# Versioning / Replication WG

## „Kick-Off"

# Versioning / Replication

- Basic Requirements / Assumptions

- Initial prototype
- Extension of publication procedure
- Publication policies and their enforcement
- Replication

- Next Steps / Roadmap

# Initial Requirements

- Bring versioning related information to end-users

- No replacement of current publishing process

- Definition of stable („core") APIs

- Enable automatic replication procedures

- Define „human ressource aware" roadmap

→ Close relation to publishing and QC working groups !

# Requirement 1: Bring versioning info to end users

**Needed:**

- Persistent identifier associated to file

- „core" metadata attached to identifier
  - ref to file, newer/older version, replica, checksum, date, ..

- Stable REST API to register/change/resolve PIDs and PID metadata

- Operational/scalable resolver system for PIDs

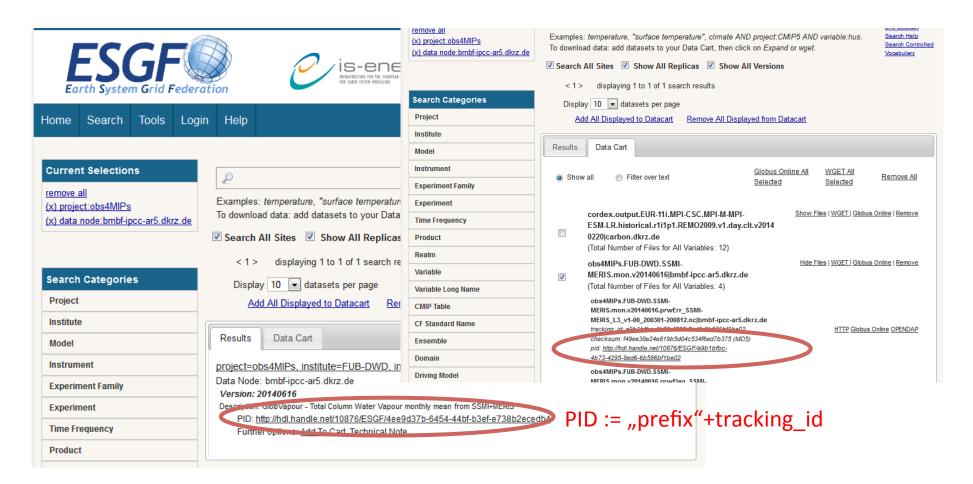# Requirement 1: Bring versioning info to end users

**Implementation options:**

A) Develop own solution and integrate with ESGF publisher

B) Take existing solution and integrate with ESGF publisher

→ Initial prototyping done following B)

**handle.net PID system:**
- stable API
- production ready, distributed, scalable resolution system
- existing large scale deployments (e.g. DOI system)

# Initial prototyping

„by hand" PID assignment for a smaller obs4MIPs project published at DKRZ:



PID := „prefix"+tracking_id

# Initial protoyping

**Tools**
- Publisher
- Replication service

**PID API**

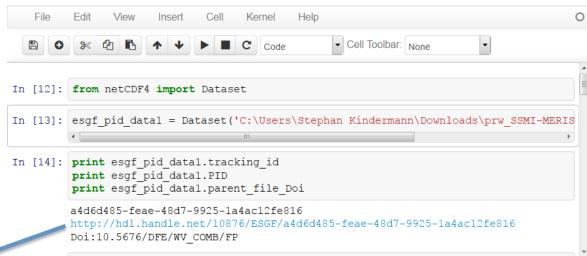**Resolver system**
(handle proxy server sytem)
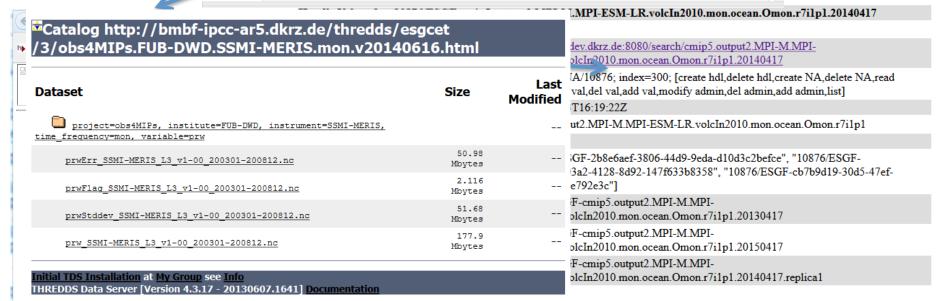
**PID Metadata**
- Data Url(s)
- checksum
- …

redirect

**Web Landing Pages**
- Data download page
- Versioning Info
- Replica Info

IP[y]: Notebook  pid1

File    Edit    View    Insert    Cell    Kernel    Help

Code          Cell Toolbar: None

In [12]:  from netCDF4 import Dataset

In [13]:  esgf_pid_data1 = Dataset('C:\Users\Stephan Kindermann\Downloads\prw_SSMI-MERIS

In [14]:  print esgf_pid_data1.tracking_id
          print esgf_pid_data1.PID
          print esgf_pid_data1.parent_file_Doi

          a4d6d485-feae-48d7-9925-1a4ac12fe816
          http://hdl.handle.net/10876/ESGF/a4d6d485-feae-48d7-9925-1a4ac12fe816
          Doi:10.5676/DFE/WV_COMB/FP

# Initial protoyping

# Initial protoyping

# Next steps

Integration in publishing process

# The (modified) publication process

„intelligent" file services:
- newer versions
- nearest replica
- „persistent" wget scripts ..

Solr Index

PID
Resolver system

Thredds catalogues

Publisher DB

(+ persistent publishing history
+ annotation links)

*data publisher*

Publish/
Unpublish

PID
API

PID
Registration
Agency
(for „prefix")

„curated ESGF files"

„ESGF files"  (+ „actionable tracking id" := prefix + tracking_id)

*modeling center*  cmor  (+ config option for PID prefix)

raw files

# Next steps: Agreements / Policies

- Restrict user definable publication options
  - publish/unpublish only, automatic versioning and replication option settings

- PID assignment is part of an atomic publication: pid assignment failure – publication failure

- Assignment of PID prefix to modeling centers

- Commitment of some sites to run PID Handle servers or establish liaison with existing PID sites
  - Long term commitment → careful planing !

# Next steps: Implementation

## We propose to follow the prototype:

- Integrate handle PIDs in publication process
- PID metadata includes basic versioning and replication information
- Versioning and publication history storage at data nodes → see prototype in QC WG

Enforcement of PIDs for data entities inline with other intitiatives:
- DataOne (https://mule1.dataone.org/ArchitectureDocs-current/design/PIDs.html)
- EarthCube
- RDA (http://rda.org)
- ANDS (http://ands.org.au/guides/persistent-identifiers-working.html)

# Replication: Next

→ Agree on one replication tool to work on

- syncrodata !?

→ Requirements list and priorities

- supported tranfer mechanisms (+ globus!?)
- monitoring ( ←→ icnwg working group)
- notification hooks and mechanism to enable automatic procedures

→ Definition of roadmap

# Summary

- Start from end user perspective: How to bring versioning (and replica) info to end users (and later data-evaluation wflows - for provenance tracking)
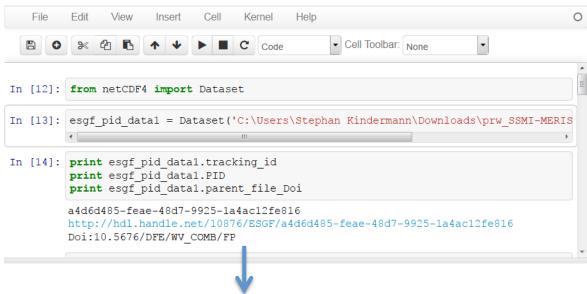
**~ Plan:**

- February/March detailed work plan with ressource estimation
- March: ESGF PID scenario presentation at RDA meeting in San Diego



- June: working publication add on prototype – additional sites running Handle service
- August/September: Intensive testing – tuning
- November / December: integrate as optional part in ESGF publisher
  (requirement for ESGF projects with strong data curation requirements)

# Initial protoyping



PID → DOI transition strategy: Later step