# Automating Sensor Data Management using the GCE Data Toolbox for MATLAB

#### Wade Sheldon

Georgia Coastal Ecosystems LTER
Department of Marine Sciences
University of Georgia
Email: sheldon@uga.edu







### Introduction

- Continual improvements in sensor technology allow scientists to observe the environment at unprecedented spatial and temporal scales
- Advances in IT infrastructure allow increasing numbers of sensors to be deployed, managed and accessed remotely
- Volume of sensor data that many researchers manage is increasing rapidly
- Data collection no longer limiting factor now it's data validation, quality control and documentation
- Funding agencies and journals now require data and metadata archiving, sharing
- Few generalized, pre-built software tools are available to meet these data needs



### Sensor Data Challenges

- Complex landscape of sensor systems, software tools and repositories
- Numerous data formats, syntax conventions and standards
- Disruptive to research groups without trained informatics experts on staff





























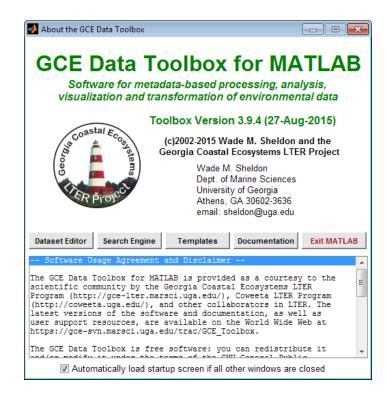






### GCE Data Toolbox

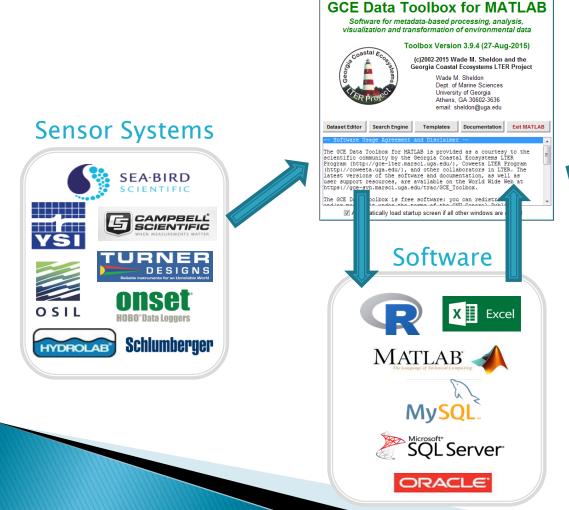
- Software developed at the Georgia Coastal Ecosystems LTER (GCE Data Toolbox) proven effective for sensor data management
- Open source MATLAB code library (toolbox)
  - Integrates data processing, quality control and analysis with metadata capture, creation and management to support the entire data lifecycle
  - GUI forms and command line API interfaces
  - Can be used as a complete stand-alone solution for sensor data and metadata management
  - Can be used in conjunction with other data management tools ("tool chaining")
- Data can be retrieved from USGS, NOAA,
   DataONE and other archives for analysis and comparison
- Generates structured data files, metadata to meet new archiving mandates from funding agencies and publishers



### GCE Data Toolbox

Effectively integrates sensor data resources and data repositories

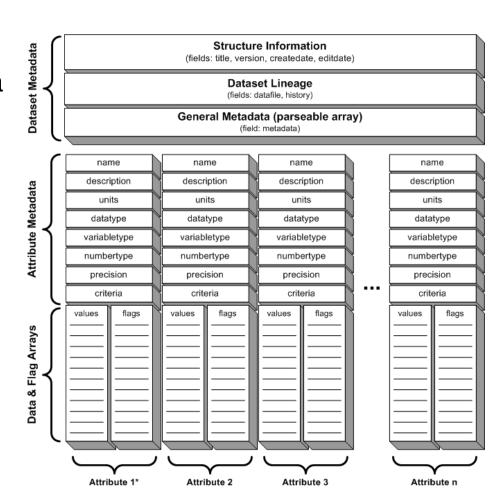
About the GCE Data Toolbox



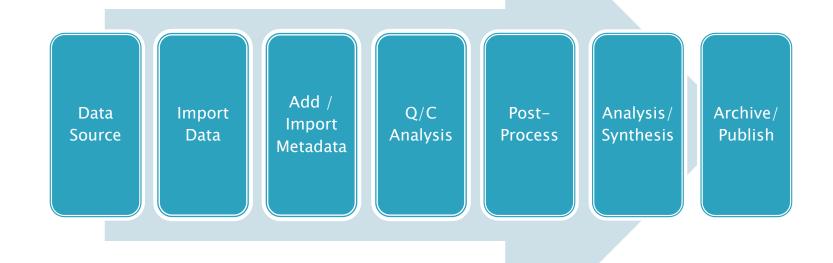


### **Data Model**

- Toolbox data model stores all contextual info with data
  - Documentation
  - Variable descriptions
  - Data values
  - Quality control "rules"
  - Qualifier flags
  - Processing history
- Supports "semantic processing" by tools
- Maintains context/validity throughout analysis



# Data Management Cycle

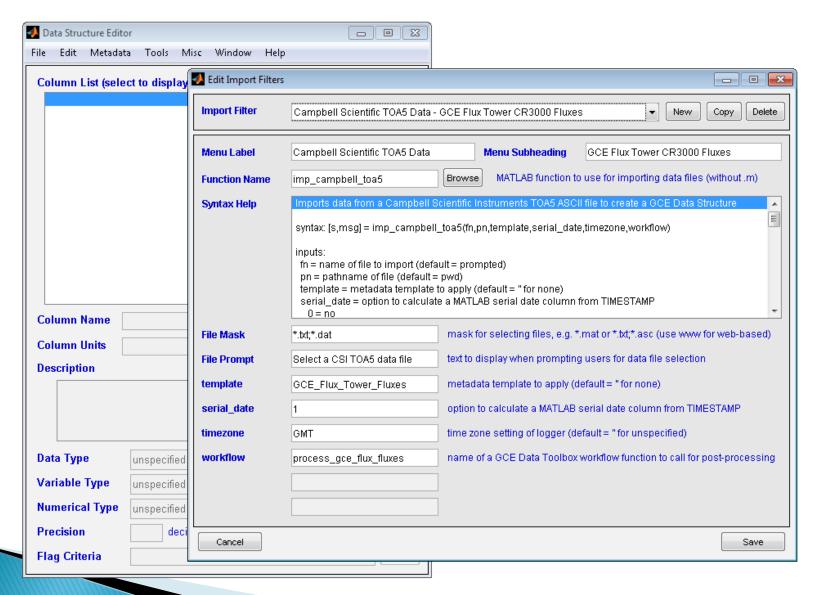


### Importing Data

- Generic parsers
  - Delimited text (CSV, comma, 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, AquaTroll, OSIL, ...)
  - Network data sources
    - SQL database queries (JDBC, ODBC)
    - Federal databases (USGS NWIS, NOAA NCEI/GHCN, NOAA HADS)
    - LTER ClimDB/HydroDB
    - Ecological data repositories (LTER NIS/PASTA, KNB, DataONE)
    - Data Turbine servers
- Custom parsers can be added



### Example - Adding a Parser



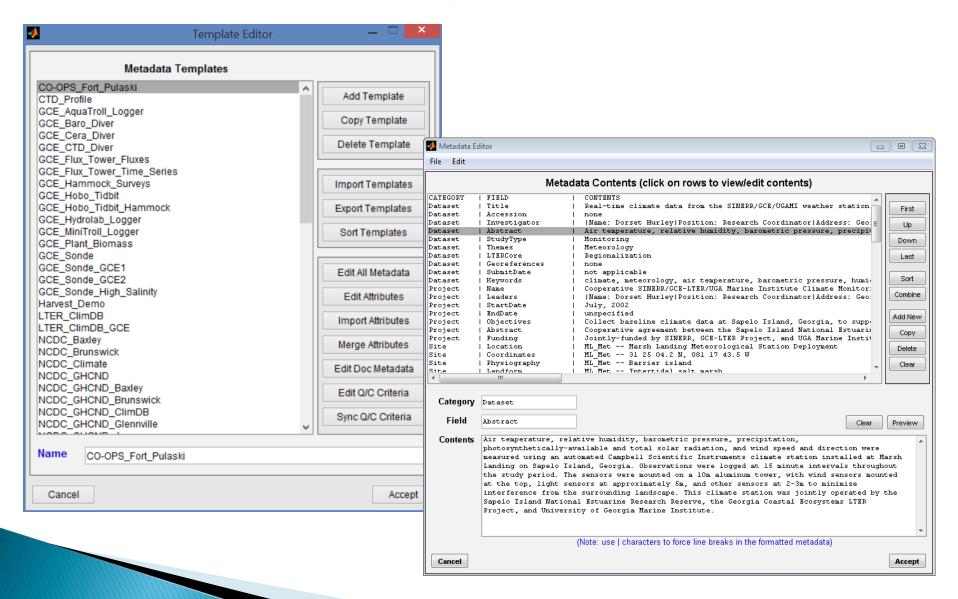
### Adding/Importing Metadata

Data documentation (metadata creation) 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 web sites
    - Tokenized headers from Data Submission Template
  - Metadata imported from other GCE Data Toolbox data structures
  - Metadata imported from the source data repository
  - Metadata added from reusable "Templates"
    - Detailed variable (attribute) metadata matched to raw data fields
    - Boilerplate documentation
    - GUI tool for creating/managing templates



# Adding/Importing Metadata



### Q/C 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

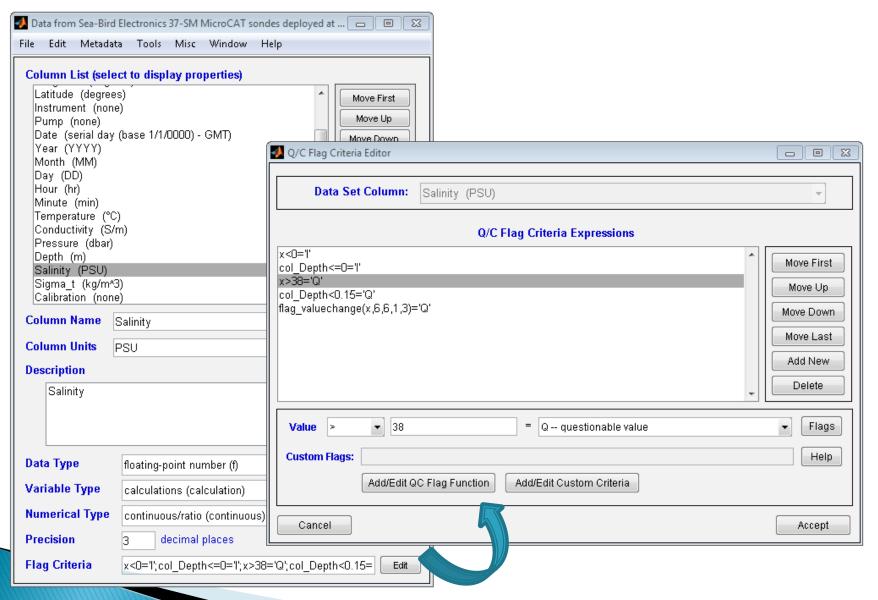
#### Automatic Documentation of Q/C Steps

- 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

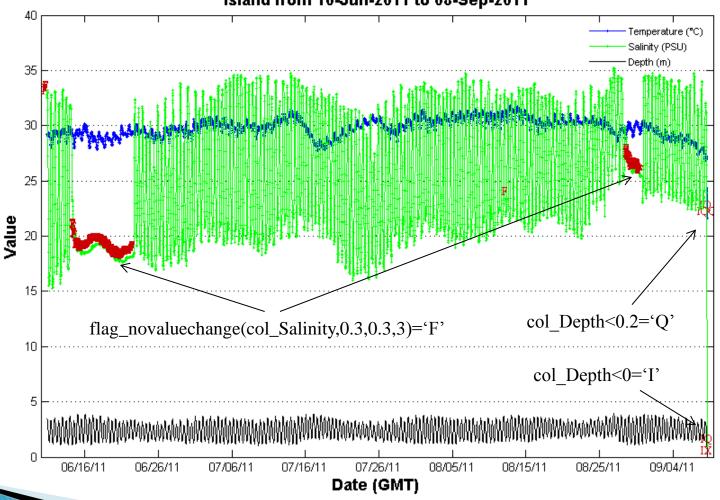
- Qualified values can be filtered, summarized, visualized during analysis
- Statistics about missing/qualified values tabulated, used to qualify derived data

### QA/QC Analysis Framework



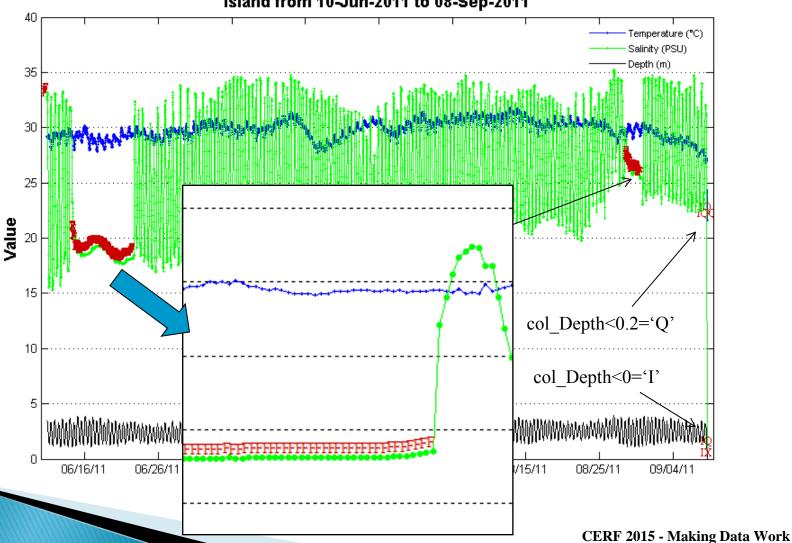
# Example - CTD mooring

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

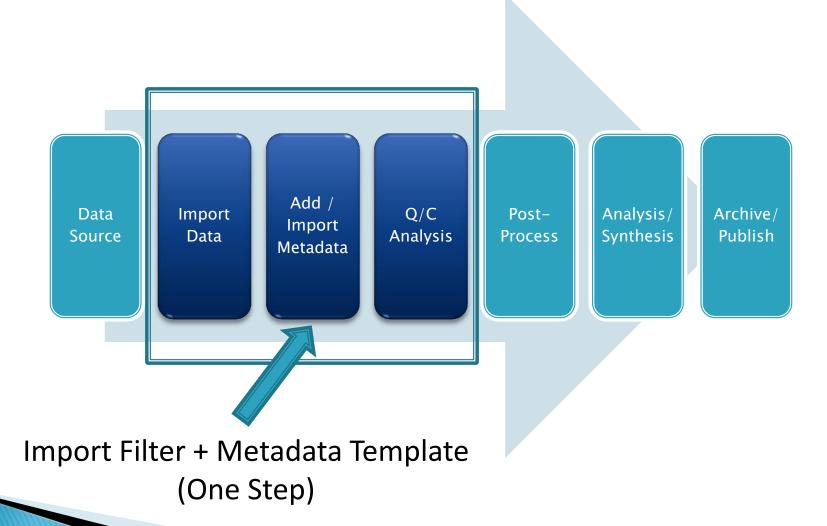


### Example - CTD mooring

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

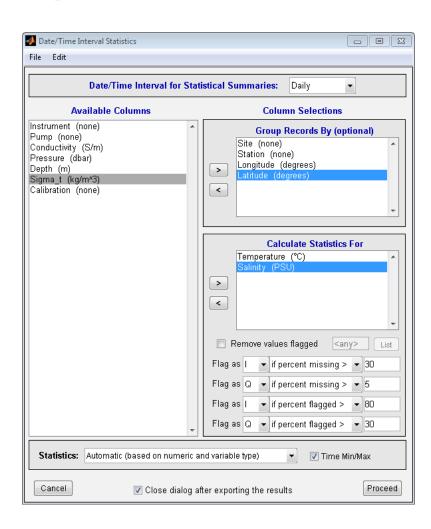


# Data Management Cycle

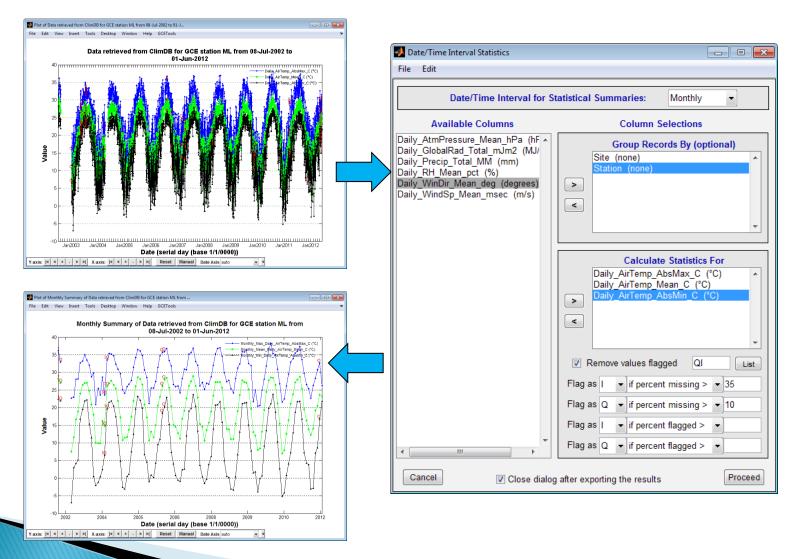


# Post-Processing and Synthesis

- Calculated columns can be generated using mathematical formulas, functions
- Data can be gap-filled, drift-corrected
- 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
- Q/C rules can be generated for derived data columns automatically based on number or percent missing/flagged values

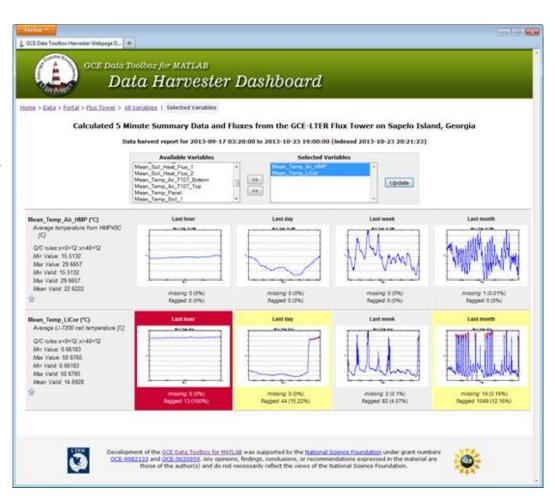


# Example - Date/time Aggregation



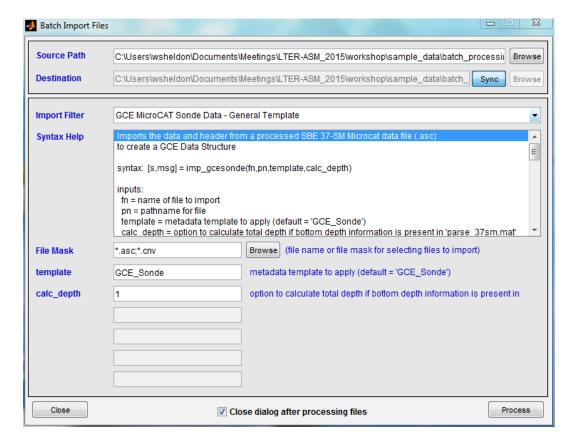
### Archive/Publish Data

- Data and metadata can be exported in many formats
  - Text, MATLAB, XML/KML/HTML
  - RDBMS tables
- EML-described data packages can be generated/published in data archives
  - LTER NIS/PASTA
  - DataONE Member Nodes (e.g. KNB)
- Data can be refactored and published in CUAHSI ODM database, Hydroserver
- Data can be displayed on MATLAB-generated web pages and dashboards



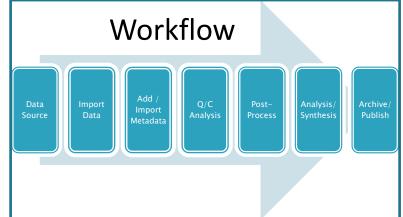
### Automation - Batch Processing

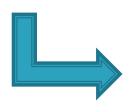
- Entire directories of raw data files can be processed at once using an import filter and metadata template
- Multiple GCE Data Structure files in a directory can be batch-converted to TXT, CSV
- Multiple files can be merged via metadata-based union to create an integrated data set
- All directories containing GCE Data Structures can be indexed and searched, then merged, joined and exported



# Automation - Data Harvesting



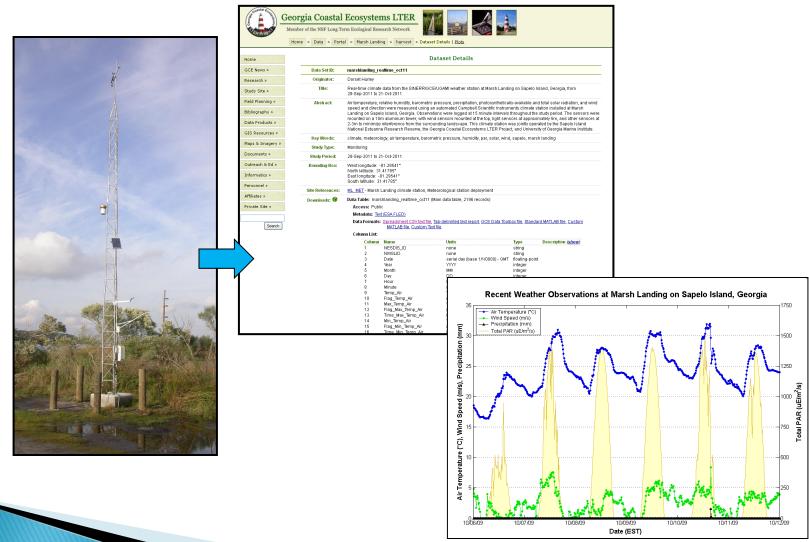




Products

Reports

# Example - Real-time Harvesting



# **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, complete 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/publish data
- 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. DEIMS, ODM Database)

#### 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)

#### 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, Tutorial, 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/)