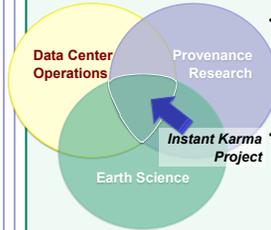# Instant Karma Status Update: Provenance at the AMSR-E SIPS

Helen Conover[1], Beth Plale[2], Mehmet Aktas[2], Bruce Beaumont[1], Dawn Conway[1], Sara Graves[1], Scott Jensen[2], Harsh Joshi[2], Ajinkya Kulkarni[1], Yuan Luo[2], Robert Ping[2], Prajakta Purohit[2], Rahul Ramachandran[1], Kathryn Regner[1], Cara Stein[1]

[1]University of Alabama in Huntsville    [2]Indiana University

## Approach



- Collaboration among
  - AMSR-E SIPS (MSFC Earth Science Office and UAHuntsville ITSC)
  - Provenance researchers at Indiana University's Data to Insight Center
  - AMSR-E Sea Ice science team (GSFC)
- Primary goal is to improve the collection, preservation, utility and dissemination of **provenance information** within the NASA Earth Science community
  - Using Karma provenance tool
  - Initial focus on Sea Ice processing

The Instant Karma project will integrate Karma, a provenance collection and representation tool developed at Indiana University, into the AMSR-E Science Investigator-led Processing System (SIPS) production environment, managed jointly by NASA/MSFC and UAHuntsville. The AMSR-E SIPS generates Level 2 and Level 3 data products from AMSR-E observations. An initial focus on Sea Ice processing will allow the project to engage the Sea Ice science team and user community in customizing Karma for NASA science data.

## Provenance Collection and Storage

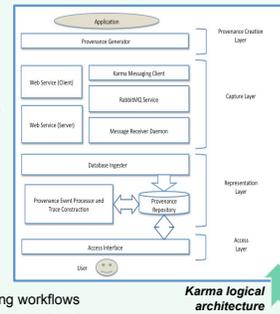AMSR-E SIPS processing workflow for Sea Ice instrumented in the testbed environment.

- Provenance information is captured in experiment run log files.
- Log files are parsed to generate provenance notifications.
- These notifications are then imported into the Karma database.
- The Karma Service Query API is used to generate OPM-compatible XML graphs, each corresponding to a processing run.

*Note that several of the services in the sea ice workflow are housekeeping and processing automation scripts, which are part of the processing workflows for other AMSR-E daily products.*

## Karma Provenance Collection Tools



*Karma logical architecture*

- Efficient and lightweight tools that support provenance collection, representation, and use
- Modular and programmable
  - Support diverse workflow architectures that consist of web services, java classes, message bus listeners
- Capture provenance in streaming workflows
  - No need to know workflow structure in advance
- Support interoperability
  - Implement Open Provenance Model (OPM) v1.1* to represent provenance graph (access interoperability)
  - OPM enables provenance information exchange with other OPM-compliant tools
- Recent redesign of internal database schema and data structures represents Earth science relevant provenance more efficiently

* http://eprints.ecs.soton.ac.uk/16148/1/opm-v1.01.pdf



*Processing Testbed*
*Provenance Collection*
*Provenance Browser*
*Query API*
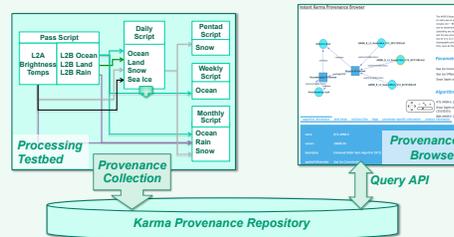*Karma Provenance Repository*

## Provenance and Context Information

*Lots of information already available, but scattered across multiple locations*

- Processing system configuration
- Dataset and file level metadata
- Processing history information
- Quality assurance information
- Software documentation (e.g., algorithm theoretical basis documents, release notes)
- Data documentation (e.g., guide documents, README files)

*Instant Karma project aims to collate and organize information from multiple sources*

## Defining and Collecting Science-Relevant Provenance and Context
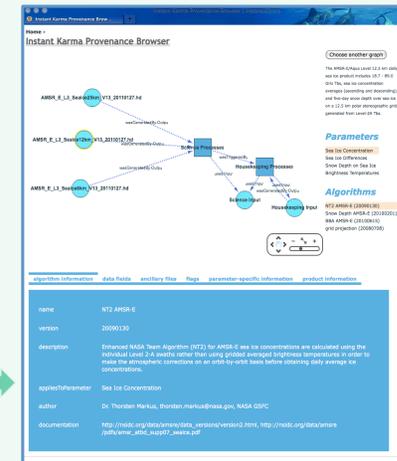
- Harvesting granule information from ECS metadata
- Also recording processing location associated with each data granule
- Working with AMSR-E Science Computing Facility to identify algorithm and data product information
  - Algorithm versions and descriptions
  - Parameters and data fields
  - Ancillary files
  - Flag values and explanations
  - Pointers to full documentation
- Defining how to harvest, transmit and display this information

## Science Use Cases

- ✓ Browse provenance graphs : convey rich information about final data granule details [Use case 1]
  - Spatial location, time of observation, algorithms employed, input data and ancillary files
  - Provenance bundle to include pointers to relevant documentation
- ✓ Answer "Something isn't right" question [Use case 1 variant]
  - E.g., did not receive data for several days so snow melt mask may be inaccurate.
- Compare two data granules [Use case 2]
  - Query system to get list of provenance differences (e.g., versions of software, number and versions of input files)
- ✓ General provenance graph for a given science process, e.g., Sea Ice processing [Use case 3]
  - Current algorithms and versions, nominal number and versions of input files, pointers to relevant documentation
- Embed provenance information as annotations in HDF files

## Browsing Provenance Information

- Interactive web application allows users to view the provenance graph for a specified data product
- Click on a node to display the full description of the product or process
- Trace full lineage of a data product by viewing the provenance information for each input file
- Access relevant information for the data product
  - Algorithm Theoretical Basis Document
  - README files
  - Product and inventory level metadata
- Uses Karma Service Query API to extract provenance graphs from the Karma provenance repository.



*Browser prototype showing provenance graph and related information for generation of a daily 12.5 km Sea Ice product from AMSR-E Brightness Temperatures.*