#### From Scientist and Sensor to Synthesis: Overview of the GCE Data Toolbox for MATLAB

# Wade Sheldon Georgia Coastal Ecosystems LTER John Chamblee & Richard Cary Coweeta LTER

## Background & Motivation

- Georgia Coastal Ecosystems LTER project started in Sept 2000
  - Large data collection effort (cruises, moorings, met stations, water quality, field surveys, ...)
  - NSF & LTER require data archiving and sharing
  - > LTER requires detailed "metadata" for every data set
  - 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
  - ➤ A few sites using relational databases, client/server apps proprietary, complex, unfamiliar, require constant network access
- Chose to develop custom data management software (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: (http://www.mathworks.com/products/matlab/)

"MATLAB is a programming environment for algorithm development, data analysis, visualization, and numerical computation. Using MATLAB, you can solve technical computing problems faster than with traditional programming languages, such as C, C++, and Fortran."

#### Benefits:

- Ubiquitous in engineering and many science branches
- Rapid development with lots of pre-built functionality, Java integration
- Cross-platform code, GUIs and data formats (Windows, \*nix, Mac OS/x)
- Stable: good support and backward compatibility (~30 year history)
- Scalable (netbook to cluster) great performance with huge data sets
- Broad I/O support (serial ports to web services)

#### Drawbacks:

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



## Toolbox Development

#### Started by reviewing ESA's "FLED" report

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

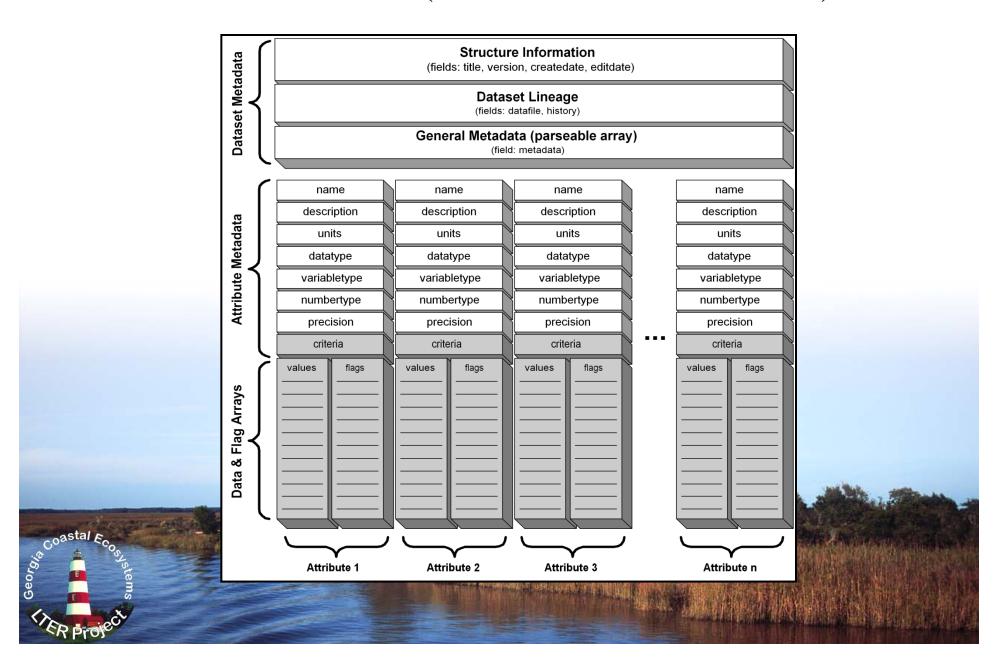
- Any number of numeric (integer, float, exponential) and text variables
- Structured attribute metadata for each variable (name, units, desc., type, precision, ...)
- Structured documentation (dataset metadata) for dynamic updating, formatting
- Versioning and processing history info (lineage)
- Quality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

#### Designed data model: "GCE Data Structure"

- MATLAB "struct" array 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

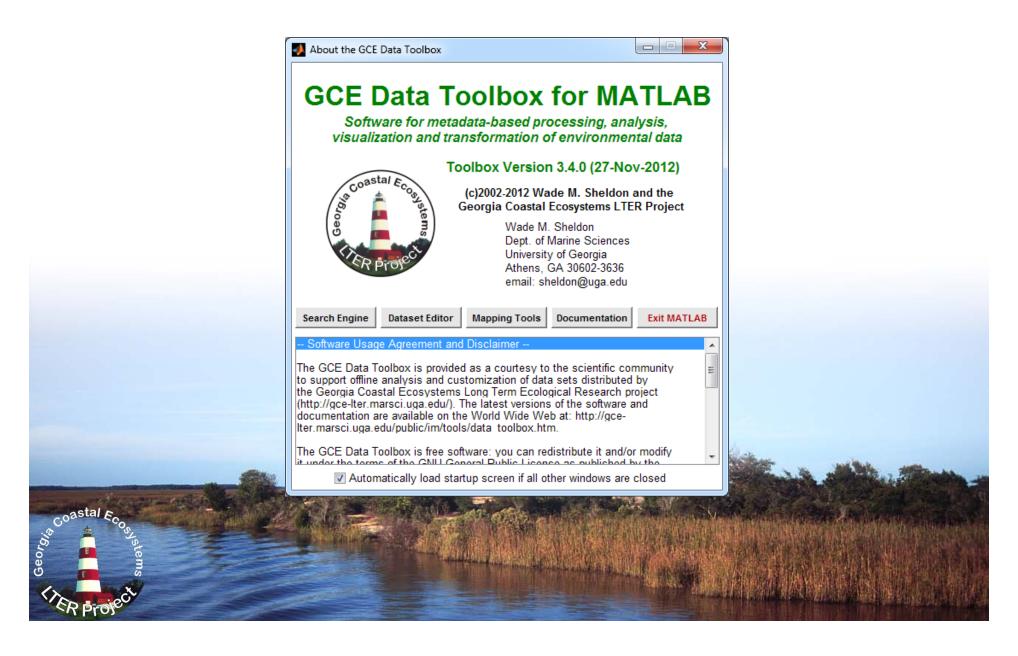
- Developed MATLAB software library to work with data structures
  - Utility functions to abstract low-level operations (API)
    - Create structure, add/delete columns, copy/insert/delete rows
    - Extract, sort, query, update data, update flags
  - Analytical functions for high-level operations
    - 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 auto-parameterize and automate operations (semantic processing)
- Developed indexing and search support (and GUI search engine)



#### **Command Line**

```
🥠 МАТLAB 7.9.0 (R2009b)
File Edit Debug Desktop Window Help
                                                                                                         ~ [...]
🌇 🚰 👗 🖣 📬 🤊 🥲 🧥 📸 😭 📦 Current Folder: c:\userfiles\wade\svn_repositories\gce_toolbox
 Shortcuts 🖪 How to Add 🔃 What's New 📣 GCE Toolbox
   >> [s,msg] = fetch_usgs('02226000','realtime',60,'USGS_Doctortown');
   s =
            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: {1x12 cel1}
          datatype: {'s' 's' 'd' 'f' 'd' 'd' 'd' 'd' 'f' 'f' 'f'}
       variabletype: {lx12 cell}
         numbertype: {1x12 cel1}
          precision: [0 0 0 8 0 0 0 0 0 2 1 2]
             values: {lx12 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
                                       Bytes Class
     Name
                      Size
                                                       Attributes
                     12x63
                                        1512 char
      ans
      discharge
                    5797x1
                                        46376 double
                    5797x1
                                       46376 double
                    0x0
                                         0 char
      msg
                      1x1
                                     1346932 struct
  fx >> |
```

# Startup Dialog



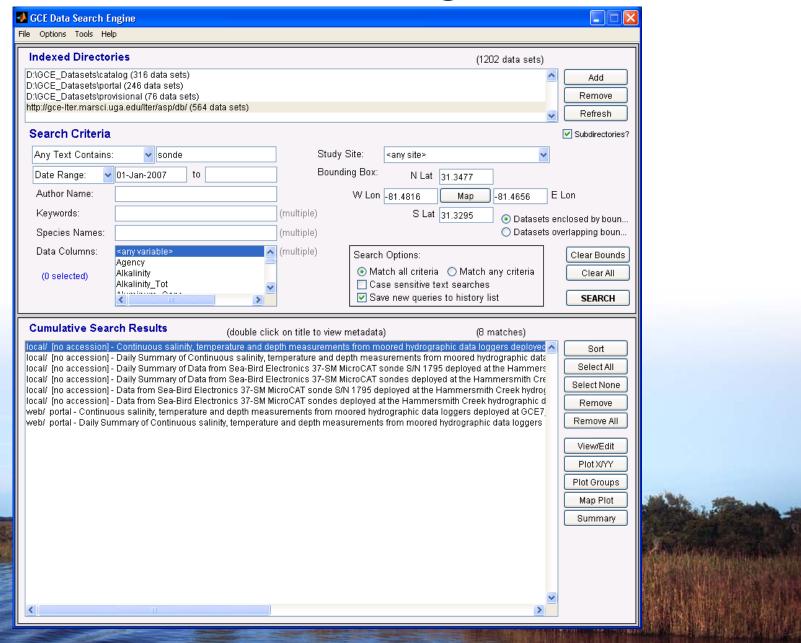
## Dataset Editor

		inity, temperature and depth n Tools Misc Window Help	neasuremen	. 🔲 🗆 🔀	
	Column List (se	lect to display properties)			
	Site (none) Longitude (degreen Latitude (degreen Latitude (none) Pump (none)	grees) es) ne) y (base 1/1/0000) - GMT) °C) S/m) )	Mo' Mo' Mo' Mo' Mi	ove First ove Up ove Down ove Last review stogram onual QC Add Delete	
	Column Name	Site			
	Column Units	none		Convert	
	Description				
	Nearest nomin	al GCE-LTER sampling site		^ _	
	Data Type	integer (d)		~	Hehola .
	Variable Type	categorical values (nominal)		~	
	Numerical Typ	e discrete/interval (discrete)		~	
Coastal Eco	Precision	O decimal places			Managaria de la companya de la comp
Coastal Eco	Flag Criteria			Edit	

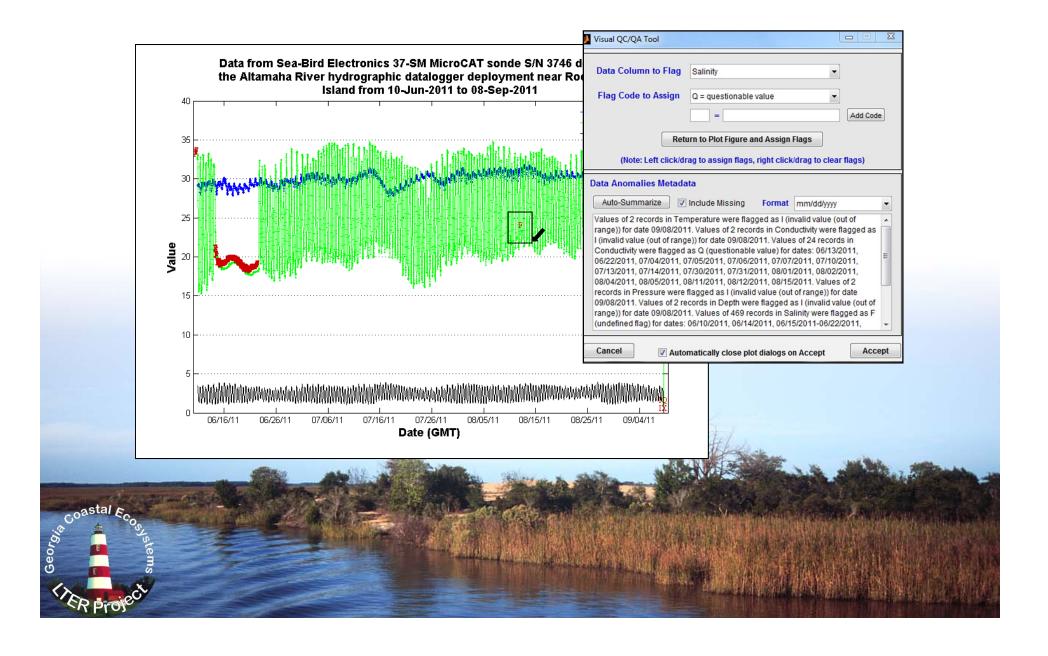
#### Data Viewer/Editor

File Edit Options													
All	Site	Longitude	Latitude	Instrument	Pump	Date	Year	Month (MM)	Day (DD)	A			
None	(none)	(degrees)	(degrees)	(none)	(none)	(serial day (base	(YYYY)						
<b>1</b>	7	-81.475500	31.338383	2398	0	733043.000000	2007	1	k S	1			
2	7	-81.475500	31,338383	2398	0	733043.020833	2007	1		1			
3	7	-81.475500	31.338383	2398	0	733043.041667	2007	1	1	1			
4	7	-81.475500	31.338383	2398	0	733043.062500	2007	1		1			
5	7	-81.475500	31.338383	2398	0	733043.083333	2007	1	ľ e	1			
6	7	-81,475500	31.338383	2398	0	733043.104167	2007	1	9	1			
7	7	-81.475500	31.338383	2398	.0	733043.125000	2007	1		1			
8	7	-81.475500	31.338383	2398	0	733043.145833	2007	1	1	1			
9	7	-81.475500	31.338383	2398	0	733043.166667	2007	1	F 3	1			
10	7	-81.475500	31,338383	2398	0	733043.187500	2007	1	1	1			
11	7	-81.475500	31.338383	2398	.0	733043.208333	2007	1		1			
12	7	-81.475500	31.338383	2398	0	733043.229167	2007	1		1			
13	7	-81.475500	31.338383	2398	0	733043.250000	2007	1		1			
14	7	-81,475500	31.338383	2398	0	733043.270833	2007	1		1			
15	7	-81.475500	31.338383	2398	0	733043.291667	2007	1		1			
16	7	-81.475500	31.338383	2398	0	733043.312500	2007	1		1			
17	7	-81.475500	31.338383	2398	0	733043.333333	2007	1	K S	1			
18	7	-81.475500	31.338383	2398	0	733043.354167	2007	1		1			
19	7	-81.475500	31.338383	2398	0	733043.375000	2007	1	ė	1			
20	7	-81.475500	31.338383	2398	0	733043.395833	2007	1		1			
21	7	-81.475500	31.338383	2398	0	733043.416667	2007	1	ř F	1			
22	7	-81,475500	31.338383	2398	0	733043.437500	2007	1		1			
23	7	-81.475500	31.338383	2398	0	733043.458333	2007	1		1			
24	7	-81,475500	31.338383	2398	0	733043.479167	2007	1		1			
25	7	-81.475500	31.338383	2398	0	733043.500000	2007	1	r s	1			

## Data Search Engine



## Interactive Plotting & Q/C Tools



## **Key Concepts**

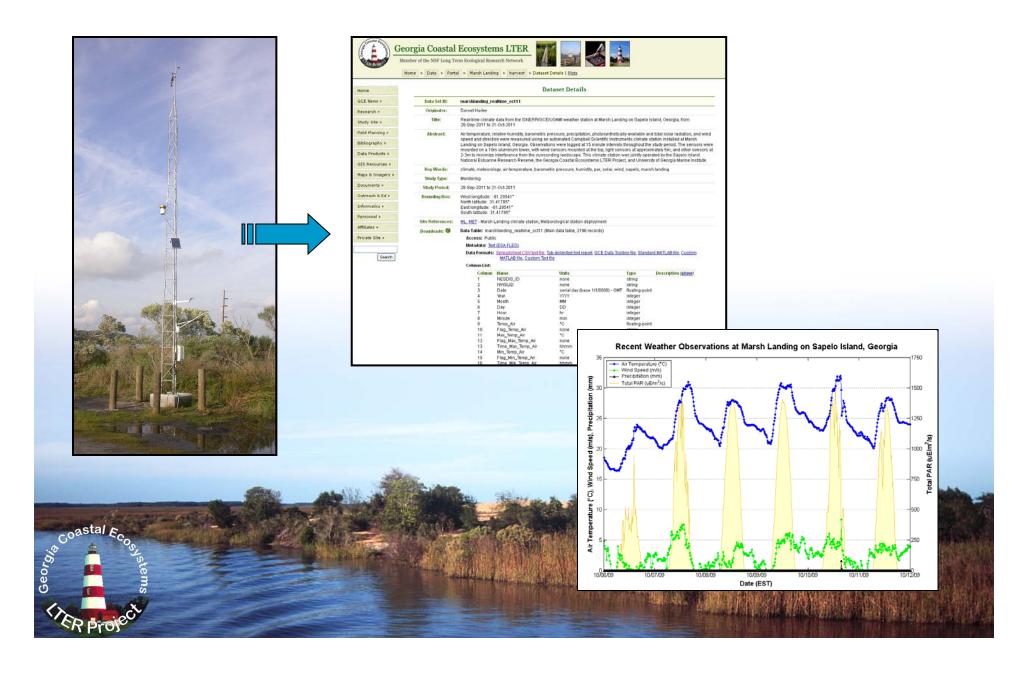
- Every operation is performed in context of a "dataset"
  - > Passing data columns to a tool transports metadata as well
  - Dataset metadata used to guide transformation, plotting, analysis
  - Metadata used to auto-parameterize functions
- Data structure instances are independent
  - > Each step along a workflow results in a complete data set with metadata
  - Intermediate datasets can be saved or overwritten in workflows
- Processing history ("lineage") information captured for all steps
  - Each tool logs operations by date/time
  - > Data revisions, deletions, flagging captured at user-specified detail
  - > Lineage reported in metadata
- Dataset metadata is "live", and updated automatically
  - Attribute changes
  - > Calculations, unit conversions
  - Code definitions

## Suitability for Real-Time Sensor Data

- Good Scalability
  - Data volumes only limited by computer memory (tested >2 GB data sets)
  - Multiple instances can be run on high-end, 64bit, clustered workstations
  - Good flag evaluation performance in use, testing with diverse rule sets
- Good scope for automation
  - Command-line API for unattended batch processing via workflow scripts
  - Timed and triggered workflow implementations easy to deploy
- Support for multiple I/O formats, transport protocols
  - > Formats: ASCII, MATLAB, SQL, specialized (CSI, SBE, NWIS RDB, HADS, ...)
  - Transport: local file system, UNC paths, HTTP, FTP, SOAP
- Already used for real-time GCE data, USGS data harvesting service (LTER HydroDB, CWT)



# Real-Time GCE Data Harvesting



## Implementation Scenarios

- End-to-End Processing (logger-to-scientist)
  - Acquire raw data from logger, file system, network (CIFS,HTTP,FTP,SOAP)
  - Assign metadata from template or using 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 their desktop

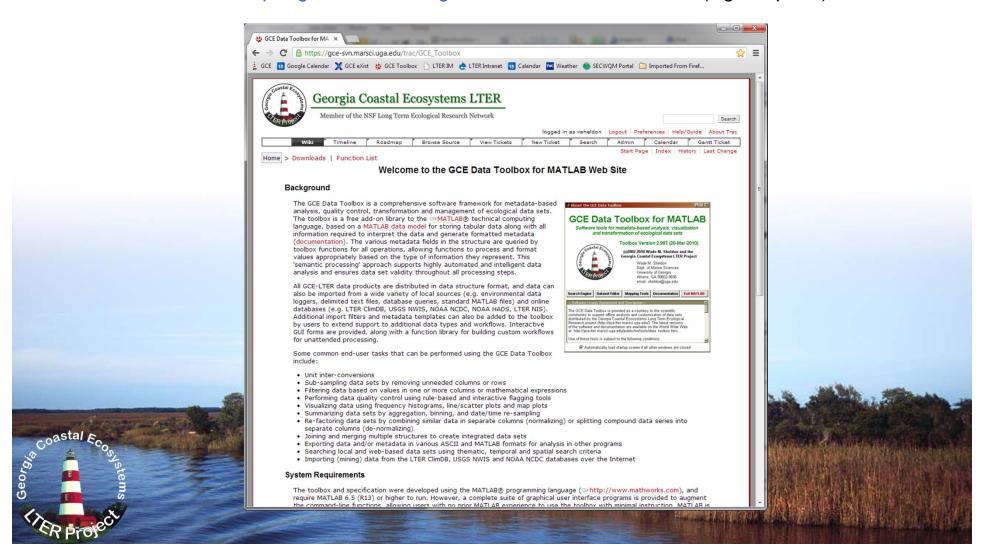
#### Data Pre-processing

- Acquire, validate and flag raw data (on demand or timed/triggered)
- Upload processed data files (e.g. csv) or value & flag arrays to RDBMS (e.g HIS)
- Workflow Step
  - Call toolbox from other software as part of workflow (e.g. LoggerNet)
  - DataTurbine via MATLAB off-ramp or Java API



## Toolbox Code & Support

- Trac support web site: https://gce-svn.marsci.uga.edu/trac/GCE\_Toolbox
- SVN address: https://gce-svn.marsci.uga.edu/svn/GCE\_Toolbox/trunk (login required)



#### Toolbox Timeline

- 2001 Initial toolbox development completed in May 2001, Metabase MMS online Sep 2001
- 2002 Added basic GUI interface, released code to GCE affiliates
- 2003 Added dynamic data harvesting support (USGS, NOAA, CSI LoggerNet);
   automated USGS harvesting service for ClimDB/HydroDB
- 2004 Added "search engine" tool for local search/integration of data
- 2005 First public distribution of "compiled" code; source code on request to LTER sites
- 2006 Added ClimDB data mining GUI
- 2007 Added enhanced data synthesis, refactoring tools
- 2008 Added GUI for managing QA/QC rules in metadata templates, additional flag tools
- 2009 Refined XML schema for formatted metadata; code moved to SVN; CWT adopts toolbox
- 2010 Toolbox released as open source (GPLv3); Trac support site established
- 2011 Expanded QA/QC tool options, GUI tools, refinements; focused on usability
- 2012 Added EML support, GUI for batch processing (import/export); ARRA funding received; first training workshop held



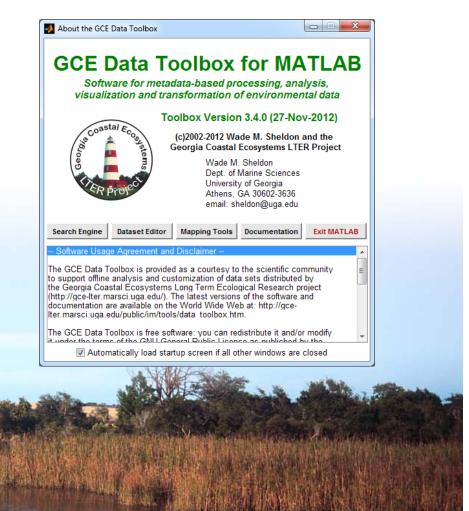
#### What's Next?

- Shameless plug: help us decide!
- Metadata model enhancements
  - ➤ Add "schema" support for better EML, Metabase alignment
  - Add EML export, enhance EML import
  - Better Metabase integration (fully bidirectional push/pull)
- Improve harvest management
  - > Add GUI tools for configuring harvest info, plot info
  - Add GUI for managing timers
  - Add data harvesting "dashboard" for monitoring activity
- Improve documentation and training materials



# Interactive Training – Day 1

- Installing and starting the GCE Data Toolbox
- Introduction to the Data Set Editor application
- Importing and exploring data
  - Generic ASCII
  - Specialized logger formats
- Metadata management
  - > Defining attribute metadata
  - Documentation metadata
  - Metadata templates
- QA/QC framework
  - Defining flags
  - Creating "rules"
  - Visual QA/QC
  - Copying and locking/unlocking flags
- Creating and exporting products
  - Batch processing raw data
  - Integrating data (join, merge)
  - Exporting



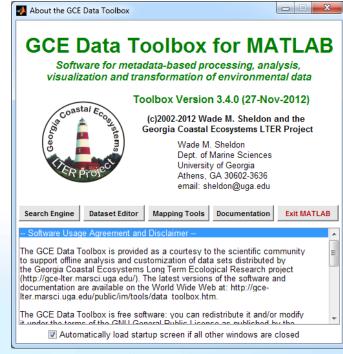
# Interactive Training – Day 2

#### Automating data processing

- Intro to command line API
- Scripted workflows
- Data harvesting support
- Harvest timers

#### Working with your own data

- > Choosing an import filter
- Defining metadata, templates
- Post-processing, analysis
- Harvesting scenarios
- > ...





# Interactive Training – Day 3

#### Metabase MMS & GCE Toolbox

- Intro to the Metabase MMS
- Services and applications
- > EML & NIS support
- Metabase-Toolbox integration
- > Future plans

#### Open discussion

- Training feedback
- Next steps
- "Missing links"
- **>** ...

