# A Web-based Tool for Translating Unstructured Data from Dataloggers into Standard Formats

Sean C. Arms [1]

Jennifer Oxelson Ganter [1]

Jeff Weber [1]

Mohan K. Ramamurthy [1]

[1]UCAR/Unidata

# Overview

- The Problem: Data Friction

- The Logger Problem

- Standard Formats to the Rescue (...or not)

- ρζητα

  - Architecture

  - Workflow

  - Current Status

Advanced Cooperative Arctic Data & Information Service

ACADIS

unidata

ρζητα

# Data Friction

- The Problem: Data Friction

- Examples:

  - The data I need is stuffed into a netCDF file, but I don't know how to use it and don't have the time to learn.

    - I might like to learn, but I don't have the time.

    - Just what exactly IS netCDF? Looks like Egyptian Hieroglyphics.

# Data Friction

# Data Friction

- The Problem: Data Friction

- Examples:

  - My data file is in netCDF, but I don't know how to use it and don't have the time to learn.

    - I might like to learn, but I don't have the time.

    - Just what exactly IS netCDF? Looks like Egyptian Hieroglyphics.

  - I love csv files, but the file layout always changes!

    - and…And…AND half the time there isn't enough information available for me to feel comfortable using the data!

# The "Logger Problem"

- What's the big deal?
  - Vast amount of datalogger output available
    - Typically in ASCII CSV format
    - *N* number of files, *~N* number of layouts
    - Not in a format to easily enable search, subset capabilities, other services
  - Value* added by placing into spreadsheet
- The "Logger Problem" is Huge for the Advanced Cooperative Arctic Data and Information Service (ACADIS) Project

# Standard Formats to the Rescue!

- Solution: Standard Formats!

# Standard Formats to the Rescue...
## or not...

- Solution: Standard Formats ☹
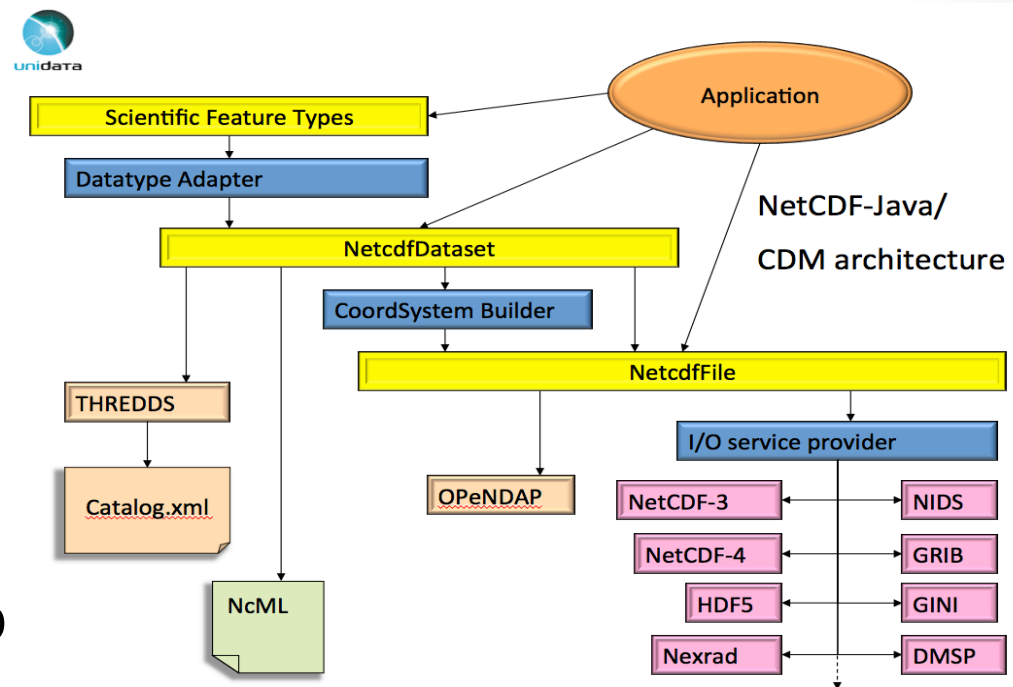
# Standard Formats to the Rescue… or not…

- Solution: Standard Formats ☹

- Why this Does Not Work:
  - Data come in a format that works for the PIs
  - Why put effort into transforming into a a new format with very specific conformance constraints (e.g. CF)
  - Usually code must be written to convert "useful" dataset into "standard" dataset
    - Even more work to do _before getting to the science stuff_!

8

# ρζητα

- What is ρζητα?
  - Vision: General Purpose Data Format converter
  - Goal: Get data into standard format *while providing data users with the format they want*
  - Why?
    - Enable services such as search, subsetting, aggregation, etc. for observational datasets, without writing readers for each "flavor" of ASCII data

# ρξητα - Architecture

- Basic Idea
  - Common Data Model (CDM)
    - netCDF-Java
  - I/O service provider (IOSP)
  - Web interface for collection of metadata
    - as needed
- The Idea is to get into CDM, then use IOSPs



10

# ρζητα - Architecture

- Java WebApp
  - Java
  - Apache Tomcat
  - Spring MVC
  - netCDF-Java(CDM)
  - JavaScript
    - jQuery, SlickGrid, jWizard

# ρζητα - Workflow

- Workflow
  - Define Discrete Sampling Geometry

# Define Discrete Sampling Geometry

# ρζητα - Workflow

- Workflow
  - Define Discrete Sampling Geometry
  - Upload CSV, XLS(X)

14

# Example Input

# ρζητα - Workflow

- Workflow
  - Define Discrete Sampling Geometry
  - Upload CSV, XLS(X)
  - Define Parsing Information

16

# Define Parsing Information

# ρζητα - Workflow

- Workflow
  - Define Discrete Sampling Geometry
  - Upload CSV, XLS(X)
  - Define Parsing Information
  - Define Variable Metadata

# Define Variable Metadata



19

# Define Variable Metadata

# ρζητα - Workflow

- Workflow
  - Define Discrete Sampling Geometry
  - Upload CSV, XLS(X)
  - Define Parsing Information
  - Define Variable Metadata
  - Define Global Metadata

# Define Global Metadata

# ρζητα - Workflow

- Workflow
  - Define Discrete Sampling Geometry
  - Upload CSV, XLS(X)
  - Define Parsing Information
  - Define Variable Metadata
  - Define Global Metadata
  - Transform

23

# ρζητα - Workflow

- Workflow
  - Define Discrete Sampling Geometry
  - Upload CSV, XLS(X)
  - Define Parsing Information
  - Define Variable Metadata
  - Define Global Metadata
  - Transform
  - Return netCDF and Transaction Receipt (NcML file)

# netCDF file (CF-1.6 Compliant)



```
1. lesserwhirls@micromac: /Users/lesserwhirls/Desktop (less)
netcdf ilu01_07_10 {
dimensions:
        time = 1036 ;
        name_strlen = 3 ;
        station_id_strlen = 3 ;
variables:
        float lat ;
                lat:units = "degrees_north" ;
                lat:long_name = "latitude" ;
                lat:standard_name = "latitude" ;
        float lon ;
                lon:units = "degrees_west" ;
                lon:long_name = "longitude" ;
                lon:standard_name = "longitude" ;
        float alt ;
                alt:units = "meters" ;
                alt:long_name = "height above mean sea-level" ;
                alt:positive = "up" ;
                alt:axis = "Z" ;
                alt:standard_name = "height" ;
        char station_id(station_id_strlen) ;
                station_id:cf_role = "timeseries_id" ;
                station_id:long_name = "station name" ;
                station_id:standard_name = "station_id" ;
        float time(time) ;
                time:standard_name = "time" ;
                time:long_name = "Time from datalogger" ;
                time:units = "days since 1970-01-01" ;
                time:coordVar = "yes" ;
                time:_columnId = "1" ;
        float soil_temperature_1(time) ;
                soil_temperature_1:missing_value = "-999" ;
                soil_temperature_1:standard_name = "soil_temperature" ;
                soil_temperature_1:valid_max = "-50" ;
                soil_temperature_1:long_name = "Soil Temperature at 25 cm" ;
                soil_temperature_1:source = "Thermometrics SN101" ;
                soil_temperature_1:comment = "Probe was exposed!" ;
```

Advanced Cooperative Arctic Data & Information Service
ACADIS
unidata
ρξητα

# Transaction Receipt (NcML file)



```xml
<?xml version="1.0" encoding="UTF-8" standalone="no"?>
<netcdf xmlns="http://www.unidata.ucar.edu/namespaces/netcdf/ncml-2.2">
  <dimension length="1036" name="time"/>
  <attribute name="Conventions" value="CF-1.6"/>
  <attribute name="featureType" value="timeSeries"/>
  <attribute name="title" value="Greenland Borehole Dataset"/>
  <attribute name="institution" value="University of Alaska, Fairbanks"/>
  <attribute name="processor" value="Vladimir E. Romanovsky"/>
  <attribute name="title" value="Greenland Borehole Dataset"/>
  <attribute name="comment" value="READ ME STUFF HERE :-)"/>
  <variable name="lat" type="float">
    <attribute name="units" value="degrees_north"/>
    <attribute name="long_name" value="latitude"/>
    <attribute name="standard_name" value="latitude"/>
    <values>69.290</values>
  </variable>
  <variable name="lon" type="float">
    <attribute name="units" value="degrees_west"/>
    <attribute name="long_name" value="longitude"/>
    <attribute name="standard_name" value="longitude"/>
    <values>51.0623</values>
  </variable>
  <variable name="alt" type="float">
    <attribute name="units" value="meters"/>
    <attribute name="long_name" value="height above mean sea-level"/>
    <attribute name="positive" value="up"/>
    <attribute name="axis" value="Z"/>
    <attribute name="standard_name" value="height"/>
    <values>10</values>
  </variable>
  <dimension length="3" name="name_strlen"/>
  <variable name="station_id" type="string">
    <attribute name="cf_role" value="timeseries_id"/>
    <attribute name="long_name" value="station name"/>
    <attribute name="standard_name" value="station_id"/>
    <values>ILU</values>
  </variable>
```

1. lesserwhirls@micromac: /Users/lesserwhirls/Desktop (less)

Advanced Cooperative Arctic Data & Information Service

ACADIS

unidata

26

# ρζητα - Status

- Current Status
  - Import Single-block CSV, XLS(X)
  - Produce netCDF with NcML Transaction Receipt
    - Prepared for data portal submission
  - Will soon be on GitHub
- Next Steps
  - Enable Mining of Header block
  - Import Multi-block CSV, XLS(X) files
  - Allow Upload of CF1.6 netCDF Discrete Geometry (A.K.A. point) Files
  - Return "Standard" CSV or XLS(X) format
- Somewhat Larger Goals
  - Enable desktop use for easy subsetting on CDM files
    - E.g. Easy grid point times series extraction from netCDF or GRIB files (returned as what the user would like, of course)

# ρζητα - Questions

- Unidata Funded by NSF 0833450 (AGS)

  - http://www.unidata.ucar.edu

- The Advanced Cooperative Arctic Data and Information Service (ACADIS) Funded by NSF 1016034 (ARC)

  - http://www.aoncadis.org/

ACADIS
Advanced Cooperative Arctic Data & Information Service

unidata

ρζητα