

Enhancing Data Sets for Accelerated Wind Energy Development

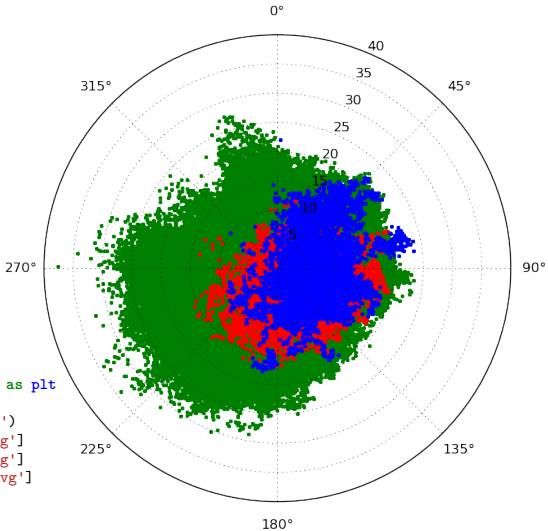
Advancing the state of the art in data set dissemination

Erik Quaeghebeur

Wind Energy Group – Delft University of Technology

OEEC 2017 – EUROS for wind energy

11 October 2017



```
import netCDF4 as nc
import matplotlib.pyplot as plt
```

```
f = nc.Dataset('Fino1.nc')
d = f['WV/Wr'][:, 3]['avg']
v = f['CA/Wg'][:, 6]['avg']
T = f['HTT/Lt'][:, 4]['avg']
freezing = T <= 0
warm = T > 20
```

```
sp = plt.subplot(1, 1, 1, projection='polar')
sp.set_theta_direction(-1)
sp.set_theta_zero_location('N')
plt.plot(np.radians(d), v, '.g')
plt.plot(np.radians(d[warm]), v[warm], '.r')
plt.plot(np.radians(d[freezing]), v[freezing], '.b')
```

Basic Claims

Problem

Current data dissemination practice *hampers* research.

Basic Claims

Problem

Current data dissemination practice *hampers* research.

Solution

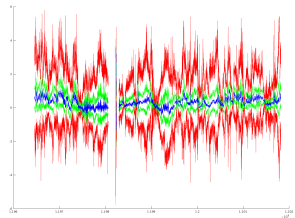
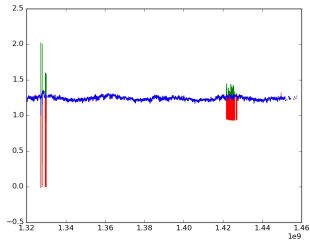
We have the tools to instead *facilitate* research.

Problem: Issues *accessing* the data

- scarcity of open data sets

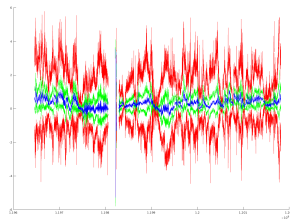
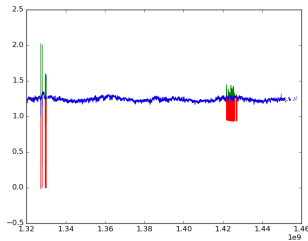
Problem: Issues *with* the data

- still faulty data included despite best-effort tests



Problem: Issues *with* the data

- still faulty data included despite best-effort tests

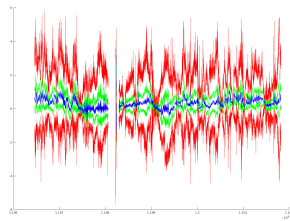
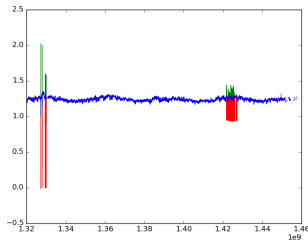


- no principled significance information

1.23456789 or 1.235 or 1.2?

Problem: Issues *with* the data

- still faulty data included despite best-effort tests





- no principled significance information

1.23456789 or 1.235 or 1.2?

- no information about missingness mechanisms
 - instrument defect? type?
 - eliminated during quality control?

Solution: Binary formats for enhanced data access

- established binary data formats

(HDF5 , netCDF4 )

Solution: Binary formats for enhanced data access

- established binary data formats

(HDF5 , netCDF4 )

- support in data analysis tools, network support (OPeNDAP)

Solution: Binary formats for enhanced data access

- established binary data formats

(HDF5 , netCDF4 )

- support in data analysis tools, network support (OPeNDAP)
- allows structured storage & various data types

Solution: Binary formats for enhanced data access

- established binary data formats

(HDF5 , netCDF4 )

- support in data analysis tools, network support (OPeNDAP)
- allows structured storage & various data types
- arbitrary key-value metadata (but use conventions!)

Metadata

```
<class 'netCDF4._netCDF4.Variable'>
compound Wr(time, level_UA)
  accuracy_abs: 1.0
  long_name: wind direction
  valid_range: [ 0. 359.]
  standard_name: wind_from_direction
  units: °
  sampling_frequency: 50.0
  cell_methods: ['avg'] time: mean (interval: 10 minutes ...)
                ['std'] time: standard_deviation (inter...)
  accuracy_propagated_avg: 0.0057735
  standard_error_avg: 0.0057735
  standard_error_std: 0.00408255
  accuracy_propagated_std: 0.0057736
compound data type: {'names':['avg','std','flag'], ...}
path = /UA
unlimited dimensions:
current shape = (683856, 3)
```

Solution: Binary formats for enhanced data access

- established binary data formats

(HDF5  , netCDF4 )

- support in data analysis tools, network support (OPeNDAP)
- allows structured storage & various data types
- arbitrary key-value metadata (but use conventions!)
- integrated compression & check-summing

Solution: Ensuring quality with improved processes

- add to existing set of automated tests
 - leverage available metadata more fully
 - more statistical analyses
 - build a shared database of issues and tests

- automatically calculate and encode per-value error bounds

Solution: Ensuring quality with improved processes

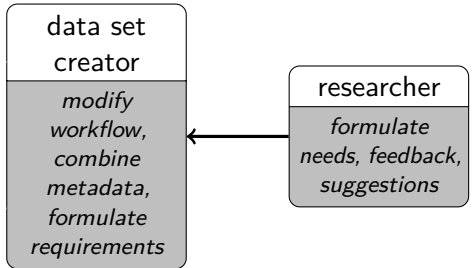
- add to existing set of automated tests
 - leverage available metadata more fully
 - more statistical analyses
 - build a shared database of issues and tests
- automatically calculate and encode per-value error bounds
- rigorous documentation of data set creation process
- versioning and unique identification of data sets

What must we do to make this vision come true?

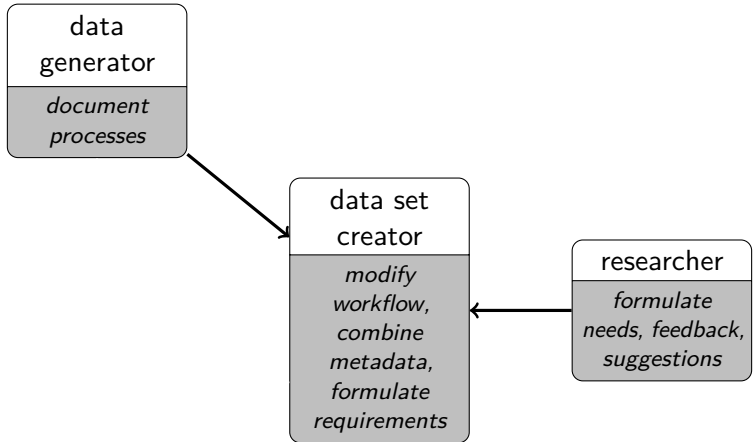
data set
creator

*modify
workflow,
combine
metadata,
formulate
requirements*

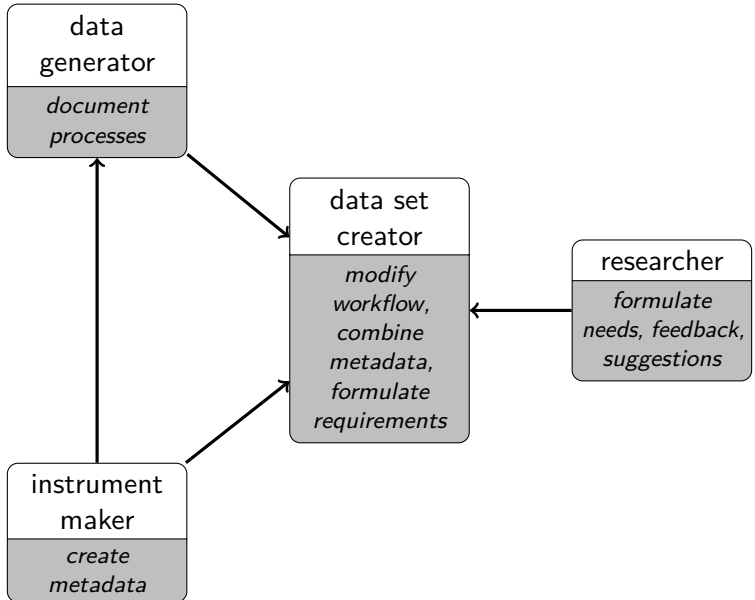
What must we do to make this vision come true?



What must we do to make this vision come true?



What must we do to make this vision come true?



Looking beyond

- role of data set (warehouse) managers?
- not offshore wind-specific
 - share experiences with other communities
 - join forces with them
- influence metadata conventions
- provide feedback to binary format designers

Basic Claims Revisited

Problem

Current dissemination practice *hampers* research.

- unnecessarily difficult to access data
- documentation detached from data sets
- still too much faulty data goes unnoticed

Basic Claims Revisited

Problem

Current dissemination practice *hampers* research.

- unnecessarily difficult to access data
- documentation detached from data sets
- still too much faulty data goes unnoticed

Solution

We have the tools to instead *facilitate* research.

- binary data formats such as HDF5 and netCDF4
- metadata attached to data sets
- improved processes (testing, encoding, documentation, . . .)

Thanks

- to the *people at ECN* for the help they provided
- to *Michiel Zaaijer* for focusing discussion
- and **to you for your attention**

Enhancing Data Sets for Accelerated Wind Energy Development

Advancing the state of the art in data set dissemination

Erik Quaeghebeur

Wind Energy Group – Delft University of Technology

OEEC 2017 – EUROS for wind energy

11 October 2017