The GCE Data Toolbox for MATLAB

Wade Sheldon

Georgia Coastal Ecosystems LTER
University of Georgia
email: sheldon@uga.edu







Background & Motivation

- Georgia Coastal Ecosystems LTER started in 2000
 - Large and varied data collection effort
 - NSF & LTER require data and metadata archiving, sharing
 - Needed to standardize data processing, quality control, documentation
- No ready-to-use software for LTER data management
 - Lots of great papers and reports, no tools to download
 - Most LTER sites were using "flat files" limiting
 - Relational databases emerging proprietary, complex, require constant network access
- Developed custom data management framework using MATLAB
 - Experienced using MATLAB for automating data processing, GUIs
 - Better code-reuse potential than database/web solution
 - Best compromise: file-based but supports fully dynamic operations

What is MATLAB?

From Mathworks (www.mathworks.com):

"MATLAB is a programming environment for algorithm development, data analysis, visualization, and numerical computation."

Benefits

- Ubiquitous in engineering and many science branches
- Rapid development, pre-built functionality, Java integration
- Code, GUIs and data formats multi-platform (Windows, Unix, Mac)
- Stable: good support and backward compatibility (~30 year history)
- Scalable (netbook to HPC) good performance with huge data sets
- Broad I/O support (serial ports to web services)

Costs

- Commercial ("licensed source") limits flexibility, costs \$-\$\$\$
- Some programming required for maximum use

Toolbox Development

Used ESA's "FLED" report to guide approach

Gross, Katherine L. and Catherine E. Pake. 1995. Final report of the Ecological Society of America Committee on the Future of Long-term Ecological Data (FLED). Volume I: Text of the Report. The Ecological Society of America, Washington, D.C.

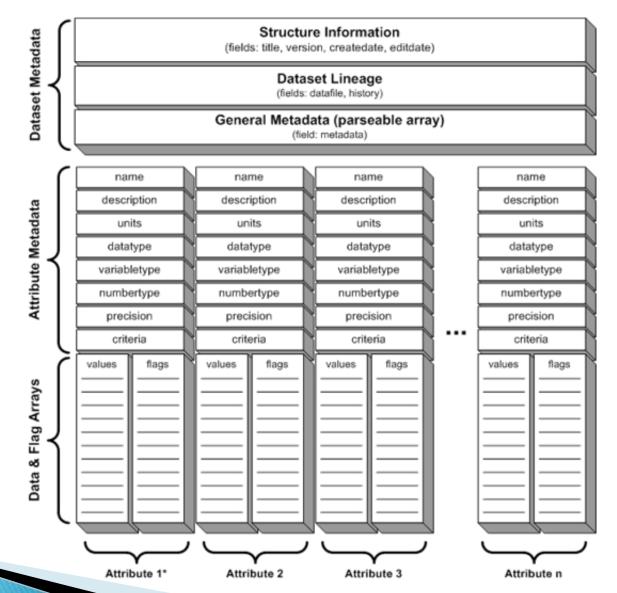
Identified information storage requirements

- Any number of numeric and text variables
- Attribute metadata for each variable (name, units, description, type, precision, ...)
- Structured documentation metadata for dynamic updating, formatting
- Versioning and processing history info (lineage)
- Quality control rules for every variable, qualifier flags for every value

Designed tabular data model in MATLAB

- MATLAB object ("struct") with named fields for each class of information
- Detailed specifications for allowed content in each field
- "Virtual table" design based on matched arrays for linking attribute metadata, data, flags
- Same philosophy as relational database table plus additional descriptors

Data Model (GCE Data Structure)



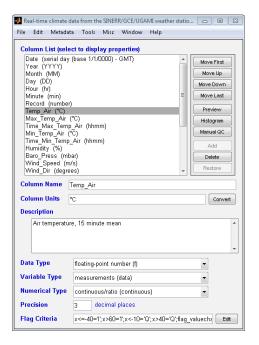
Toolbox Development

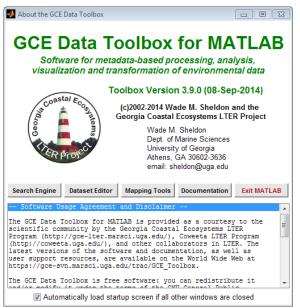
- MATLAB software library ("toolbox")
 - Utility functions to abstract low-level operations (API)
 - Create structure, add/delete columns, copy/insert/delete rows
 - Extract, sort, query, update data, update flags
 - Analytical functions for high-level operations
 - Statistics, visualizations, geographic & date/time transformations
 - Unit inter-conversions, aggregation/re-sampling, joining data sets
 - GUI interface functions to simplify using the toolbox
 - All functions use metadata, data introspection to autoparameterize and automate operations (semantic processing)
- Indexing and search support (and GUI search engine)
- Data harvest management tools

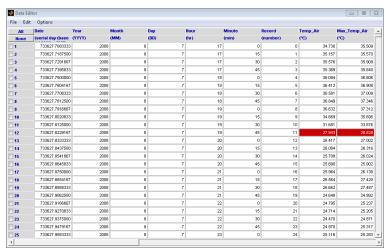
Command Line Interface (API)

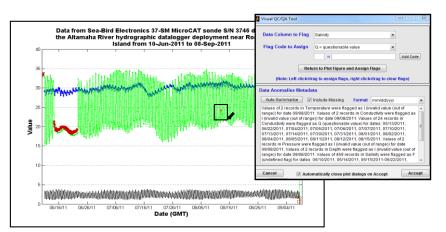
```
📣 МАТLAB 7.9.0 (R2009b)
File Edit Debug Desktop Window Help
🛅 👸 🤚 🖺 🤚 🥙 (*) 🚵 📸 😭 🖹 🕡 Current Folder: c:\userfiles\wade\svn_repositories\gce_toolbox
                                                                                                        ~ [...]
 Shortcuts 7 How to Add 7 What's New 📣 GCE Toolbox
   >> [s,msg] = fetch usgs('02226000', 'realtime',60, 'USGS Doctortown');
   >> s
   3 =
            version: 'GCE Data Structure 1.1 (29-Mar-2001)'
             title: 'Data from USGS Station 02226000 (ALTAMAHA RIVER AT DOCTORTOWN, GA) for 05-Feb-2010 through 06-Apr-2010
           metadata: {87x3 cell}
           datafile: {'usgs 02226000 realtime 20100406 1130 mod.txt' [5797]}
         createdate: '06-Apr-2010 11:30:48'
           editdate: '06-Apr-2010 11:30:50'
           history: {16x2 cell}
              name: {lx12 cel1}
              units: {'none' 'none' 'serial day (base 1/1/0000) - GMT' 'YYYY' 'MM' 'DD' 'hr' 'min' 'm' 'm^3/sec' 'mm'}
         description: {lx12 cel1}
           datatype: {'s' 's' 'd' 'f' 'd' 'd' 'd' 'd' 'f' 'f' 'f' }
        variabletype: {1x12 cel1}
         numbertype: {1x12 cel1}
          precision: [0 0 0 8 0 0 0 0 0 2 1 2]
            values: {1x12 cel1}
           criteria: {1x12 cel1}
              >> listcols(s)
    ans =
     1: Agency -- string
     2: StationID -- string
     3: Provisional -- integer
     4: Date (serial day (base 1/1/0000) - GMT) -- floating-point
     5: Year (YYYY) -- integer
     6: Month (MM) -- integer
     7: Day (DD) -- integer
     8: Hour (hr) -- integer
     9: Minute (min) -- integer
    10: GageHeight (m) -- floating-point
    11: Discharge (m^3/sec) -- floating-point
    12: Precipitation (mm) -- floating-point
    >> dt = extract(s,'Date'); discharge = extract(s,'Discharge');
    >> whos
     Name
                                       Bytes Class
                                                       Attributes
                     12x63
                                       1512 char
     discharge
                    5797x1
                                       46376 double
                    5797x1
                                       46376 double
                    0x0
                                          0 char
                                     1346932 struct
     3
  fx >> |
📣 <u>S</u>tart
```

GUI Application Interface

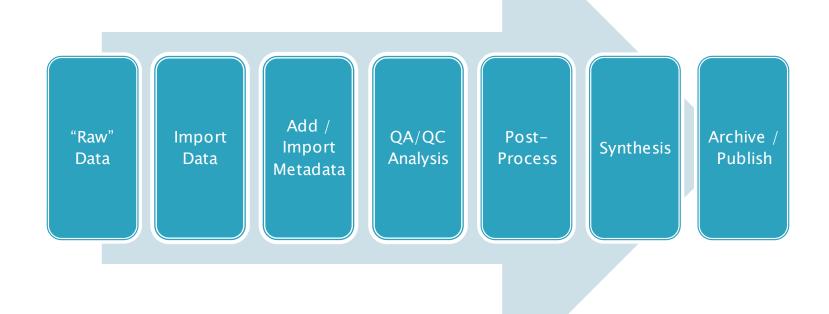




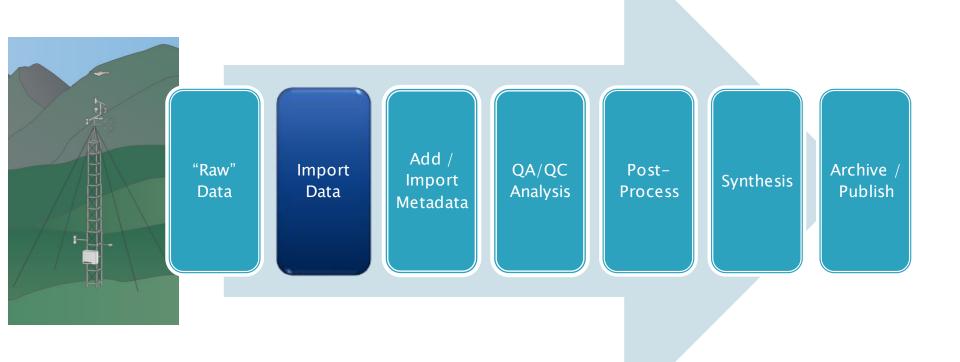




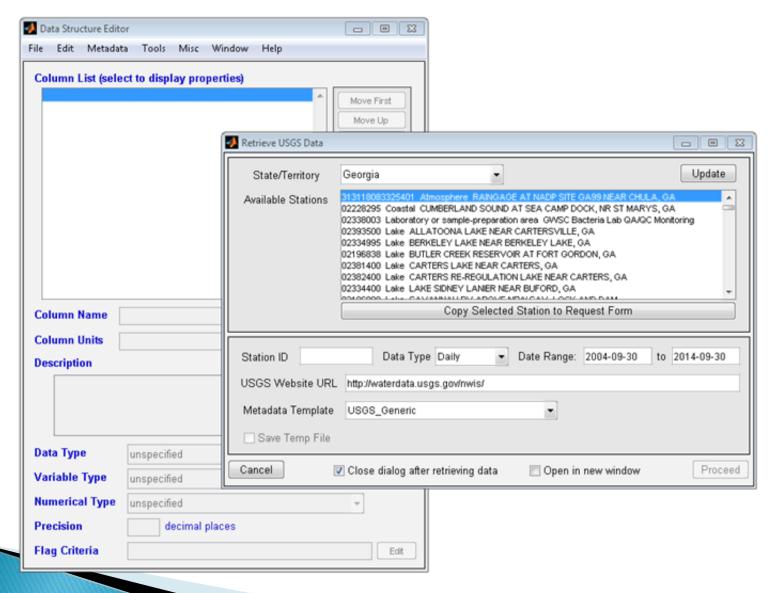
Data Management Cycle

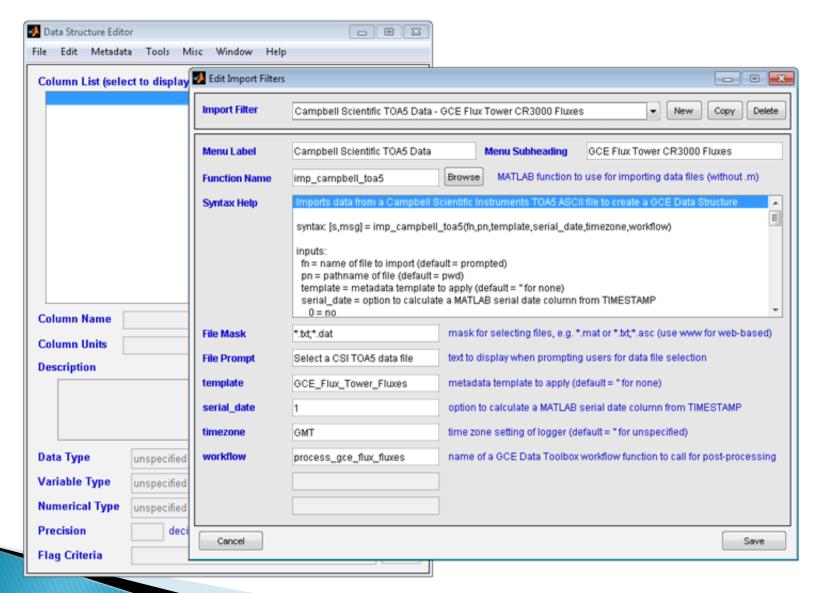


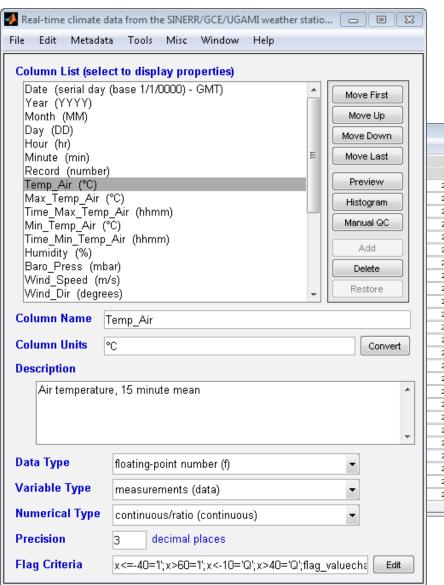
Data Management Cycle

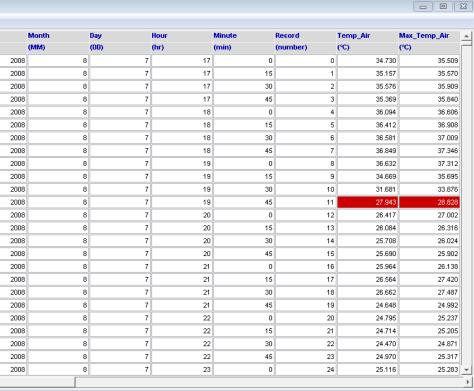


- Generic parsers
 - Delimited text (CSV, space, tab)
 - MATLAB variables (arrays, matrices, structs)
- Specialized parsers
 - Vendor–specific logger files
 - Campbell Scientific Instruments (tables, arrays)
 - Sea-Bird CTD, sondes
 - Others (Hobo, Schlumberger, OSIL, ...)
 - Network data sources
 - SQL database queries (JDBC)
 - Federal databases (USGS NWIS, NOAA NCDC, NOAA HADS)
 - LTER ClimDB/HydroDB
 - EML repositories (LTER NIS/PASTA, KNB)
 - Data Turbine servers
- Custom parsers can be added

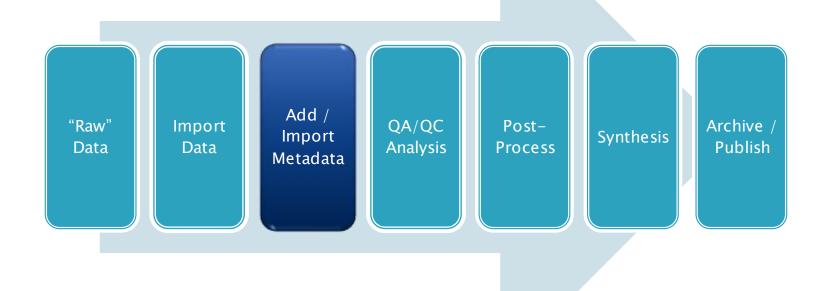








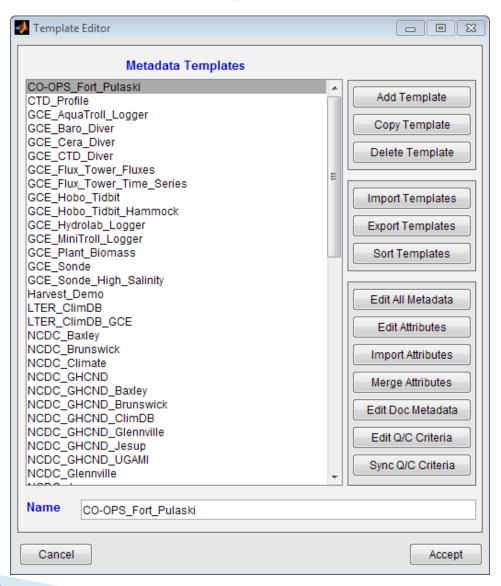
Data Management Cycle



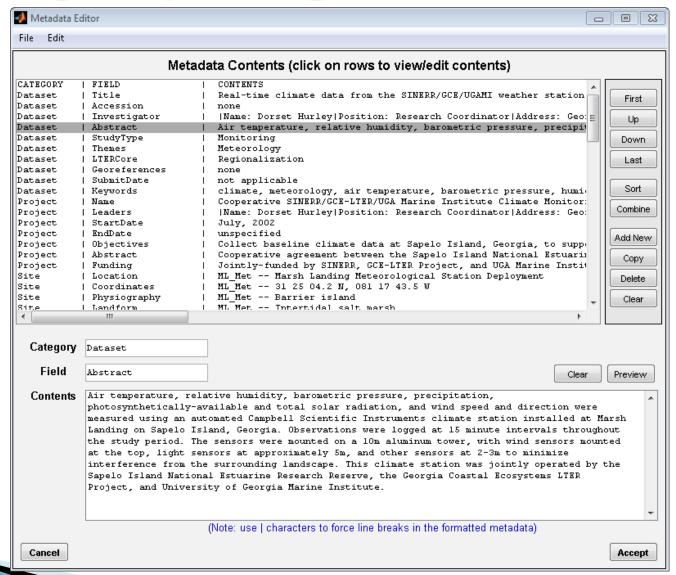
Adding/Importing Metadata

- Metadata entry is time-consuming (and tedious)!
 - Metadata capture, re-use emphasized
 - Metadata auto-generated whenever possible
- Many pathways for building Metadata
 - Metadata can be imported along with data
 - Logger file headers (Campbell, Sea-bird)
 - Station, parameter information from USGS, NOAA
 - Tokenized headers from Data Submission Template
 - Metadata imported from other data structures
 - Metadata imported from data repository (EML/XML)
 - Metadata added from user-defined "Templates"
 - Column (attribute) metadata matched to "variables"
 - Boilerplate documentation
 - GUI tool for creating/managing templates

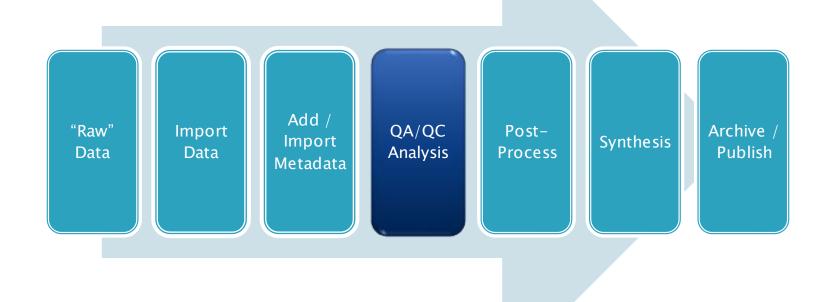
Adding/Importing Metadata



Adding/Importing Metadata



Data Management Cycle



QA/QC Analysis Framework

Programmatic Q/C Analysis

- "Rules" (i.e. criteria) define conditions in which values should be flagged
- Unlimited Q/C rules for each variable
- Rules evaluated when data loaded and when data or rules change
- Rules predefined in metadata templates to automate Q/C on import

Interactive Q/C Analysis and Revision

- Qualifiers can be assigned/cleared visually on data plots with the mouse
- Qualifiers can be propagated to dependent columns
- Qualifiers can be removed or edited (search/replace) if standards change

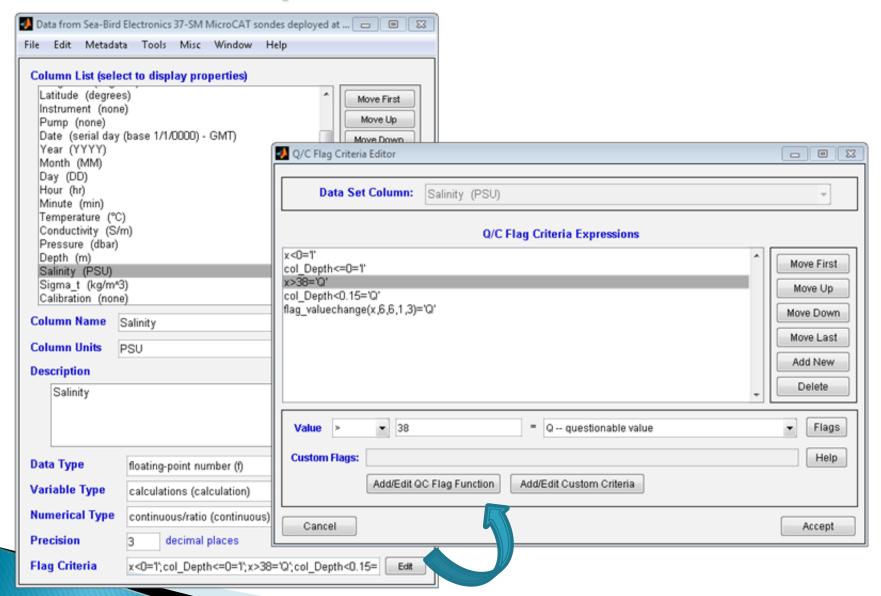
Automatic Documentation of Q/C Steps

- Q/C operations (including revisions) logged to processing lineage
- Data anomalies reports can be auto-generated and annotated

Data analysis, synthesis tools Q/C-aware

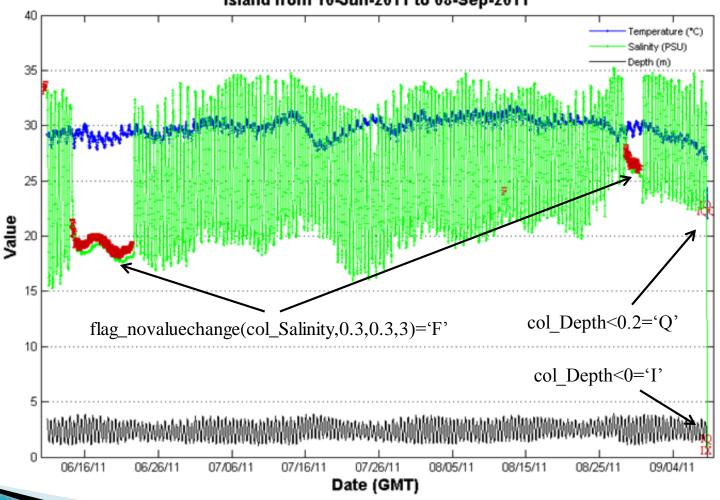
- Qualified values can be filtered, summarized, visualized during analysis
- Statistics about missing/qualified values tabulated, used to qualify derived data

QA/QC Analysis Framework

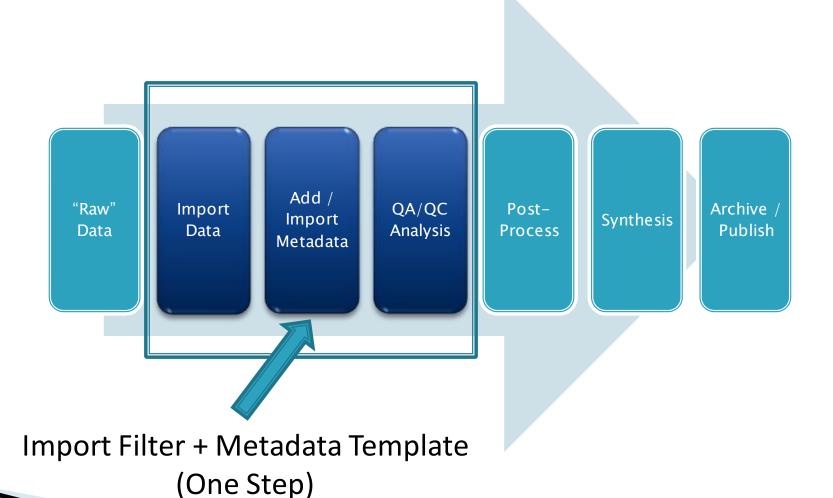


Example - CTD mooring

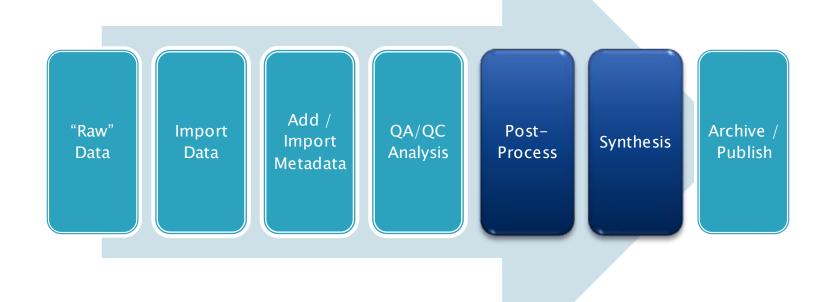
Data from Sea-Bird Electronics 37-SM MicroCAT sonde S/N 3746 deployed at the Altamaha River hydrographic datalogger deployment near Rockdedundy Island from 10-Jun-2011 to 08-Sep-2011



Data Management Cycle

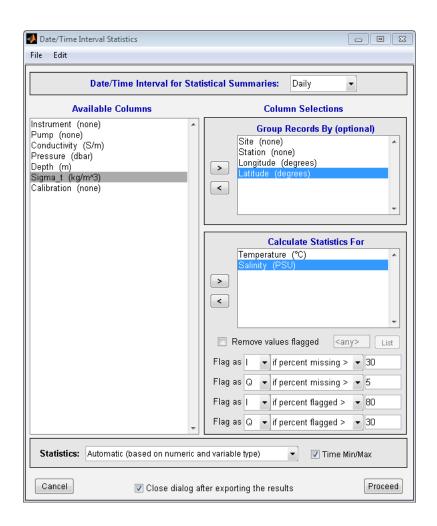


Data Management Cycle

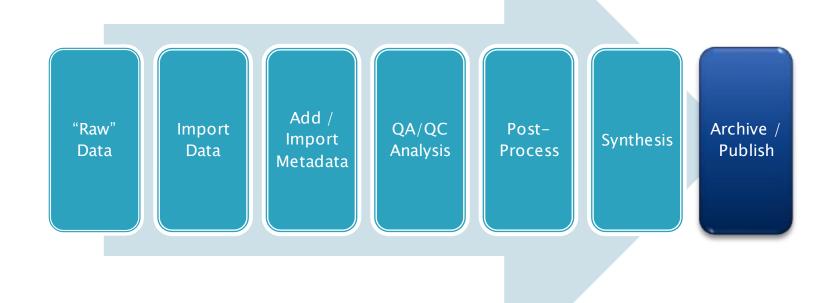


Post-Processing and Synthesis

- Derived data sets can be created by filtering values or refactoring data table structure (e.g. combining or splitting columns)
- Data can be re-sampled or summarized by aggregation, binning and date/time scaling
- Multiple data sets can be combined by merging (union) and joining on key columns
- All derived data contain complete metadata describing the entire processing history
- QA/QC rules can be generated for derived data columns automatically based on number or percent missing/flagged values

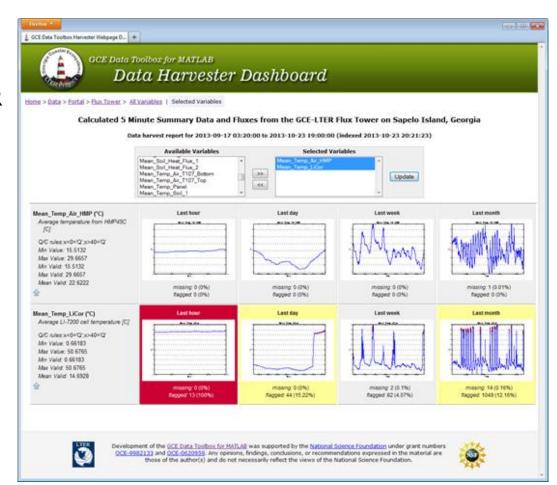


Data Management Cycle

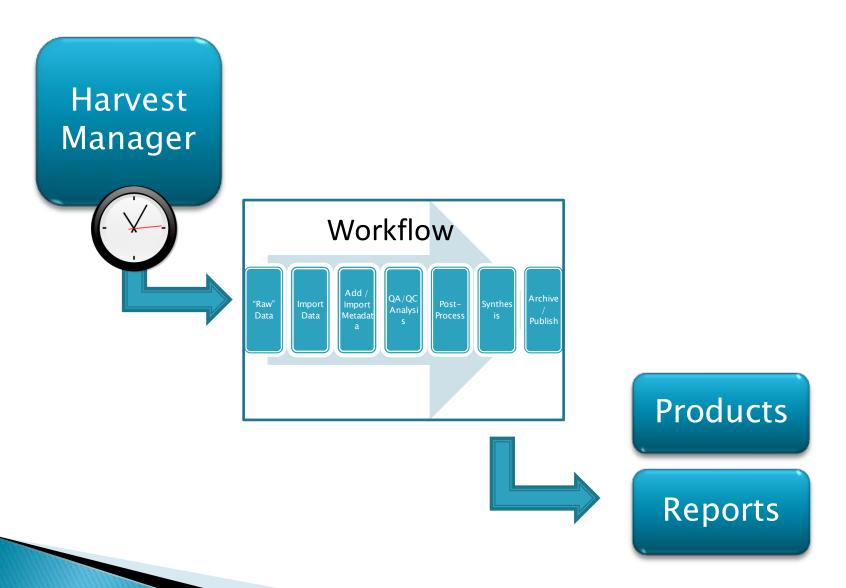


Archive/Publish Data

- Data and EML metadata can be published in scientific data archives, including LTER NIS, DataONE Member Nodes (e.g. KNB), Dryad, ...
- Data can be refactored and published in CUAHSI ODM database, Hydroserver
- Harvested data can be displayed on MATLABgenerated web pages and dashboards



Data Harvesting Workflows



Key Concepts

Operations performed in context of a "dataset"

- Passing data columns to a tool transports metadata as well
- Dataset metadata used to guide transformation, plotting, analysis
- Metadata used to auto-parameterize functions

Workflow steps generate new datasets

- Each step along a workflow results in a complete data set with metadata
- Intermediate datasets can be saved or overwritten in workflows

Processing history ("lineage") tracked

- Each tool logs operations by date/time
- Data revisions, deletions, flagging captured at user-specified detail
- Lineage reported in metadata

Metadata are "live", and updated automatically

- Data column metadata and data revisions
- Calculations and unit conversions
- Code definitions
- Metadata merged when multiple data sets are joined

Implementation Scenarios

End-to-End Processing (logger-to-scientist)

- Acquire raw data from logger, file system, network
- Assign metadata from template or forms to validate and flag data
- Review data and fine-tune flag assignments
- Generate distribution files & plots, archive data, index for searching
- Scientists can use toolbox on desktop to analyze, integrate data

Data Processing and Q/C

- Acquire, validate and flag raw data (on demand or timed/triggered)
- Upload processed data files (e.g. csv) or values & flags to RDBMS (e.g. ODM)

Workflow Step

- Call toolbox from other software as part of workflow (tool-chaining)
- Use toolbox as middleware between other systems (e.g. Data Turbine & ODM, Kepler, Taverna)

Real-Time Sensor Data

Good Scalability

- Data only limited by computer memory (tested >1 GB data sets)
- Multiple instances can be run on high-end workstations
- Good flag evaluation performance in use, testing

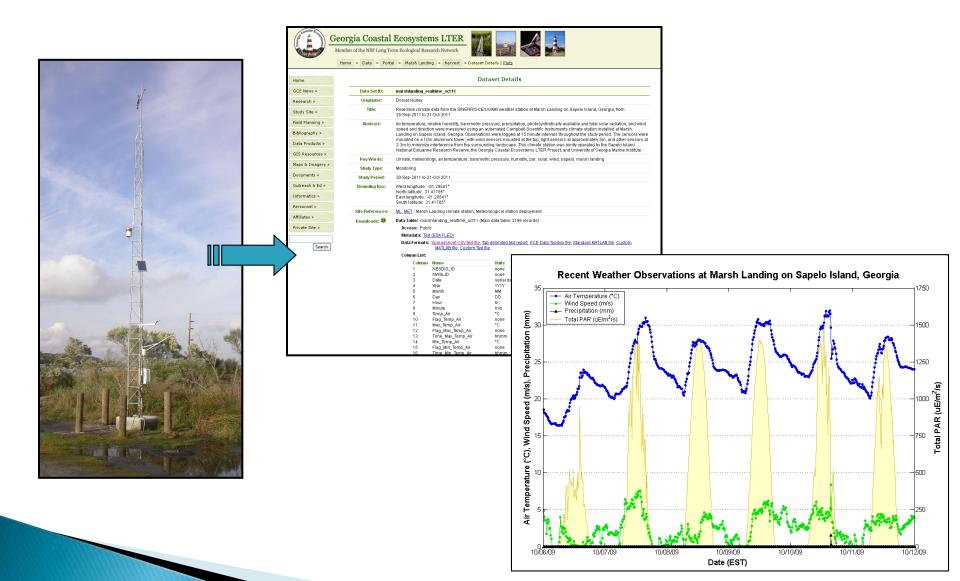
Good scope for automation

- Command-line API for unattended batch processing
- Timed and triggered workflow implementations easy to deploy

Support for multiple formats, protocols

- Formats: ASCII, MATLAB, SQL, specialized (CSI, SBE, NWIS, HADS, ...)
- Transport: local file system, UNC paths, HTTP, FTP, SOAP
- Already used at 6 LTER sites (5 sites evaluating)
- USGS data harvester for HydroDB since 2003

GCE Data Harvesting (>12 years)



Ongoing Development

Metadata model enhancements

- Add "schema" support for better EML alignment
- Improve metadata editor with dynamic form fields

Improve harvest management

- Add GUI tools for configuring harvest info, plot info
- Add GUI for managing timers
- Add data harvesting "dashboard" for monitoring activity
- Workflow editor
- "Rebranding" as a community tool

Resources

MATLAB

- Website: http://www.mathworks.com/products/matlab/
- Version R13 (2002) or higher required (full or student version)

Software Distribution

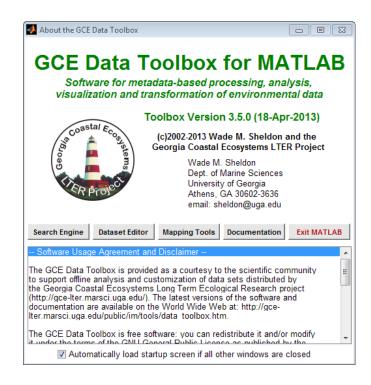
- Website: https://gce-svn.marsci.uga.edu/trac/GCE_Toolbox
 - Documentation, FAQ
 - Distribution Downloads (stable, beta)
 - Bug reporting (tickets)
- SVN: https://gce-svn.marsci.uga.edu/svn/GCE_Toolbox/trunk
- Code is open source (GPLv3) and cross-platform

User Support

- Peer-to-peer model
- Sporadic training opportunities (LTER)
- Email: gcetoolbox-l@listserv.uga.edu (http://www.listserv.uga.edu/)

Interactive Tutorial

- Introduction to the Dataset Editor application
- Importing and exploring data
 - Generic ASCII
 - Native logger files (Campbell, SeaBird)
 - Data services (Data Turbine, USGS NWIS)
- Metadata management
 - Defining attribute metadata
 - Documentation metadata, templates
- QA/QC framework
 - Defining flags and "rules"
 - Visual QA/QC
 - Managing flagged data
- Creating and exporting products
 - Batch processing raw data
 - Integrating data (join, merge)
 - Summarizing data
 - Exporting data sets and plots



Questions?

