagging criteria for this data set, the flag fields will all be blank. You can also perform manual edits to the data in these columns from this screen.

The 'Assign QA/QC Flags' window displays a table with 25 rows of data. Each row has a checkbox, a temperature value, and a flag field. All flag fields are currently blank.

	MaxTemp (C)	Flag_MaxTemp (flags)
<input type="checkbox"/> 1	21.7176	
<input type="checkbox"/> 2	21.7176	
<input type="checkbox"/> 3	21.9201	
<input type="checkbox"/> 4	21.7107	
<input type="checkbox"/> 5	21.4717	
<input type="checkbox"/> 6	21.3361	
<input type="checkbox"/> 7	21.3698	
<input type="checkbox"/> 8	21.3695	
<input type="checkbox"/> 9	21.2010	
<input type="checkbox"/> 10	21.3023	
<input type="checkbox"/> 11	21.3698	
<input type="checkbox"/> 12	21.3001	
<input type="checkbox"/> 13	20.9963	
<input type="checkbox"/> 14	20.8623	
<input type="checkbox"/> 15	20.7287	
<input type="checkbox"/> 16	20.5938	
<input type="checkbox"/> 17	20.0540	
<input type="checkbox"/> 18	19.8211	
<input type="checkbox"/> 19	19.5520	
<input type="checkbox"/> 20	19.0510	
<input type="checkbox"/> 21	18.8156	
<input type="checkbox"/> 22	18.5804	
<input type="checkbox"/> 23	18.1135	
<input type="checkbox"/> 24	17.4155	
<input type="checkbox"/> 25	16.6412	

Manual QC tool with maximum temperature data loaded.

Viewing and Editing Data Values

- Once a data set is loaded in to the Data Structure Editor you can view the data values by going to **"Edit > View/Edit Data"**. This will bring up the Data Editor window, which displays data values in a spreadsheet-like grid with horizontal and vertical scrollbars (as necessary). Column names and units are displayed at the top of each column. Headings for text columns are green and values a left-aligned, and numeric columns are blue and right-aligned as a visual aid.

	All	None	TIMESTAMP (TS)	RECORD (RN)	ATemp (C)	MaxTemp (C)	MaxTempT (hhmm)	MinTemp (C)	MinTempT (hhmm)	RH (%)	BP (mb)
1			2010-11-11 11:45:0	18060	21.5076	21.7176	11:30	21.2789	11:42	45.0759	1024.83
2			2010-11-11 12:00:0	18061	21.4795	21.7176	11:55	21.2452	11:52	45.4034	1024.73
3			2010-11-11 12:15:0	18062	21.5487	21.9201	12:3	21.2774	12:9	45.7673	1024.68
4			2010-11-11 12:30:0	18063	21.5083	21.7107	12:23	21.3032	12:30	47.6175	1024.61
5			2010-11-11 12:45:0	18064	21.1958	21.4717	12:31	20.9313	12:41	51.1464	1024.43
6			2010-11-11 13:00:0	18065	21.1895	21.3361	12:45	20.9987	12:49	51.9988	1024.37
7			2010-11-11 13:15:0	18066	21.0610	21.3698	13:0	20.8627	13:9	53.4345	1024.34
8			2010-11-11 13:30:0	18067	21.0355	21.3695	13:15	20.7613	13:19	52.1587	1024.28
9			2010-11-11 13:45:0	18068	21.0353	21.2010	13:43	20.8294	13:31	53.1862	1024.17
10			2010-11-11 14:00:0	18069	20.9542	21.3023	13:59	20.6950	13:52	51.5955	1024.12
11			2010-11-11 14:15:0	18070	21.1893	21.3698	14:0	20.8969	14:9	51.0806	1024.09
12			2010-11-11 14:30:0	18071	20.9479	21.3001	14:15	20.6255	14:27	50.5315	1024.04
13			2010-11-11 14:45:0	18072	20.7822	20.9963	14:39	20.5912	14:34	52.2331	1023.99
14			2010-11-11 15:00:0	18073	20.8515	20.8623	14:46	20.4923	14:54	55.0617	1023.99
15			2010-11-11 15:15:0	18074	20.5554	20.7287	15:0	20.4251	15:11	57.4286	1023.99
16			2010-11-11 15:30:0	18075	20.2289	20.5938	15:15	19.9865	15:29	58.5433	1024.09
17			2010-11-11 15:45:0	18076	19.8641	20.0540	15:30	19.6491	15:45	60.6480	1024.08
18			2010-11-11 16:00:0	18077	19.6217	19.8211	15:56	19.4130	15:47	60.8144	1024.02
19			2010-11-11 16:15:0	18078	19.1391	19.5520	16:0	18.8473	16:11	62.1177	1024.15
20			2010-11-11 16:30:0	18079	18.8824	19.0510	16:16	18.7143	16:29	64.2480	1024.19
21			2010-11-11 16:45:0	18080	18.6345	18.8156	16:31	18.4454	16:44	65.5880	1024.20
22			2010-11-11 17:00:0	18081	18.3188	18.5804	16:45	18.0105	16:59	67.5586	1024.20
23			2010-11-11 17:15:0	18082	17.7915	18.1135	17:2	17.3480	17:14	69.4988	1024.28
24			2010-11-11 17:30:0	18083	16.9765	17.4155	17:15	16.5399	17:29	72.6598	1024.25
25			2010-11-11 17:45:0	18084	16.0813	16.6412	17:30	15.6356	17:44	76.8499	1024.25

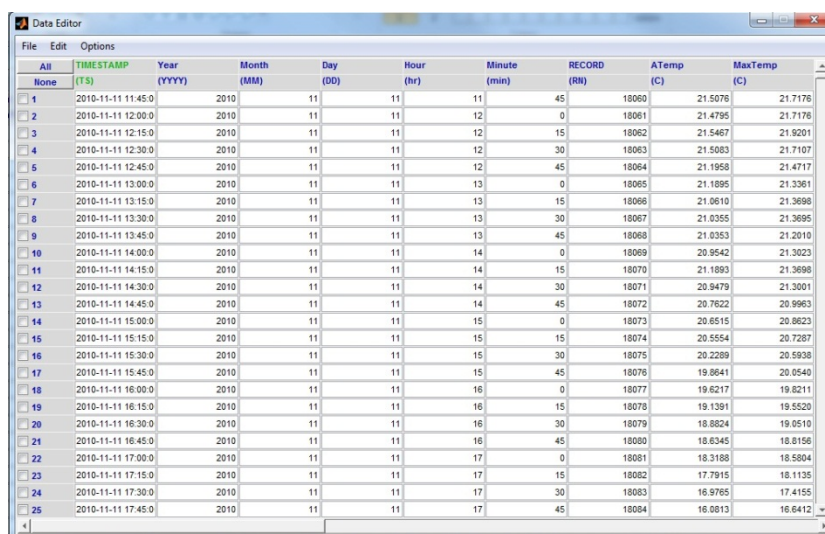
Data loaded in the Data Editor.

- The dataset values can be edited from this screen. Scroll to the record you want to change, click in the appropriate column cells and make any necessary changes, then go to **"File > Return to Editor"** to return the modified data set to the Data Set Editor window. Note that if you close the Data Editor window without returning the data to the Data Set Editor (e.g. clicking on the "X" window close button) the changes will be discarded. Value changes are automatically logged to the metadata, and attribute metadata (data type, precision) are used to validate all edits.
- Clicking on the checkbox next to the row number selects the entire row for copying or deletion. For example, click on the first 3 rows, then select **"Edit > Delete Selected Rows"**.
- Clicking on the **"All"** button selects all rows in the data set, and clicking on **"None"** clears all row selections. Note that row selections are retained as you scroll through the data set, so it is good practice to always click **"None"** before selecting rows to delete.
- Cells for any values assigned QA/QC flags are highlighted in red, and hovering over the cell with the mouse pointer will display the flag(s).
- You can specify criteria to control which data records are displayed to simplify data review or navigating large data sets. For example, use **"Options > Record View > Only Flagged Records"** to filter the list of records to only those containing one or more flagged values. Use **"Options > Record View > All Records"** to return to the standard view.

Editing Date/Time Formats

In addition to basic manual editing of data sets, there are a variety of dataset-based, automated editing functions available in the Data Set Editor under the **"Edit"** menu. Some of the most commonly used are the **Date Functions**, since date and time information are present in nearly every environmental data set but date/time formats used by data loggers and software systems vary widely. Internally, MATLAB uses floating-point serial dates that start at 0 for 00-Jan-0000 00:00:00, but a wide variety of string date formats are supported as well. Various automatic or manual date/time inter-conversions are available under **"Edit > Date Functions"**.

For example, to generate separate numeric columns for year, month, day, hour, etc. from a date column, use **"Edit > Date Functions > Date Components from Date Column > Automatic"**.



All	TIMESTAMP (TS)	Year (YYYY)	Month (MM)	Day (DD)	Hour (hr)	Minute (min)	RECORD (RN)	ATemp (C)	MaxTemp (C)
1	2010-11-11 11:45:0	2010	11	11	11	45	18060	21.5076	21.7176
2	2010-11-11 12:00:0	2010	11	11	12	0	18061	21.4795	21.7176
3	2010-11-11 12:15:0	2010	11	11	12	15	18062	21.5467	21.9201
4	2010-11-11 12:30:0	2010	11	11	12	30	18063	21.5063	21.7107
5	2010-11-11 12:45:0	2010	11	11	12	45	18064	21.1956	21.4717
6	2010-11-11 13:00:0	2010	11	11	13	0	18065	21.1895	21.3361
7	2010-11-11 13:15:0	2010	11	11	13	15	18066	21.0610	21.3698
8	2010-11-11 13:30:0	2010	11	11	13	30	18067	21.0355	21.3695
9	2010-11-11 13:45:0	2010	11	11	13	45	18068	21.0353	21.2010
10	2010-11-11 14:00:0	2010	11	11	14	0	18069	20.9542	21.3023
11	2010-11-11 14:15:0	2010	11	11	14	15	18070	21.1893	21.3698
12	2010-11-11 14:30:0	2010	11	11	14	30	18071	20.9479	21.3001
13	2010-11-11 14:45:0	2010	11	11	14	45	18072	20.7622	20.9963
14	2010-11-11 15:00:0	2010	11	11	15	0	18073	20.6515	20.8623
15	2010-11-11 15:15:0	2010	11	11	15	15	18074	20.5554	20.7287
16	2010-11-11 15:30:0	2010	11	11	15	30	18075	20.2289	20.5938
17	2010-11-11 15:45:0	2010	11	11	15	45	18076	19.8641	20.0540
18	2010-11-11 16:00:0	2010	11	11	16	0	18077	19.6217	19.8211
19	2010-11-11 16:15:0	2010	11	11	16	15	18078	19.1391	19.5520
20	2010-11-11 16:30:0	2010	11	11	16	30	18079	18.8824	19.0510
21	2010-11-11 16:45:0	2010	11	11	16	45	18080	18.6345	18.8156
22	2010-11-11 17:00:0	2010	11	11	17	0	18081	18.3188	18.5804
23	2010-11-11 17:15:0	2010	11	11	17	15	18082	17.7915	18.1135
24	2010-11-11 17:30:0	2010	11	11	17	30	18083	16.9765	17.4155
25	2010-11-11 17:45:0	2010	11	11	17	45	18084	16.0813	16.6412

Results from the "Automatic" Date/Time function.

If the Automatic function does not work properly (i.e. an error is displayed), the most likely cause is that the date column (TIMESTAMP in this case) was not properly classified as a date/time variable in the attribute metadata. For the example data set, check the Variable Type designation of TIMESTAMP and change it to **"date or time (datetime)"** if necessary, then run the "Automatic" date function again. Once you make the correction and rerun the function, you will see that five new fields (Year, Month, Day, Hour, and Minute) are added to the dataset and now appear in the dataset editor window.

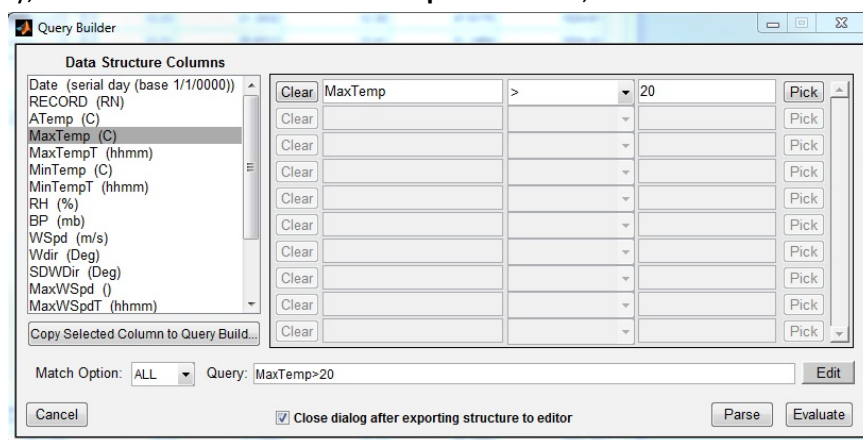
Date Padding

Chronological gaps often occur in time-series data sets when logging is interrupted or instruments are swapped, resulting in discontinuous data with uneven time steps. The GCE Toolbox can pad these gaps with appropriate missing values to create a continuous (monotonic) time series data set. Date/time values are automatically generated, and values in non-data columns (e.g. instrument or site codes) can be replicated, if desired. To pad date gaps, select **"Edit > Date Functions > Fill in Date Gaps (time series data) > Do Not Replicate Values."** This function will add empty records with only Date/Time values to the existing dataset. Note that this function also adds a serial date field, recording the date as fractional

serial day based on 1 being equal to midnight on January 1, 0000, or the beginning of the Common Era (C.E.) according to the Gregorian calendar.

Filtering Data

1. Data sets can be filtered based on queries to create a subset. From the Data Structure Editor, go to **“Tools > Filtering > Filter/Subset Data by Column Values”** to bring up the query builder.
2. For this example, we want to create a new data set that contains records where the maximum temperature was over 20 degrees Celsius. In order to do this, double-click on the MaxTemp column in the Data Structure Column window. Select, the greater than sign from the drop down menu. We want Data from records higher than 20 degrees, so enter **“20”** into the third field of the query, so that it will look like **“MaxTemp > 20”**. Next, click the Evaluate button.



The filtering query builder.

3. A new Data Structure Editor window will pop up containing just the data from the query. You can view the data in the new data set by going to **“Edit > View/Edit Data.”**

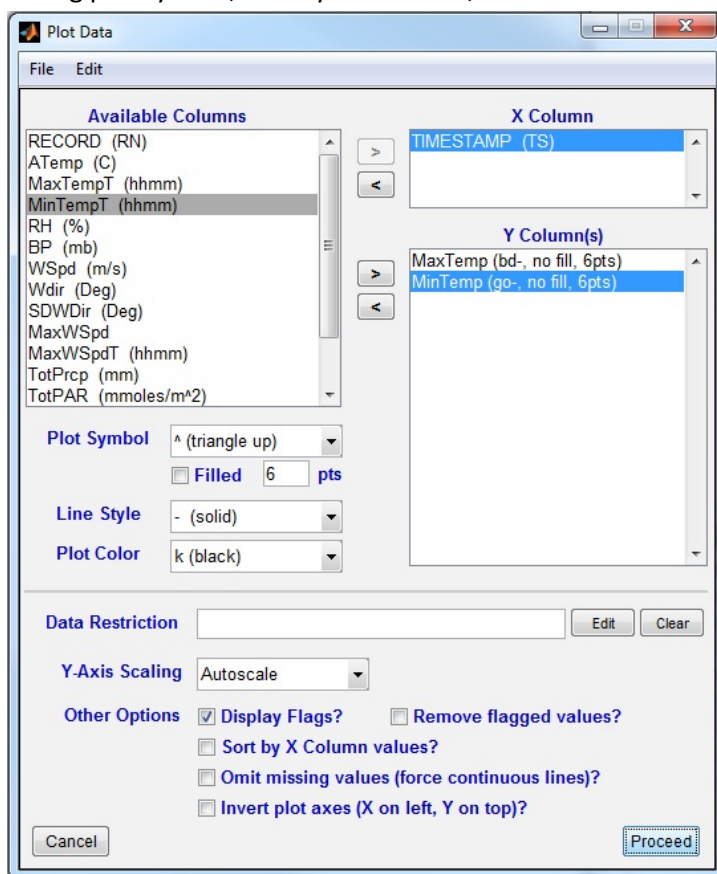
Viewing Data Statistics

1. A quick summary of dataset statistics by column can be viewed from the Data Structure Editor by going to **“Tools > Statistics > View Column Statistics > Include Flagged Values”** or **“... > Exclude Flagged Values”**. A scrolling window will pop up displaying general statistics for each column in the data set.
2. To generate a formatted report of basic column statistics, use **“Tools > Statistics > Column Statistics Report”** and specify the filename, format and options.
3. You can also generate derived statistical summary data sets by Grouping, Binning, Date/Time Interval or Moving Date Interval by selecting the corresponding option on the **“Tools > Statistics”** menu and filling in options on the form that is displayed.

Plotting Data

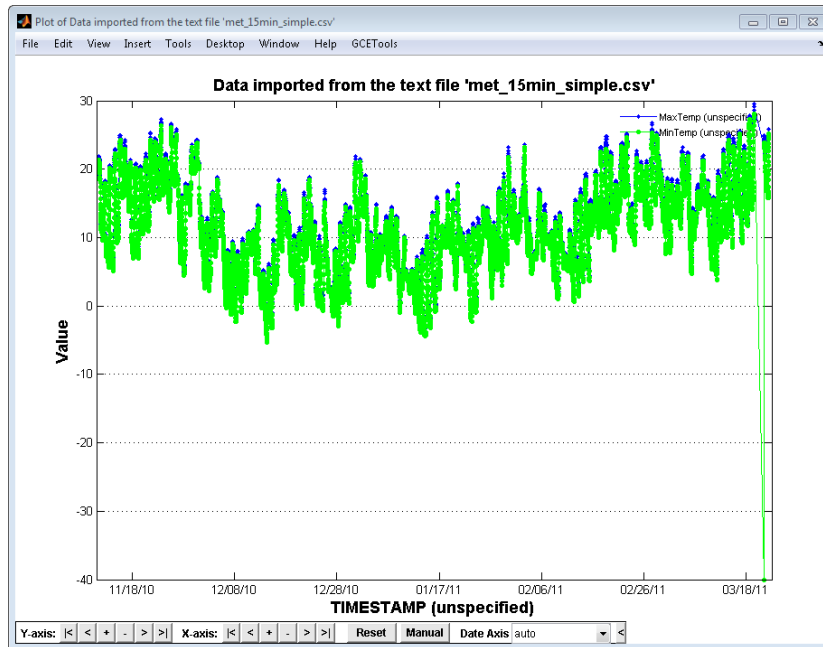
The GCE toolbox has several plotting tools to quickly generate graphs of data sets.

1. Using the met_15min_toa5.dat dataset you already loaded into the Data Structure Editor, go to **“Tools > Plotting > 2D line/Symbol (Multiple Y).”**
2. The new window that pops up allows you to choose which columns to plot and select the corresponding plot symbol, line style and color, as well as Y-Axis scaling and other options.



Plot Creation tool.

3. For this exercise, select the Date column generated during the date padding operation and move it to the X column box using the “>” button. Next, choose the MaxTemp and MinTemp columns and move them to the Y columns box in the same manner. Click the **“Proceed”** button. The graph will display in a new window.



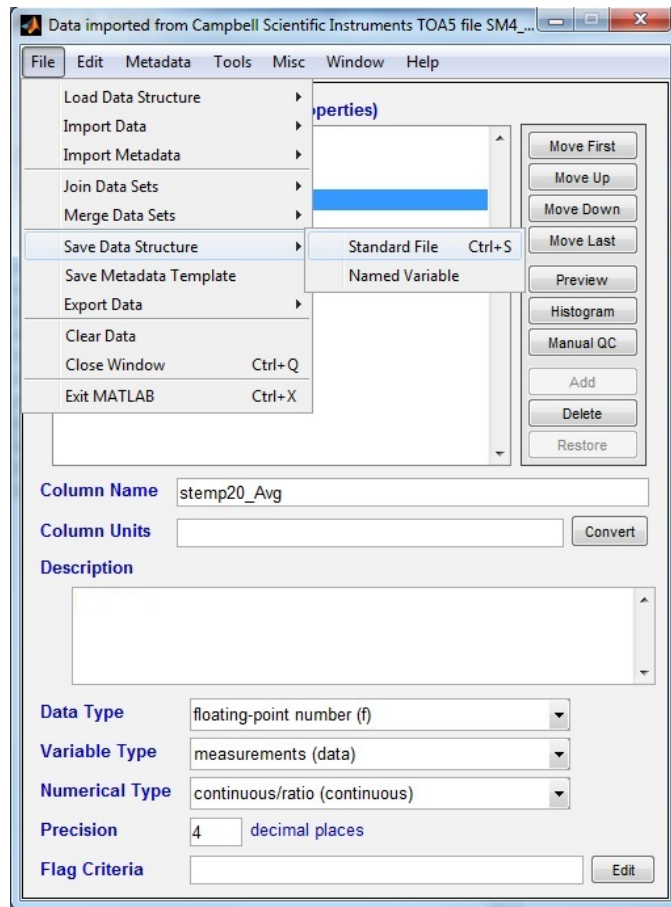
Example plot of MaxTemp and MinTemp data.

4. Toolbar buttons below the plot can be used to zoom and pan through the dataset and change the format of date ticks. Numerous other options for customizing the plot and doing simple curve fitting are available from the menu bar.
5. The graph can be exported by going to **"GCETools > Export Plot"** and specifying the format and resolution desired (note: **"File > Save As"** can also be used, but the toolbar and other graphical controls will export or print as well)

Saving and Loading GCE Toolbox (.mat) Files

Once you have imported a data set and made changes to it, you can save your work. Go to **"File > Save Data Structure > Standard File"** then navigate to the directory where you wish to save the file and specify a filename. This will save the data set as the variable "data" in a MATLAB binary file (.mat). Save the imported data as "met_15min_simple.mat" in the workshop_products directory for this exercise.

From this point on it will not be necessary to import the raw data from the original text file into the Toolbox unless you wish to start over.



Saving the data as a .mat file.

Loading GCE Toolbox (.mat) Files

Once you have created a .mat file for a dataset you can load it into the Data Structure Editor by going to **"File > Load Data Structure > Load Structure from File"** then selecting the .mat file you wish load. Data sets saved using the "Standard File" option should load automatically; otherwise a variable selection dialog will be displayed if more than one data structure variable is present in the file.

Identifying Empty Columns, Rows, and Duplicate Dates

The GCE Toolbox has built in functionality to check data sets for numerous errors and missing values. In this short exercise, we will load a data set with duplicate dates and empty columns, remove them, and save the new data set.

1. Begin by loading the met_15min_simple_dupes.mat file. Once the file is loaded into the Data Structure Editor, go to the Data Editor screen.
2. From the Data Editor screen, select **"Options > Record View > Only Duplicate Records > Date Time Columns Duplicated."** This will now display only the records with duplicated date-time values in the Data Editor.

3. In this case, all columns of the duplicated records are the same, so we can delete either one of the duplicates. Select one of the duplicated records, then go to ***“Edit > Delete Selected Rows”*** to delete the record. (Note that both records disappear, because removing the duplicate excludes the remaining record from the duplicate record display filter)
4. Once all the duplicate records have been removed, save the edited data set by going to ***“File > Return to Editor”*** and then saving the file normally from the Data Structure Editor. Again, not using Return to Editor will cause your edits to not be saved.
5. While it is usually preferable to remove duplicate records from the Data Editor, because you are able to see the records that will be removed, you can also remove them all at once using ***“Edit > Remove Duplicate Records”*** and then either selecting ***All Columns Duplicated*** or ***Non-data Columns Duplicated***. The total number of records removed will be displayed in a message box.
6. Additionally, you can remove records or entire columns that do not contain any data values. Columns are removed by going to ***“Edit > Remove Empty Columns”***, while removing empty records is done by going to ***“Edit > Remove Empty Records”*** and then choosing the appropriate option. ***All Columns Empty*** will remove records where no columns have any valid data values (i.e. other than NaN or an empty string), ***All Data Columns Empty*** will remove records where no columns with Variable Type of 'data' or 'calculation' have any valid data values (e.g. to remove empty records added by "Pad Date Gaps"), and ***Selected Columns*** removes records where none of the selected columns in the Column List have any valid values (note: use Ctrl-click, Shift-click, or Command-click, as appropriate, to select multiple columns).

Loading a Single TOA5 .dat File

The GCE Toolbox contains a number of specialized import filters for commonly encountered file formats from environmental data loggers and online databases. These filters are m-file functions that contain source-specific logic for analyzing and parsing the data, generating appropriate attribute metadata and performing common post-processing. Import filters are therefore pre-built data processing workflows that can greatly simplify data processing using the GCE Data Toolbox software.

The met_15min_toa5.dat file that was used in the custom ASCII import example is actually in the Campbell Scientific Instruments table-oriented ASCII (TOA5) format. A generalized Campbell TOA5 import filter comes pre-installed with the toolbox, so we will now repeat the data import processing using this filter instead of the generic delimited ASCII filter.

1. Bring up the Data Structure Editor window and select ***“File > Import Data > Campbell Scientific TOA5 Data > Any Station (generic template)”*** and then navigate to the met_15min_toa5.dat file and click the Open button. The file will now be imported into the Data Structure Editor as before.
2. Note that a MATLAB "Date" column was automatically generated from the timestamp, column names and units were automatically parsed from the header, the missing value code was

automatically recognized, and basic column descriptions were generated from ancillary information in the Campbell header (e.g. whether the measurement was totaled, instantaneous, from a vector product, etc.).

Loading a Campbell Scientific CR10X File

If the data that you are importing is in the Campbell Scientific Instruments array format (e.g. CR10x data), you will need to first create a specialized template file for the csi2struct.m import filter in order to successfully load the data. Individual arrays will be split from the logger file automatically, then processed and documented using information in the template to produce a GCE Data Structure for each array.

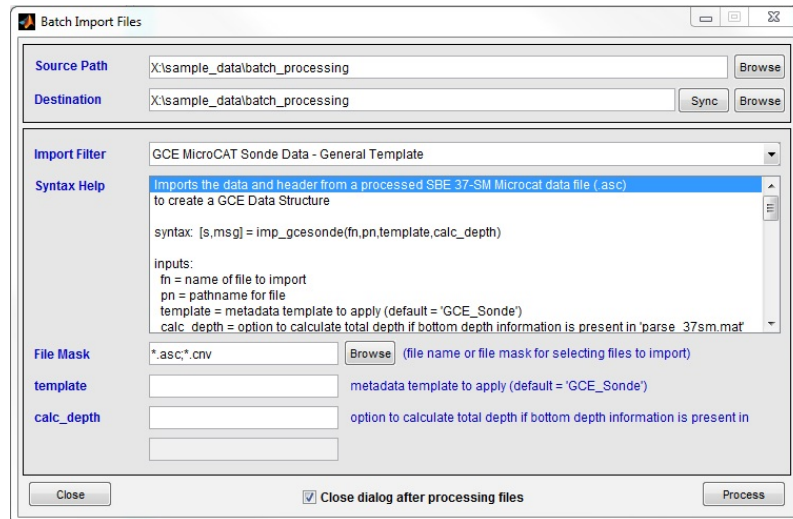
A csi2struct.m template is actually a stand-alone GCE Data Structure that contains details for each array column that may be present in a file (matched by array ID and column position) as data set columns. Boilerplate documentation metadata in the template is applied to each processed data set. An example file is included in the toolbox distribution (\userdata\csi2struct.mat).

Due to the specialized nature of this process no exercises are planned, but it may be demonstrated if time and interest permits.

Batch Importing Files

If multiple files need to be imported into the GCE Toolbox, the batch import function can process all of them in a single operation. Batch importing will create a corresponding MATLAB .mat file for each raw data file with the selected import filter and metadata template already applied. If you wish to use the batch processing tool, the files that will be processed need to be in the same directory.

1. This exercise will use files located in the Batch Processing subdirectory of the workshop products directory. These files represent periodic downloads from two different moorings with similar sensor packages installed. Note that logger data formats vary over time due to differences in sensor firmware versions (i.e. .cnv and .asc), but the import filter used for this data source is able to import both formats automatically.
2. From the Data Structure Editor Window, go to **"File > Batch Import Files"**.
3. In the new window that pops up, browse to the Batch Processing directory in the Source Path field.
4. Select the destination directory for the imported files in the Destination field. For illustration purposes, we recommend creating a batch_products subdirectory in the workshop_products directory. However, the destination directory can be the same as the source, and you can copy the source path to the Destination field using the "Sync" button.
5. Choose which import filter to use with the data that will be imported. The data files for this exercise can be imported using the **GCE MicroCAT Sonde Data – General Template** filter entry, so choose that from the dropdown menu. Note that this filter uses a file mask consisting of .asc and .cnv files, so it will not process files that have other extensions.



The Batch Processing window.

6. Hit the **Process** button. The files will now be imported and the resulting data structures will be saved in MATLAB .mat format in the destination directory. A text report describing the results will also be generated and displayed.
7. To load the imported files into the Data Structure Editor, use "**File > Load Data Structure > Load Structure from File**" as usual

Basic Metadata

This section will introduce the basic process of managing data set metadata using the GCE Data Toolbox. In addition to being a good practice to fully document environmental data, the GCE Toolbox uses attribute metadata to control and automate all data processing and analysis. Generating correct and complete metadata is therefore crucial to successful use of the toolbox.

Creating Attribute Metadata

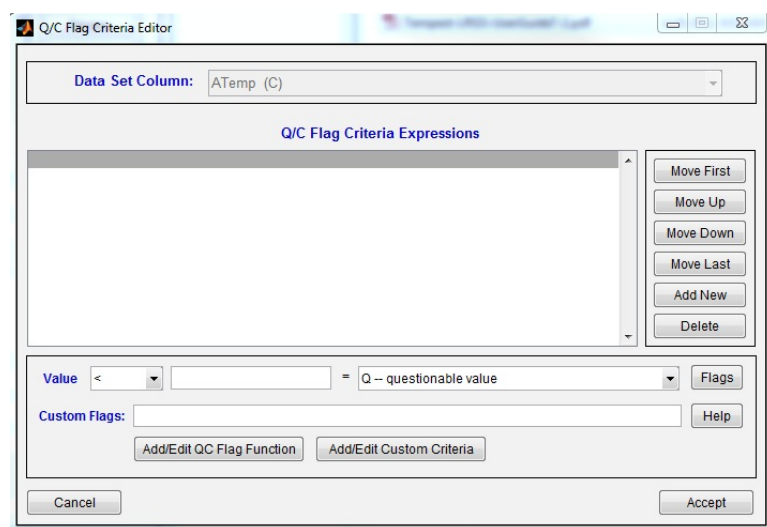
You can edit the attribute metadata for a data set from the Data Structure Editor. During a data import, the GCE Toolbox will assign information to each attribute based on the import criteria or filters that were used. It may be necessary to change the column units, data type, variable type, numerical type, or add to descriptions or make other changes. To do this, simply load the data into the Data Structure Editor, select the attribute that you wish to edit from the column list, and make the desired changes in the fields and drop down menus.

When a generic import filter is used (e.g. Delimited Text File, MATLAB Data File) It is particularly important to check the contents of attribute metadata fields for accuracy and to set an appropriate Variable Type for non-data columns (e.g. dates, geographic coordinates, coded columns). Tools and functions in the GCE Toolbox use attribute metadata to configure dialogs and process and display the data values, so inappropriate Variable Type settings can lead to unexpected errors (e.g. unrecognized date/time columns or calculation of inappropriate statistics in summary data).

Creating Basic QA/QC Rules

The GCE Toolbox has an interactive tool for designing QA/QC rules for attributes. These rules can include numeric conditional checks, column cross-reference checks, mathematical comparisons, and statistical checks. This example will cover how to create a simple limit check rule using the ATemp attribute in the GCE structure file you derived from the met_15min_toa5.dat data set.

1. From the Data-Structure Editor, with the met_15min_toa5.dat-derived structure file loaded, select the ATemp attribute then click the “**Edit**” button to the right of the Flag Criteria field at the bottom of the screen. This will bring up the Q/C Flag Criteria Editor.



The Q/C Flag Criteria Editor.

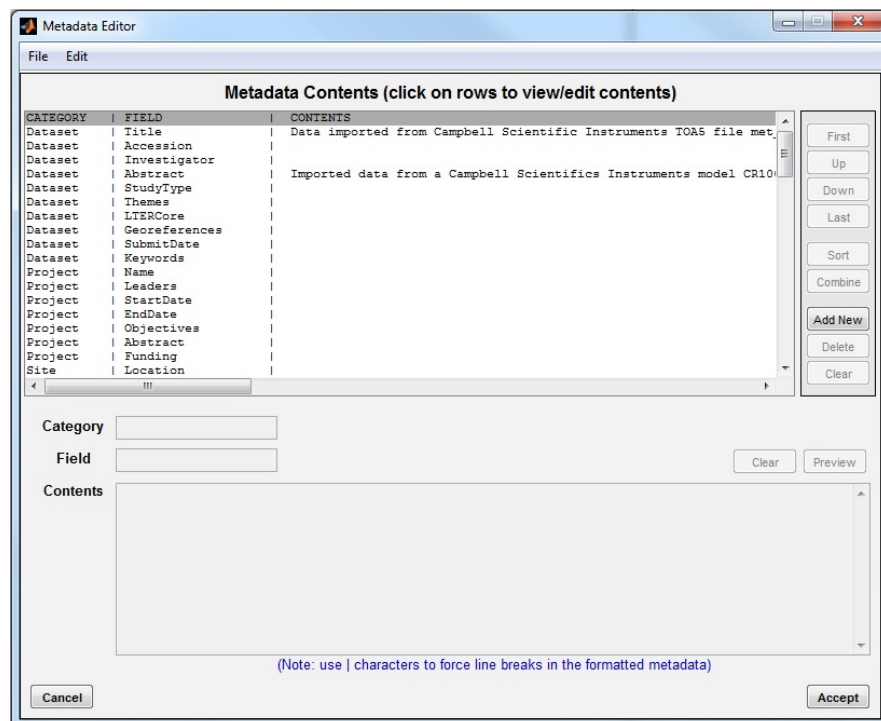
2. This editor will allow you to create basic QC rules similar to the data filter query builder. In the Value drop-down menu, select the ">" (greater than) sign, and in the field next to that, enter a threshold value of 40. Now, select the “**Q – questionable value**” flag from the drop down menu next to the equals sign. Click in the Expressions window to transfer the flagging rule to the Expression list. This new rule will now flag all ATemp values that are higher than 40 with a flag of “Q”.
3. Repeat these steps using the "<" (less than) sign and a threshold of -10 to create a rule for flagging values that are below -10.
4. Hit the “Accept” button, and the new flagging rules will be added to the attribute metadata for the ATemp column in the Data Structure Editor. These new rules will be evaluated automatically when the data set is saved or the data are viewed or plotted.
5. You can define additional flags to assign by hitting the “Flags” button on the right side of the window.
6. There are a wide variety of custom flagging criteria that can be created in the Q/C Criteria Editor. Hit the “Help” button to bring up documentation that further explains the functionality of Custom Flags.

Using the Metadata Editor

Each data set created as a MATLAB structure in the GCE Toolbox has both attribute and documentation metadata associated with it and stored as part of the structure. If you export data to other formats from the Toolbox, you can export the metadata in various formats as well. Documentation metadata fields are generally empty if the data set has just been imported, although some fields may be filled out if a specialized import filter has been used. In order to create a complete metadata file, the metadata will either have to be manually entered, a metadata template will need to be applied, or metadata will need to be imported from another data structure file.

Important Note: The documentation metadata schema used by the GCE Data Toolbox is currently being updated for improved parity with EML 2.1 and the Metabase Metadata Management System. The new schema will provide better control over field contents, sub-fields for specific elements (e.g. name fields, method fields, instrument fields), and selective repeatability of fields. The exercises below are based on the current "loose" schema where metadata fields are operationally defined based on user-editable style definitions.

1. Attribute metadata can be edited by loading the data set into the Data Structure Editor then defining the name, units of measurement, description, data type, variable type, numerical type and precision of each column as previously described.
2. To edit the documentation metadata from the Data Structure Editor, go to “*Metadata > View/Edit Metadata*” to bring up the Metadata Editor window.



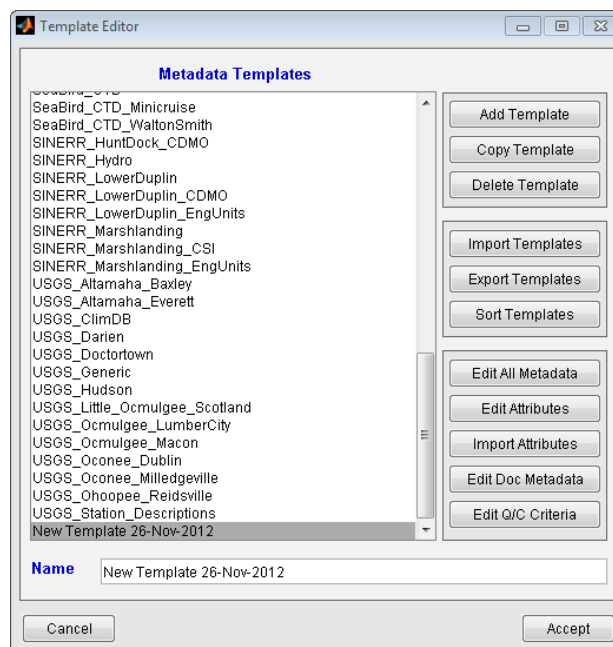
Metadata Editor with initial metadata from TOA5 import.

3. From the Metadata Editor, you can select various metadata fields (organized by category) and add text in the Contents box. By default, fields from the "LTER-FLED" style are added when a new structure is created, but additional metadata fields can be added using the **"Add New"** button to the right, or fields can be deleted using the **"Delete"** button.
4. Fields can be reorganized using the **"First"**, **"Up"**, **"Down"**, and **"Last"** buttons to simplify editing, but note that metadata field position is not critical when metadata are styled for display or export.
5. Metadata fields and content can also be imported from pre-defined metadata templates or existing data structures using **"File > Import Fields"** and **"File > Import Metadata"**, resp.
6. Once you have made the desired changes to the metadata, click the **"Accept"** button to accept the changes and return to the Data Structure Editor. You will still need to save the structure file in order to permanently save the metadata updates.

Creating and Applying Metadata Templates

Metadata templates can be created to apply attribute and documentation metadata to similar data sets automatically. Templates can be applied to data sets at various points in the workflow, but are particularly effective when combined with source-specific import filters, allowing documented, quality-controlled data sets to be produced simply by importing a raw data file. The following steps cover how to create a metadata template from a new data set.

1. Begin by importing an undocumented data set into the Data Structure Editor and completing the data set metadata in the Metadata Editor as described above.
2. Once the metadata are completed, from the Data Structure Editor go to **"File > Save Metadata Template"**, which brings up the Metadata Template Editor with a new entry.



Metadata Template Editor

3. Enter a name for the new template and click accept to create and save the new template.
4. Once a metadata template has been created, it can be applied to a new dataset. To do this, import or load a raw data set into the Data Structure Editor, then select "**File > Import Metadata > Standard Template**" and select the new template from the list.
5. To import documentation metadata without changing attribute metadata, open the Metadata Editor and select "**File > Import Metadata > Metadata Template > Overwrite All Fields.**" This will cause a new window to pop up where you choose the template that you wish to apply. Select the metadata template, and then click the "Ok" button. The documentation metadata from the template will now be applied to the current data set.

Managing Metadata Templates

Once a metadata template has been created, the GCE Toolbox has a number of tools for managing and updating the template contents. In addition, empty templates can be created and populated *de novo*, and templates can be derived from existing templates to re-use metadata content. Templates can also be exported to or imported from other toolbox instances. The following exercise will acquaint you with the basic features of the Metadata Template Editor application

1. Open the Metadata Template Editor from the Data Set Editor window by going to "**Misc > Add/Edit Metadata Templates**".
2. To select an existing template to edit, click on the template name in the Metadata Templates list
3. To edit all template metadata contents (i.e. attribute metadata, documentation metadata and QA/QC rules) together, click on the "Edit All Metadata" button. This will open the template contents in a Data Structure Editor instance, but with data-dependent features and menus disabled (e.g. "**Preview**", "**Histogram**", and "**Manual QC**" buttons).
 - a. Click on an entry in the "Template Variable Name and Column Name List" to edit attribute metadata and/or QA/QC criteria. The raw data file variable to match and column name to assign are combined into a "Variable==Name" field; to change the name assigned to a raw data column edit the text after "==", but take care editing the Variable name to ensure that the attribute metadata are correctly matched to the raw data column.
 - b. If the variable name to match and column name to assign are the same, a single name without "==" can be used (e.g. "Temp_Air" instead of "Temp_Air==Temp_Air")
 - c. Changes to other attribute metadata fields and QA/QC Flag Criteria are made as for imported or loaded data sets in prior exercises.

Template Attributes Editor

File Metadata Window Help

Template Variable Name and Column Name List

MaxTemp==MaxTemp (C)
 MaxTempT==MaxTempT (hhmm)
 MinTemp==MinTemp (C)
 MinTempT==MinTempT (hhmm)
 RH==RH (%)
 BP==BP (mb)
 WSpd==WSpd (m/s)
 Wdir==Wdir (Deg)
 SDWDir==SDWDir (Deg)
 MaxWSpd==MaxWSpd (m/s)
 MaxWSpdT==MaxWSpdT (hhmm)
 TotPrcp==TotPrcp (mm)
 TotPAR==TotPAR (mmoles/m²)
 AvgVolt==AvgVolt (Volts)
 CumPrcp==CumPrcp (mm)
 TotRad==TotRad (mmoles/m²)

Move First
 Move Up
 Move Down
 Move Last
 Preview
 Histogram
 Manual QC
 Add
 Delete
 Restore

Variable==Name Date==Date

Column Units serial day (base 1/1/0000)

Description
 Fractional serial day (based on 1 = January 1, 0000)

Data Type floating-point number (f)

Variable Type date or time (datetime)

Numerical Type continuous/ratio (continuous)

Precision 6 decimal places

Flag Criteria Edit

Template Attribute Metadata Editor with example data loaded.

- Changes to attribute metadata alone can be made by hitting the **“Edit Attributes”** button. This will bring up a table (i.e. Data Editor grid) containing all of the attribute metadata in rows and columns. The values can be changed by selecting the field, making the change, and then using **“File > Return to Editor”** to update the attribute metadata in the specified template.

Data Editor									
All	Variable	Name	Units	Description	Data Type	Variable Type	Numerical Type	Precision	Criteria
None	(Variable name in the source (Column name to assign))	(Column measurement units (Description of the column))	(Data storage type: f for	(Data storage type: f for	(Variable type: data,	(Numerical type: none for	(Number of digits to display)	(Quality control)	
1	Date	Date	serial day (base 1/1/0000)	Fractional serial day (based on 1 =	f	datetime	continuous	6	
2	RECORD	RECORD	hhmm	Measurement of RECORD	d	discrete	discrete	0	
3	ATemp	ATemp	C	Averaged measurement of ATemp	f	data	continuous	4	
4	MaxTemp	MaxTemp	C	Instantaneous measurement of ATemp	f	data	continuous	4	
5	MaxTempT	MaxTempT	hhmm	Instantaneous measurement of ATemp	f	nominal	none	0	
6	MinTemp	MinTemp	C	Instantaneous measurement of ATemp	f	data	continuous	4	
7	MinTempT	MinTempT	hhmm	Instantaneous measurement of ATemp	f	nominal	none	0	
8	RH	RH	%	Averaged measurement of RH	f	data	continuous	4	
9	BP	BP	mb	Averaged measurement of BP	f	data	continuous	2	
10	WSpd	WSpd	m/s	Measurement of WSpd	f	data	continuous	4	
11	Wdir	Wdir	Deg	Measurement of Wdir	f	data	continuous	3	
12	SDWDir	SDWDir	Deg	Measurement of SDWDir	f	data	continuous	4	
13	MaxWSpd	MaxWSpd	m/s	Instantaneous measurement of WSpd	f	data	continuous	4	
14	MaxWSpdT	MaxWSpdT	hhmm	Instantaneous measurement of WSpd	f	nominal	none	0	
15	TotPrcp	TotPrcp	mm	Totalled measurement of TotPrcp	f	data	continuous	5	
16	TotPAR	TotPAR	mmoles/m ²	Totalled measurement of TotPAR	f	data	continuous	2	
17	AvgVolt	AvgVolt	Volts	Averaged measurement of AvgVolt	f	data	continuous	4	
18	CumPrcp	CumPrcp	mm	Instantaneous measurement of CumPrcp	f	data	continuous	4	
19	TotRad	TotRad	mmoles/m ²	Totalled measurement of TotRad	f	data	continuous	3	

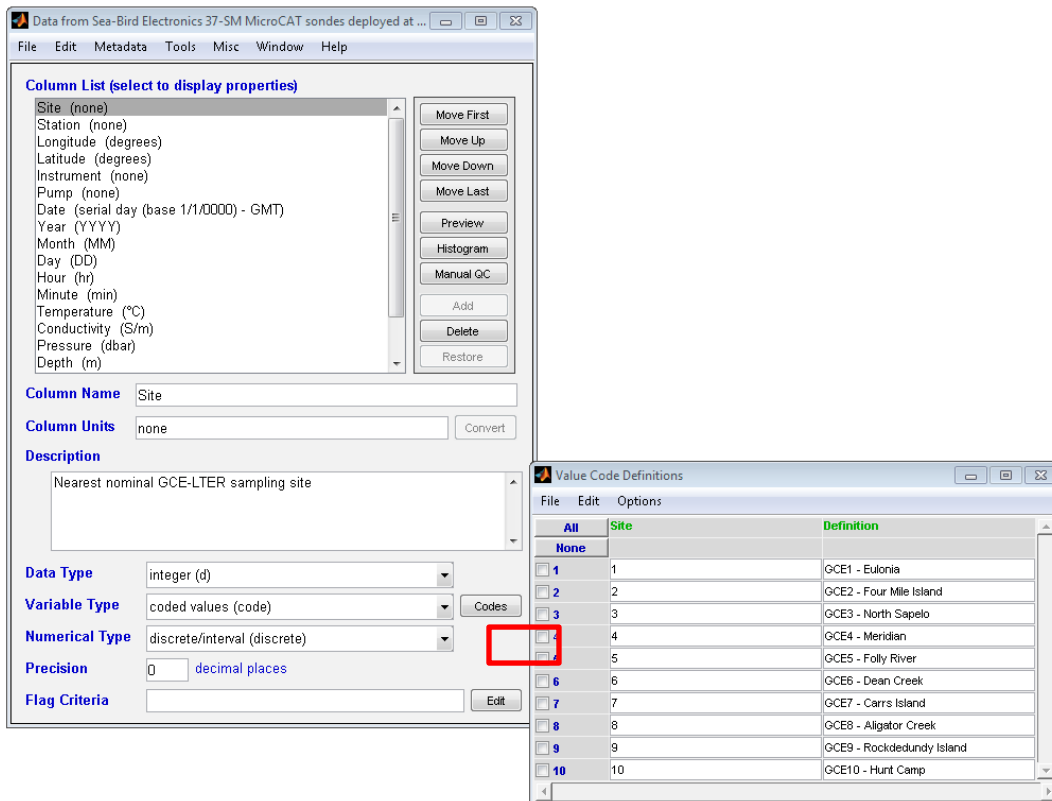
Attribute metadata editor (“expert” mode)

- Edits to the documentation metadata alone can be made by hitting the **“Edit Doc Metadata”** button. This will bring up the standard Metadata Editor window.

- Changes to Q/C criteria alone can be made by hitting the **"Edit Q/C Crit"** button. This will bring up the standard Q/C criteria window where the Q/C expressions can be edited. The drop-down menu is unlocked, unlike when the Q/C criteria editor is opened from the Data Set Editor, allowing rules for any column in the template to be edited.
- Attribute information can be imported from other data set files to the current template by hitting the **"Imp Attributes"** button. Select to the desired file using the file loading dialog and press **"Open"** to import attribute metadata and open it in a Data Editor grid for inspection. Use **"File > Return to Editor"** to save the changes to the selected template.

Handling Coded Attributes

- For coded columns, codes and code definitions are stored in the documentation metadata (i.e. Data/ValueCodes field), but code definitions can be managed as attribute metadata in the Data Structure Editor. After loading a data set in the Data Structure editor, select the column that contains the coded values. In the Variable Type drop-down menu, select Coded Values. This will cause a **"Codes"** button to appear to the right of the drop-down menu. Hit the **"Codes"** button to bring up the Value Code Definitions window (e.g. see "gce9_hydro_realtime_2012.mat").



- From this window you can see the codes present in the selected column. The definitions can be edited by selecting the Definition Field next to the code value and entering the definition. Once the definitions have been filled out, use **"File > Return to Editor"** to save the changes.

QA/QC Framework

The GCE Toolbox is capable of performing very complex QA/QC checks on data. Everything from simple limit checks through complex, parameterized models that load external reference data can be leveraged to flag data values. The Q/C Criteria Editor (see above) can be used to define many common QA/QC checks, including criteria based on custom MATLAB functions and multi-column dependency checks, but keep in mind that custom Q/C rules are essentially unlimited in scope.

In addition to Q/C rules (i.e. algorithmic checks), Q/C flags can also be assigned and cleared visually on data plots, copied from one or more columns to dependent columns, and imported from text columns. Once flags are assigned by rules or manual operations, many options are provided for managing the flagged data values.

Note that the GCE Toolbox does not impose any particular Q/C flag vocabulary or assume any flag semantics. Definition and interpretation of flags is under control of the data provider and workflow developer. Although flags are often used to qualify problematic data values, flags could also be assigned to signify "good" or reviewed data values.

The examples below provide a quick introduction to the QA/QC framework provided in the GCE Data Toolbox, but the user is referred to the toolbox documentation for more information on individual features.

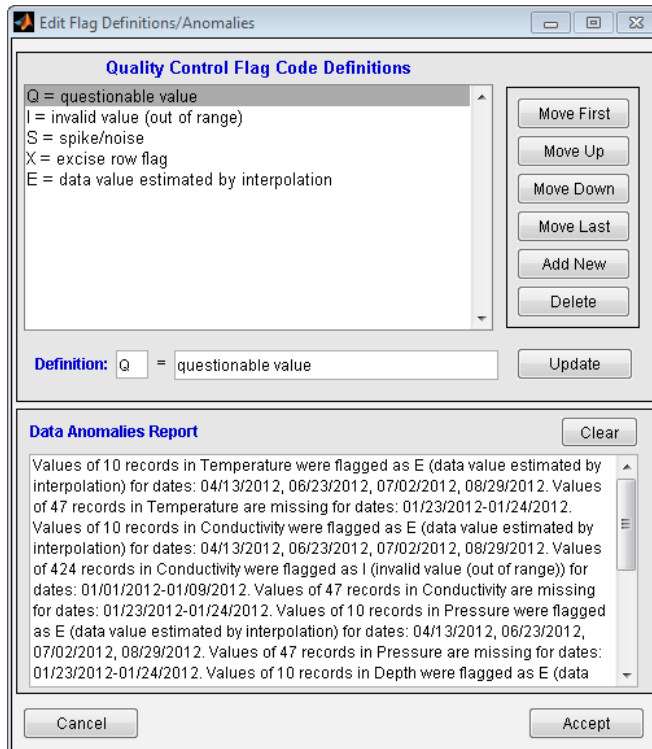
Managing Flagged Data

1. Editing Q/C Flag Definitions

Flag codes and definitions are managed on a per-dataset level using the Edit Flag Definitions/Anomalies window, which is accessed from the Q/C Criteria Editor or from the Data Structure Editor via ***Edit > Q/C Flag Functions > View/Edit Q/C Flag Functions***. In the new window that pops up, you can enter new flags and flag definitions, and manage existing ones. New flags can be entered by entering the flag code in the first Definition field, and entering the Definition in the second field, then hitting the "Update" button.

2. Generating a Data Anomalies Report

In order to generate a human-readable data quality report, go to **Metadata > Document Flagged Values as Anomalies** or **Metadata > Document Flagged and Missing Values as Anomalies** and then select grouping option and date/time format. This will bring up a new window where you select which columns you want the report to cover, or you can select all columns and click the OK button. The report will be displayed in the Data Anomalies Report field of the Edit Flag Definitions/Anomalies window for review and editing.

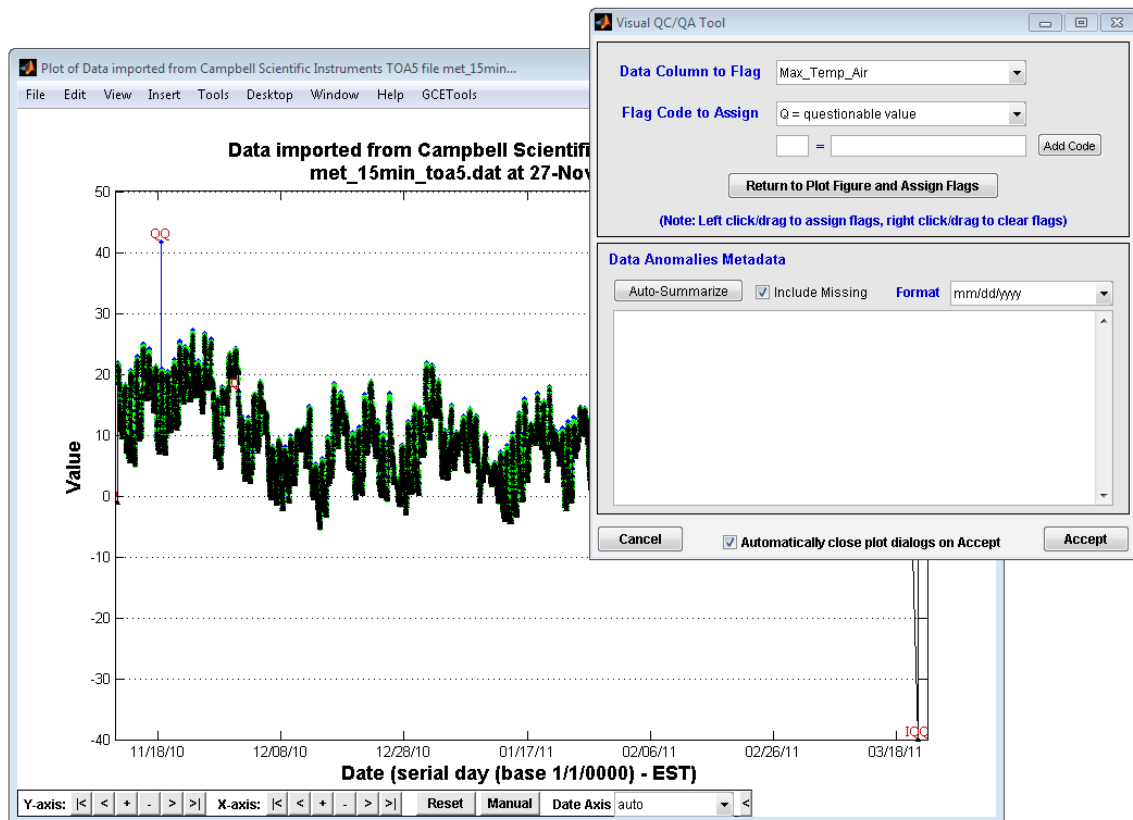


Flag Definition and Data Anomalies Editor Window

3. Visual QA/QC Flagging on Data Plots

Algorithmic flagging based on rules is a powerful QA/QC technique, but flags are often inappropriately assigned during extreme events or not assigned when subtle errors occur that are difficult to develop algorithms to detect. The GCE Toolbox allows you to review and then assign and clear QA/QC flags by plotting the data. After creating a plot (see above), click on **"GCETools > Visual QC Tool Window"** to launch the visual Q/C control panel. Select a variable to flag, choose a flag to assign (or clear), then click on **"Return to Plot Figure and Assign Flags"**. You can now left-click or drag to assign the specified flag, or right-click or drag to clear flag assignments on the chosen variable. Note that portions of an original flag may still be visible after removal (ghosting) due to graphic refresh issues, but the change will be correctly reflected in the underlying data set. If you want to revise flags for another column, return to the control panel and change your selections and then click on the return button again. If desired, you can describe rationale for assigning or clearing flags in the **"Data Anomalies Metadata"** field, and auto-summarize flag assignments (see Data Anomalies Report above).

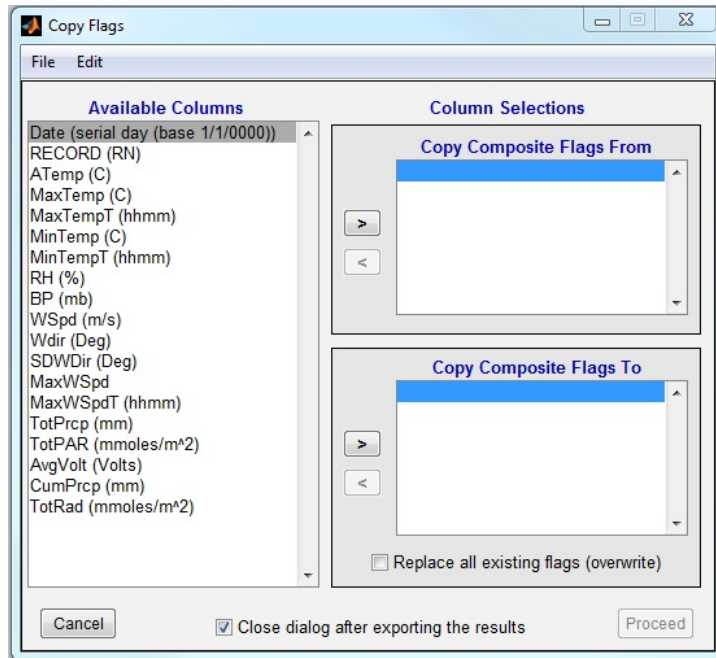
After making the necessary revisions on the plot, return to the control panel and click **"Accept"**. The control panel and plot will then be closed, and the revised data structure will be opened in a new Data Structure Editor window. You can close the original editor window and proceed with the revised version, or save both versions to disk as independent data sets.



Visual QA/QC on a Data Plot

4. Copying QA/QC Flags to Dependent Columns

The quality of a derived or calculated data column is obviously dependent on the quality of the primary measurement columns from which it is derived (e.g. Salinity is calculated from Temperature, Conductivity and Pressure, so problems with any of these measurements affects Salinity as well). As an alternative to defining complex, multi-column criteria in the dependent column, you can copy flags assigned to the primary data columns to dependent columns using **"Edit > Q/C Flag Functions > Copy Q/C Flags to Dependent Columns"**. Select the primary data columns and add them to the "Copy Composite Flags From" list using the ">" button. Next select the dependent column and add it to the "Copy Composite Flags To" list and specify the overwrite option as appropriate. Hit the proceed Button and the flags will be added to the dependent column, augmenting or overwriting existing flags.



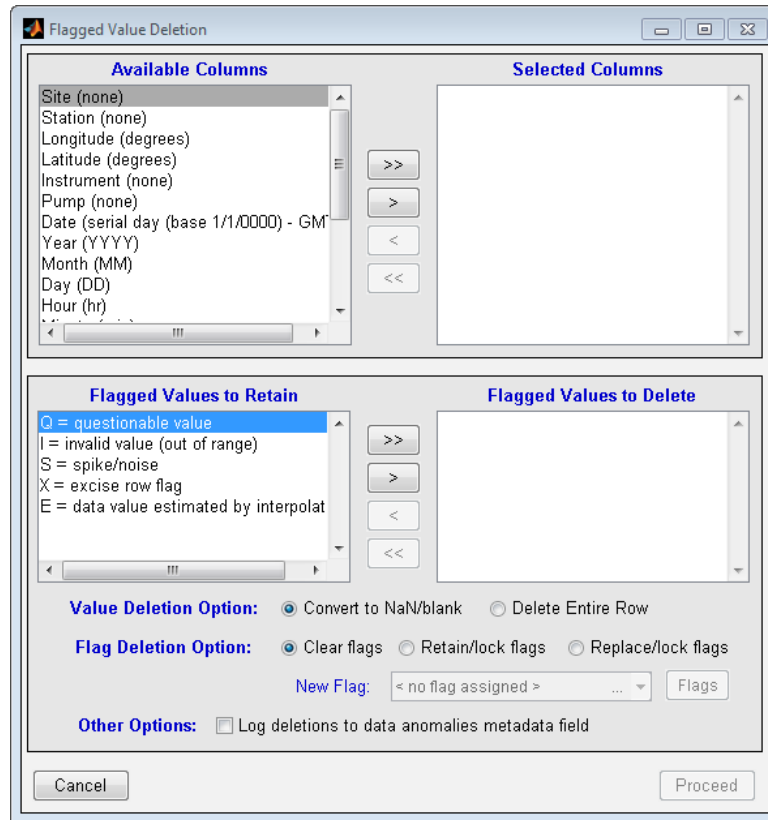
The Copy Flags to Dependent Columns Window.

5. Locking and Unlocking Flags

When raw data are first imported and Q/C rules are defined, the Q/C rules are automatically evaluated whenever rules or data values are changed to set or clear flags accordingly (i.e. flags are in an "unlocked" state). However, when flag assignments are manually edited or flags are copied to dependent columns, the term "manual" is added to the Q/C Flag Criteria field thereby "locking" the flags to prevent recalculation. If Q/C rules are subsequently changed, the rules will not be evaluated unless the flags are first unlocked. This is done by going to ***Edit > Q/C Flag Functions > Unlock Q/C Flags*** and then selecting either *All Columns*, *Data Columns Only*, or *Selected Columns Only*. This will remove the "manual" token and clear any manually-assigned or copied flags and trigger evaluation of the new rules.

6. Removing Flagged Data

Once data has been flagged, flagged values or records containing flagged values can be universally or selectively removed. Go to ***Edit > Q/C Flag Functions > Remove Data With Q/C Flags*** and then choose to ***Selectively Remove Values***, ***Null All Flagged Values***, or ***Delete All Rows with Flagged Records***. If you choose to selectively remove values, a new window will pop up that will allow you to select which flagged values will be removed.

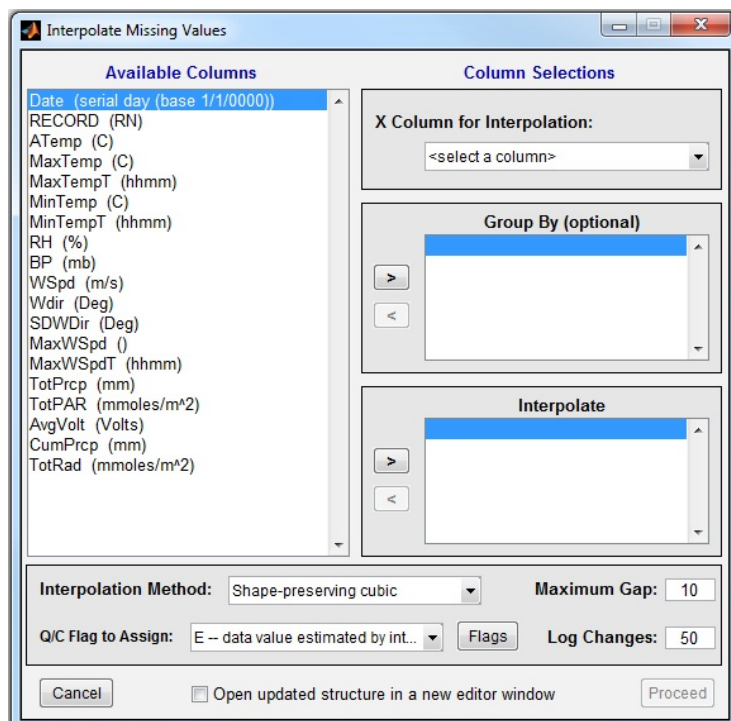


Selective Flagged Value Deletion window.

Gap Filling and Drift Correction

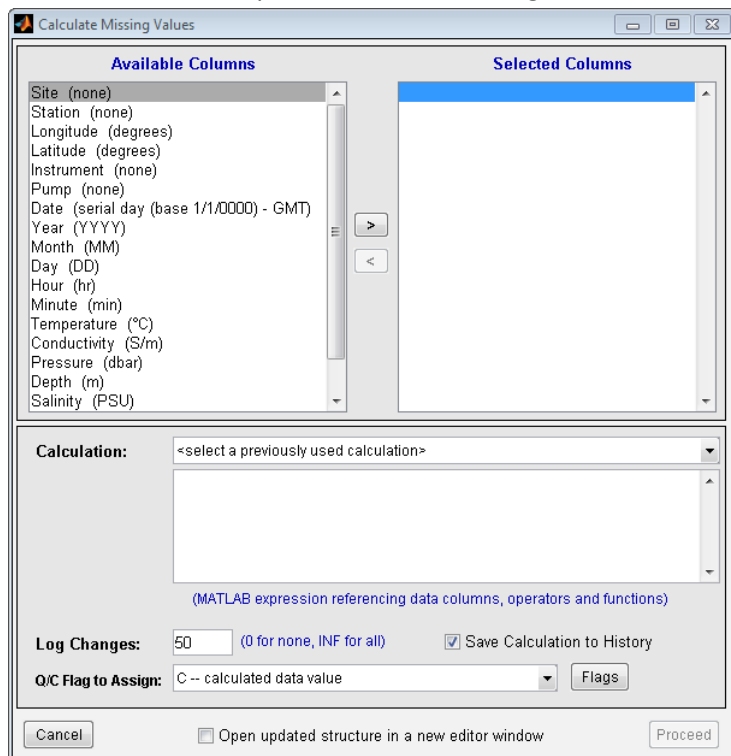
It is often necessary to fill gaps or correct for sensor drift in large time-series datasets, particularly to reduce bias when aggregating or re-sampling the data. While it is beyond the scope of this guide to recommend which specific gap filling or drift correction method to apply to a given dataset, the following exercises demonstrate the tools the GCE Toolbox provides to help with these activities.

1. Gaps can be filled using various interpolate methods by selecting ***Edit > Interpolate Missing Values*** in the Data Structure Editor. This will bring up the Interpolate Missing Values window. To use this dialog, select the X Column for interpolation from the drop down menu (a serial date column is auto-selected by default if present). Next, select columns that you wish to interpolate and add them to the "Interpolate" list using the corresponding ">" button. Optionally choose one or more columns to group by if the data set is a compound time series (e.g. multi-site data set). Select the Interpolation method, and set the Maximum Gap (maximum number of values in a gap that will be interpolated), followed by the flag that you wish to assign to the interpolated values. Hit the "Proceed" button to interpolate the values, which will then be automatically added to the data set.



The Interpolate Missing Values window

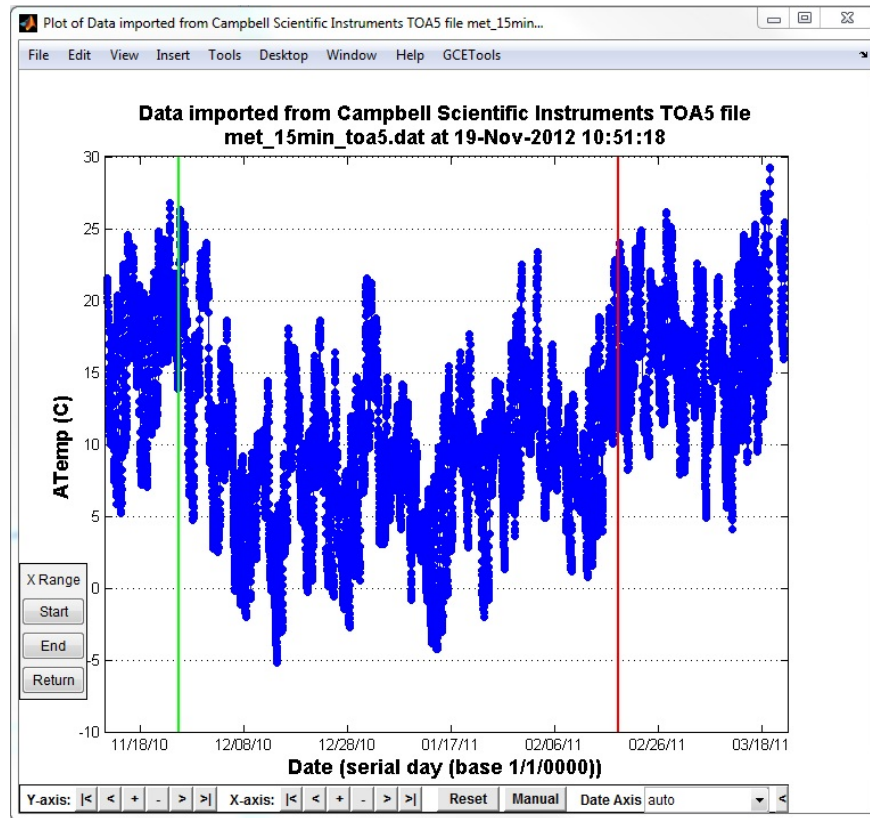
2. Gaps can also be filled using calculated values by going to "**Edit > Calculate Missing Values**". Select one or more columns to gap-fill, specify a MATLAB expression to evaluate (referencing functions, other data columns, or scalar values as appropriate), specify a flag to assign and click "**Proceed**" to evaluate the expression and fill missing values in the selected columns.



3. To correct data for sensor calibration drift, go to “**Edit > Correct for Sensor Drift**” in the Data Structure Editor. This will bring up the Correct Drift Window.

Drift Correction window.

- a. In this window, select the date column for the data set from the drop down menu. Next, select the column that you wish to correct and move it to the Columns to Correct window using the ">" button. Select the Correction Method you want to use from the drop-down menu, and then enter the offset value or weighted array of values to use for the correction as appropriate for the chosen method.
- b. Enter the date range for the values that you want to be corrected. This can be done by manually entering the date range, or by hitting the “**Pick**” button to the right of the Date Range fields. This will bring up a graph of the data for the selected column, and you can use the “**Start**” and “**End**” buttons on the left side of the graph to place markers to indicate the date range for the correction. Once the date range is set, use the return button to go back to the Correct Drift window.



Drift correction date range selection graph.

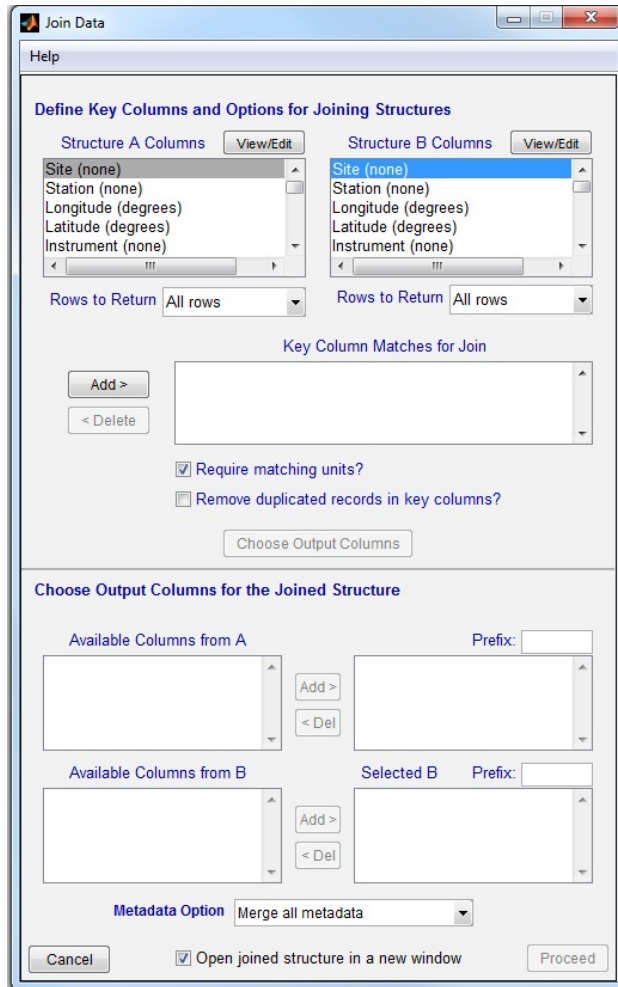
- c. Once the fields for the Correct Drift window have been filled in, click the **“Proceed”** button to apply the drift correction to the selected columns.

Creating and Exporting Data Products

Joining Data Sets

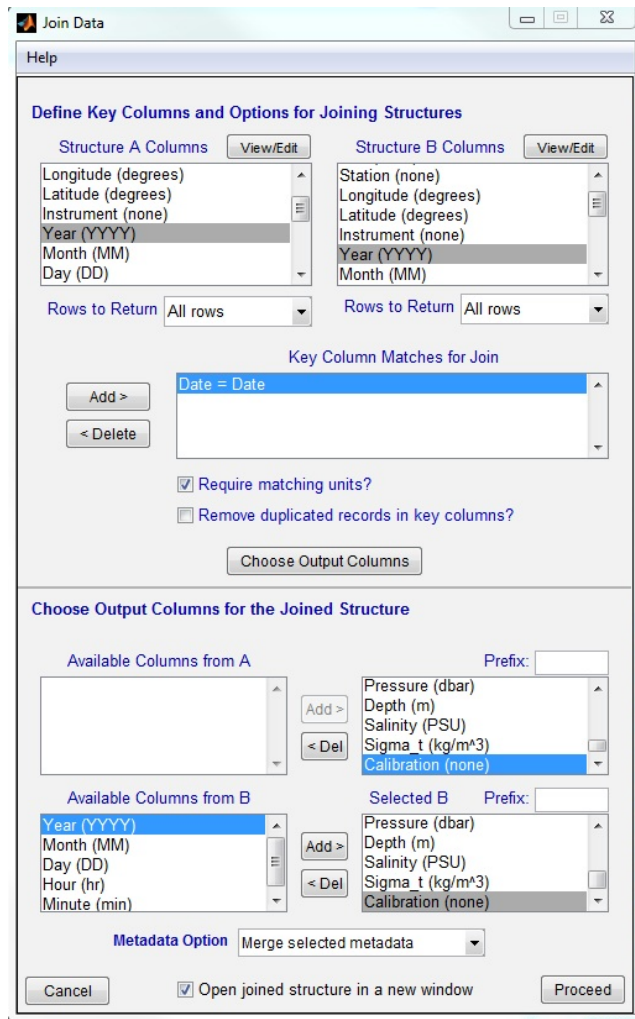
Selected columns from two different data sets can be joined together as a single file as long as they have a common unique key for each record.

1. This example will use the data from the Batch Importing exercise, which is located in the batch_products subdirectory. We will be using the 07110940.mat and 09110946.mat files from this directory. These files are from two sites with a similar data structure, and the data records cover the same date/time range.
2. Begin by loading the 07110940.mat file into the Data Structure Editor.
3. Next, go to **“File > Join Data Sets” > “Manual Key Selection” > “Data Structure File.”** Select the .mat file that will be joined to the existing data set, and then click the Open button. In this case, the file is the 09110946.mat file.
4. A new window will pop up that will allow you to choose which data column to use to match the two arrays together, and which columns you want to use in the joined dataset.



The Join Data Window.

5. For this data set, we use the “Date” column, which uses the serial date to match the records between the two data sets. Highlight “Date” in the Structure A and B Column boxes, then click the “**Add >**” button to add the Date column to the Key Column Matches for Join box.
6. Next, we need to select which columns will be joined from the two datasets by hitting the “**Choose Output Columns**” button. The “**A**” columns are from the first dataset, while the “**B**” columns are from the dataset that will be added.
 - a. In this case, we want to add all of the columns from the “Available Columns from A” box by highlighting each column, then pressing the “**Add >**” button.
 - b. Repeat this step for the columns in the “available Columns from B” box, but since we already are adding the Year, Month, Day, Hour, and Minute columns from the first data set, we don’t need to add them. Just add the columns that contain site information and sensor data.
 - c. Note that by leaving **Merge all Metadata** selected in the “**Metadata Option**” you will be able to provide complete metadata from both original data sets. You can elect to choose other options.
 - d. Hit the “**Proceed**” button.



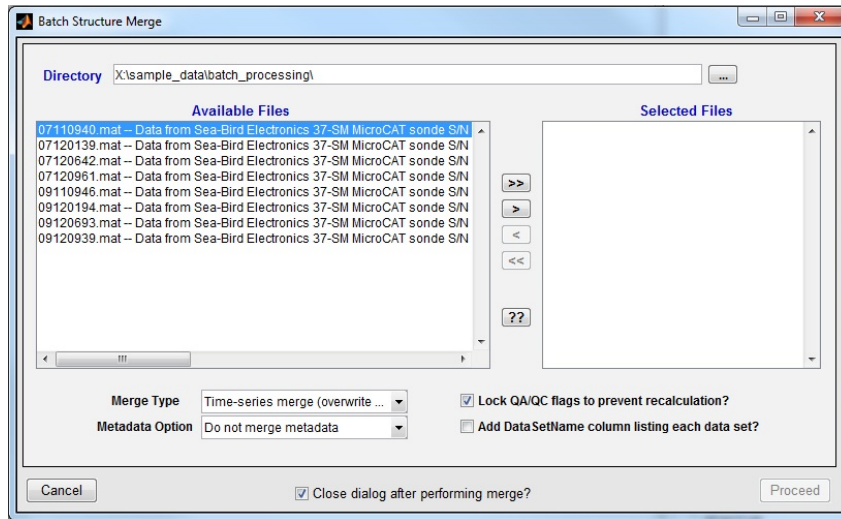
A Join Data Window that has been filled out.

7. The new merged array file will be created in the Data Structure Editor and now needs to be saved as a standard .mat file if you wish to use it again later.

Batch Merging Data Sets

Multiple related data structure files (e.g. repeated data logger downloads) can be quickly concatenated using the Data Merge Tool. In order to do this, all of the .mat files that are to be merged should have compatible data columns with the same names, units and data types. If any mismatched columns are present they will be offset and padded with missing values to create a rectangular combined data set.

1. This exercise will use the 071*.mat files from the batch_products subdirectory. These files should already be in the same directory, which is a requirement of using this function.
2. From the Data Structure Editor, go to **"Tools > GCE Data Merge Tool"** to bring up the Batch Structure Merge window.

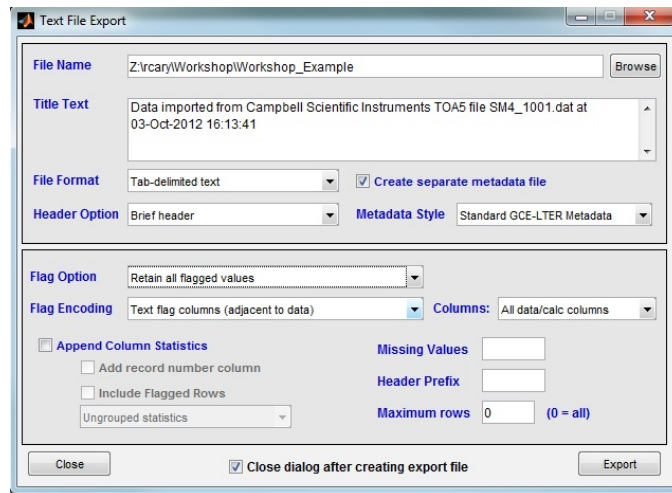


The Batch Structure Merge window.

3. In the Directory field, navigate to the directory containing the files you wish to merge, which is the batch_products directory in this case. The files will now appear in the Available Files window. Highlight the 071*.mat files that will be merged, then use the ">" button to move the files to the Selected Files window.
4. Select the merge type you wish to perform: "Append data sets in order" will append data based on the order in the Selected Files list. "Merge by study date" will append data in study-date order but keep all records even if redundancies exist. The "Time series merge" options will append the data sets based on study dates and automatically trim overlapping records to create valid time series data set, retaining all older records if "(add newer)" is specified and overwriting older records with newer observations if "(overwrite older)" is specified. For this example, we will use "Time-series merge (overwrite older)". Additionally, you can merge the metadata content from each of the data sets or just retain the metadata from the first structure. Metadata content from these data sets is the very similar, but we will choose **"Merge all metadata"** from the drop-down menu to mesh any differences and create composite documentation metadata for the integrated data set.
5. Hit the proceed button. A new data set containing all of the merged data will now be loaded into a Data Structure Editor window. Check the column list to make sure that no extra columns were added during the merge due to possible differences in attribute metadata.
6. Save the new file as a standard .mat file.

Exporting Data

1. Once the data has been processed as needed, the dataset and metadata file can be exported as a delimited text file for archiving or loading into another program. Start by going to **"File > Export Data > Text File (ASCII) > Standard Text File"**
2. In the new window that pops up, you will be able to make numerous choices about the format and content of the export file. The default format will create a tab-delimited file with brief headers, include flagged values, and create a separate file containing the metadata for the file.



Toolbox data export screen

3. Follow these steps to alter the export format:
 - a. Change the file name and location to the directory you wish to save the file to.
 - b. Change the title as you see fit.
 - c. Change the file format to comma-separated value.
4. When done, click the **“Export”** button to export the data.