# Goal

Capture provenance data relating to the "climategate" data using the W3C provenance model with Annalist

# Source

## INTERNATIONAL JOURNAL OF CLIMATOLOGY

RESEARCH ARTICLE

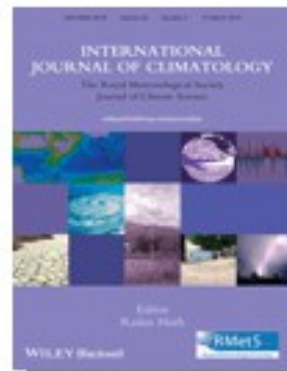### Updated high-resolution grids of monthly climatic observations – the CRU TS3.10 Dataset

I. Harris[1], P.D. Jones[1,2,*], T.J. Osborn[1] and D.H. Lister[1]

Issue

SEARCH

In this issue

Advanced

ARTICLE T

Get PDF

Save to

E-mail L

Export

Get Cita

Request

Am score   3

Additional Information (Show All)

# What we did

1. Create Annalist collection and GitHub repo

2. Created initial Annalist data definitions, based on W3C PROV

3. Glean initial high level provenance from paper

4. Refine, iterate

# Outcome 1 - online provenance capture/demo tool

http://demo.annalist.net/annalist/c/Climate_data_provenance/

# Outcome 2 – Annalist data definitions (JSON) in GitHub

https://github.com/gklyne/Climate_data_provenance

# Summary

- Manual provenance extraction is tedious

  - automated capture would be much easier!

- Have created initial set of Annalist definitions for provenance

  - next time will be so much easier

- It proved very effective to start with a very high level view of provenance, and iterate

  - be agile

- Learned more about areas where Annalist usability still needs work