# The MPO System for Automatic Workflow Documentation

## Gheni Abla
## for The MPO Team

**10th IAEA Technical Meeting on Control, Data Acquisition and Remote Participation for Fusion Research**

**20-24th April 2015**
**Ahmedabad, Gujarat, India**

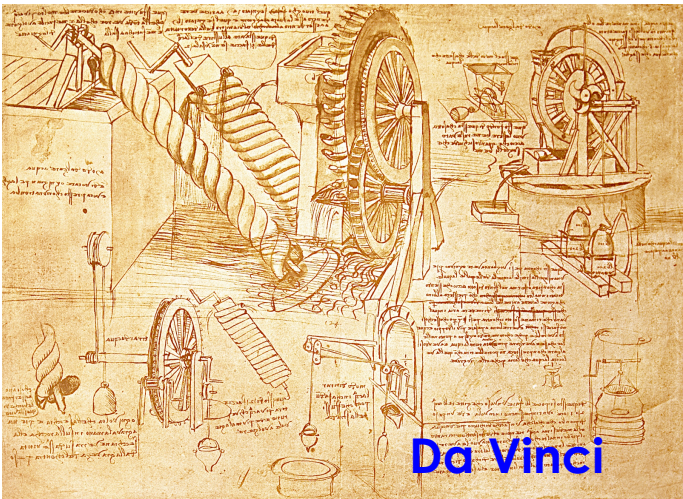DIII-D NATIONAL FUSION FACILITY SAN DIEGO
BERKELEY LAB
PSFC
GENERAL ATOMICS

# This Presentation Builds on Previous Presentations

- M. Greenwald, et al., "A Metadata Catalogue for Organization and Systemization of Fusion Simulation Data", 8[th] IAEA-TM, San Francisco, CA, June 2011

- D.P. Schissel, et al., "Automated Metadata, Provenance Cataloging, Navigable Interfaces: Ensuring the Usefulness of Extreme Scale Data", 9[th] IAEA-TM, Hefei, China, May 2013

- J.C. Wright, et al., "The MPO API: A Tool for Recording Scientific Workflows", 9[th] IAEA-TM, Hefei, China, May 2013
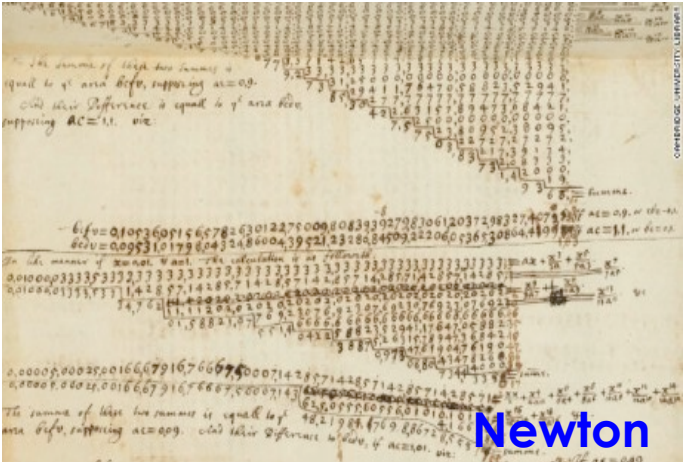
# Documenting Data and Processes is Important Aspect of Scientific Research Activities

- **Data from research activities is expensive to produce and may be critical for follow-on research**

- **It is not the mere existence of data that is important, but our ability to make use of it**

- **The context and metadata makes the data more usable**
  - Hypotheses
  - Pre-process activities
  - Experiments
  - Computational process
  - Reflections
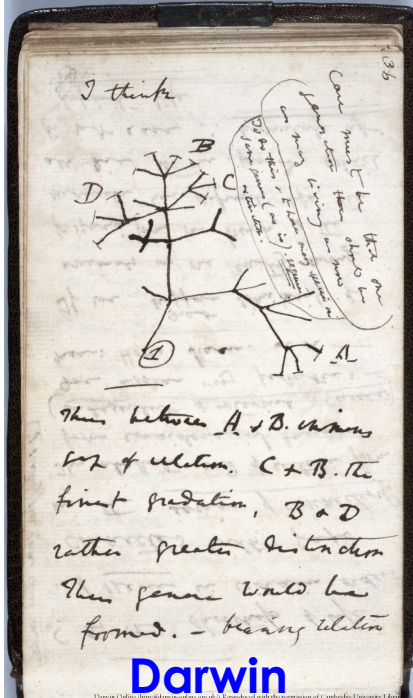
- **Documenting the process is not an easy task**

# Throughout History, Scientists Generated Handwritten Logbooks to Keep Track of Data
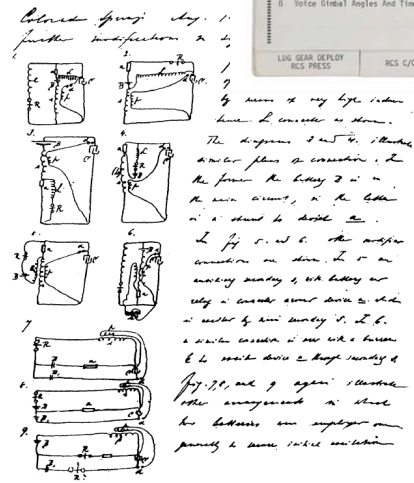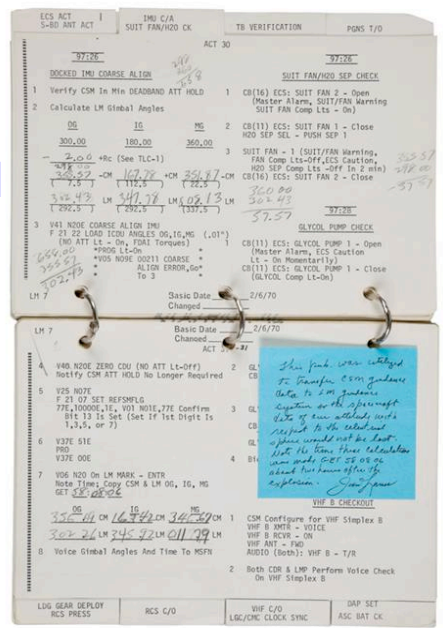


Da Vinci



Newton



Darwin



Lovell



Tesla

GENERAL ATOMICS

# In the Modern Era Documenting Research Process Made Progress and Met New Challenges

- **Personal computers and mobile devices helped electronic logbooks replace handwritten ones – brought conveniences**
  - Multi-media and hypertext support
  - Store, share and search

- **However, the content creation and log entry remained as a manual activity in the electronic logbooks**

- **As the pace of scientific research accelerated, documenting the process & data became more challenging & time consuming**
  - Increased precision of scientific instruments
  - Rise of exascale computing and arrival of Big Data
  - Result: fragmentation of data, processing, and documentation

# Metadata, Provenance and Ontology (MPO) System is for Documenting Scientific Data & Workflow

- **Provenance: Preserve meaning of data by documenting all of the steps taken to produce the data**
  - Automate metadata generation as much as possible
  - Support more systematic management of data used and resulted by analysis and simulation

- **Provide and preserve answers to two key questions:**
  - *Where did a particular piece of data come from?*
    - Assumptions, inputs, parameters used for calculation
    - The origin of inputs; reasons for assumptions & parameters
  - *Where was this data used?*
    - Other calculations
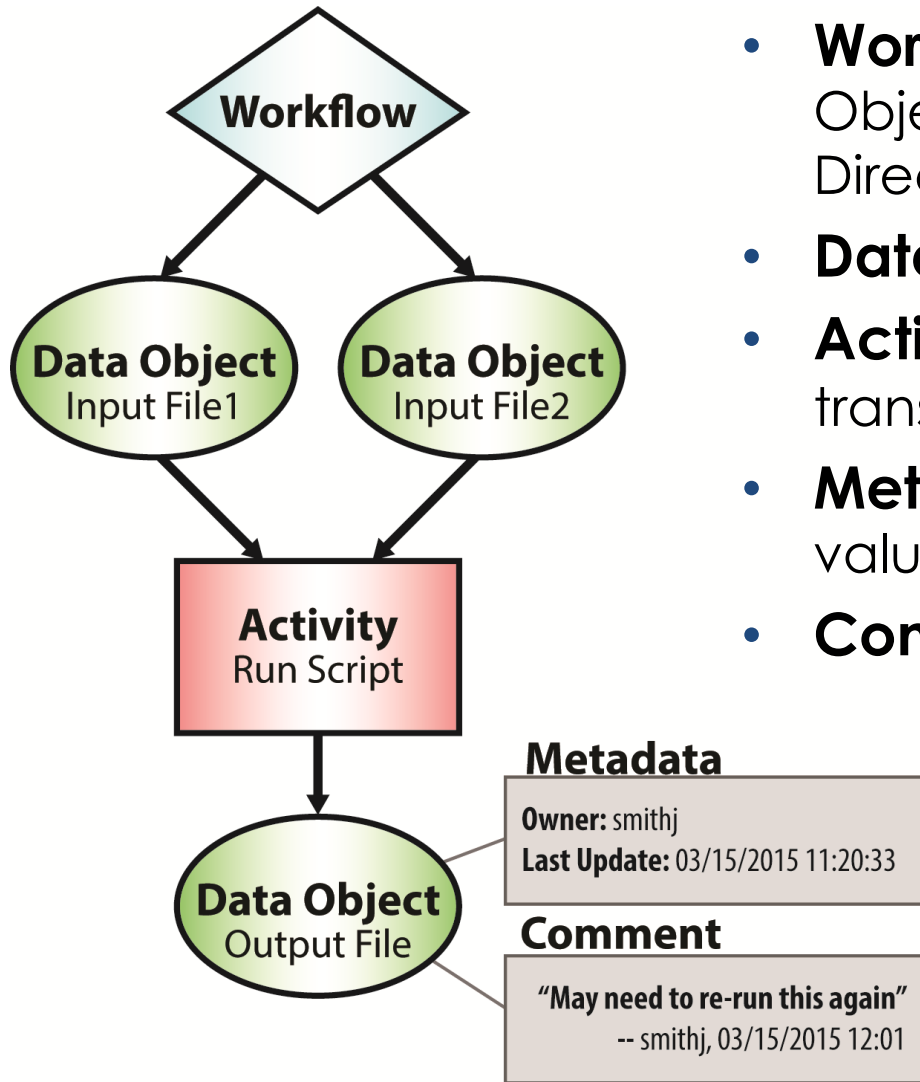    - Publication and presentation
    - Contributions to databases

# Potential Use Cases

- **How did I get the data plotted in Fig.6 of my 2014 Phys. Plasmas article?**

- **A calibration error was found in Thomson Scattering data taken during 2011 - the data has now been recalculated.**
  - Where was the old data used?
  - What publications used that data?   Were they critical for the published conclusions?
  - Did we contribute that data to a database shared by others?

- **A recently graduated PhD student left behind output from thousands of gyrokinetic simulations**
  - Which of these were used in her thesis?
  - Which might be useful in the future?
  - What were the inputs and parameters used in the interesting runs?

# Capabilities of the MPO System

- **Support all types of the scientific workflows – both experimental and computational**

- **Function in a heterogeneous environment and interoperate with workflow tools people are already using**
  - Many different computing platforms – laptop to supercomputer
  - Researchers use many different languages (Shell scripts, Python, IDL, Matlab, etc.) and tools to get their work done
  - Data is stored in different formats (MDSplus, HDF5, NetCDF, ASCII)

- **Once set up, work as automatically as possible**
  - Best suited for scripted rather than one-time use
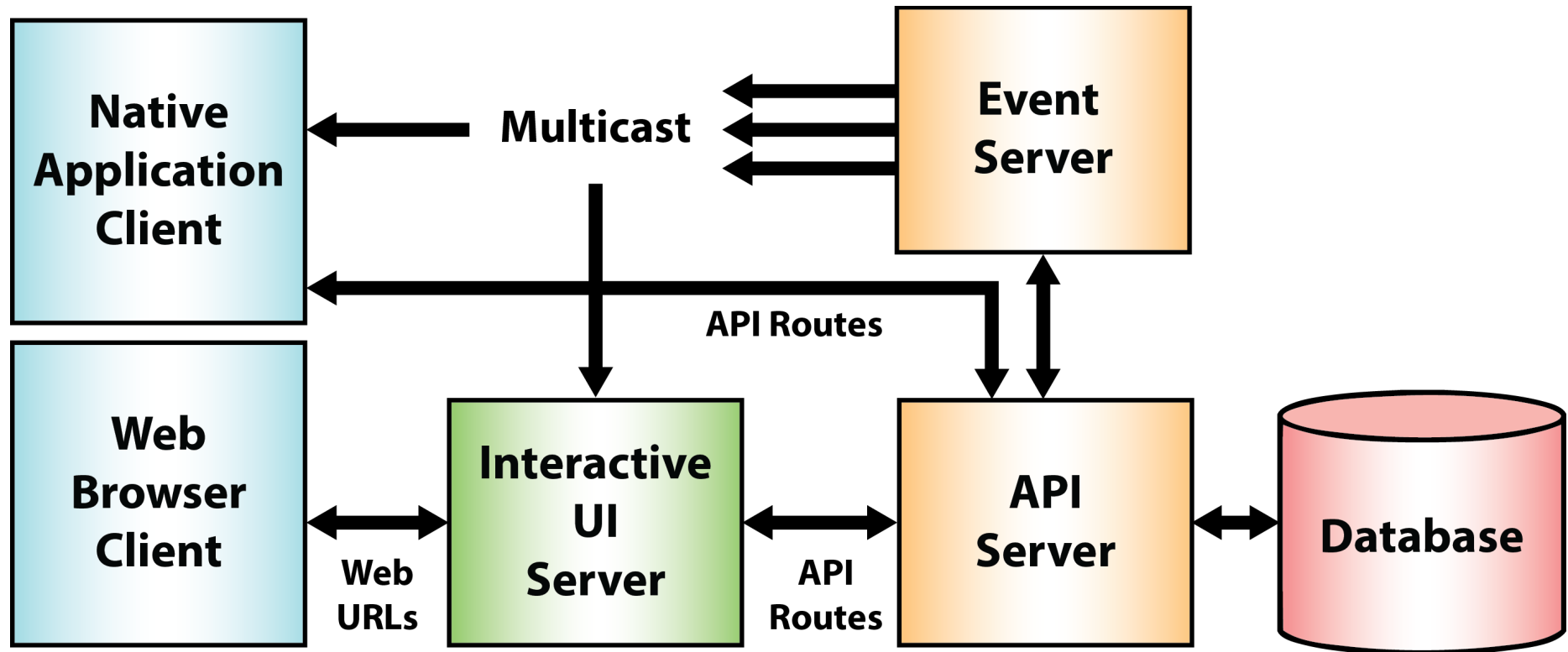
# MPO System Entities and Data Model



- **Workflow:** A series of connected Data Objects and Activities, organized as a Directed Acyclic Graph (DAG)
- **Data Object:** Structured data
- **Activity :** Process that creates, moves or transmutes data to a new form
- **Metadata:** Text-based, arbitrary name-value pairs
- **Comment:** User annotation

- **Connection:** The links between inputs, actions and results
- **Collection:** Group of Data Object, Activity, Workflow, and Collection

GENERAL ATOMICS

# MPO Entities Are Uniquely Identified

- **Each MPO entity is given a global unique numerical identifier**
  - UID – Unique ID
  - 128 bit, pseudo random numbers

- **Data object is also identified by a URI (Uniform Resource Identifier)**
  - The URI is the pointer to the data object
  - URI includes the data protocol name and the path to the data
  - Examples: MDSplus server+node path, HDF5 file +data location

- **Workflows also can be identified by composite ID**
  - Examples: doej/EFIT/52, smitha/OMFIT/1002

- **Searching is enhanced by defining a "controlled vocabulary"**
  - User-defined, hierarchical ontology

# The MPO System is Based on a Multi-Tier Software Architecture



Native Application Client

Web Browser Client

Multicast

Event Server

API Routes

Interactive UI Server

API Server

Database

Web URLs

API Routes

GENERAL ATOMICS

# MPO System Is Based on Open-Source Software

- **MPO Technology Stack**
  - "PostgreSQL" database used for current implementation
  - Both API server and Interactive UI server use "Flask", a lightweight Python web application framework
  - Twitter "Bootstrap" to create standardized Web front-end
  - DAGs rendered by "Graphviz" software
  - Authentication via x.509 certificates (currently support OSG, MIT & MPO certs)
  - MDSplus event services
  - SQLAlchemy for Object Relational Mapping

- **API is based on RESTful abstraction**
  - Services are exposed via RESTful methods (GET, POST) and URIs

**GENERAL ATOMICS**

# Interactive UI Page Example: Workflow List



**Enhanced "Controlled Vocabulary" Search:** User-defined, Hierarchical Ontology

Input a new comment

Display related comments

**Comments can be inserted/viewed directly on this listing page**

## MPO Workflows

**WORKFLOW**

- type = ALL
- time = _____ to _____

more... ▼

**ONTOLOGY**

- ACTIVITY
- GENERIC ▼
  - Status ▼
    - quality ▼
- WORKFLOW ▼
  - Type ▼
    - EFIT ▼
      - Code Characteristics ▼
        - bit_size ▼
        - documentation
        - purpose
        - source

| | CompositeID | Description | Creation Time | Comments | Quality |
|---|---|---|---|---|---|
| 1 | johnsonm / OMFIT / 31 UID: 9bd8133e-72bc-41da-8b61-c99⸱0178 | Multiple kinetic EFIT runs for 158152 | 2015-03-11 13:23:24 | 17 + | ★★ ★☆ |
| 2 | d3dauto / EFIT / 30 UID: 2f64832e-b113-4c0e-bfdc-7225c15d | | 8-11 13:23:19 | 1 + | ★★ ★★ |
| 3 | d3dauto / EFIT / 29 UID: c3cc4ab9-0d8b-4c48-84b1-343f97af | A first EFIT workflow | | | |

Submit

No MSE data after t = 5200 ms
— by carterp - 2015-04-07 14:18:23

**GENERAL ATOMICS**

# Interactive UI Page Example: Workflow Details



◇ **d3dauto / EFIT / 30**  *d3dauto* 3/11/2015, 1:23:23 PM

**Description:** EFIT02 for 161322

**UID:** `2f64832e-b113-4c0e-bfdc-7225c15d`

List of nodes and their corresponding details: UID, URI, metadata, comments, other linked workflows

**Green's Table**  ✕
*Green's table files*
2015-03-11 13:23:20
- *uri: filesys:///link/efit/rpf01.d3d*
- *uid: ad5f9396-9f7b-41ae-af9e-0232e95ee0eb*

💬 Add Comment

Snap file

Green's Table

Read Input Files

Read PTDATA

Calibrate Data

EFIT Data averaging

EFIT

Fit equilibrium

Blue nodes are Data Objects used in other workflows

**Workflow Details:**    💬 Add Comment    ⊕ Expand All

◇ d3dauto / EFIT / 30    3/11/2015, 1:23:23 PM

**Comments:**
💬 No MSE data after t = 5200 ms *by carterp 2015-04-07 14:18:23*

○ shot    2015-03-11 13:23:20

○ Snap file    2015-03-11 13:23:20

**Description:** EFIT input file mdsplus:///efit02/158025&path=\efit02:namelist

**UID:** `f3c51774-8473-4d2e-b5f9-e36170b3`

**Metadata:**
➕ last_update: 20150211 *d3dauto 2015-03-16 16:16:15.336422*
➕ owned_by: leexia *d3dauto 2015-03-16 16:16:39.372713*

○ Green's Table    2015-03-11 13:23:20

**Description:** Green's table files filesys:///link/efit/rpf01.d3d

**UID:** `ad5f9396-9f7b-41ae-af9e-0232e95e`

↓ *View linked Workflows*

▭ Read Input Files    2015-03-11 13:23:20.908132

○ Plasma Current    2015-03-11 13:23:21

# Interactive UI Page Example:  Collections List

## MPO Collections

| Name | Description | Username | Creation Time |
|------|-------------|----------|---------------|
| **OMFIT kinetic EFIT**<br>UID:<br>3f03306d-37da-4209-955e-fa13c16f | OMFIT kinetic EFIT runs for shots 158634-158640 | smitha | 2015-03-11 14:25:19.223908 |
| **Johnson IAEA talk April 2015**<br>UID:<br>4454ac1c-7c43-42e3-b67a-357ef27e | Collection of elements referenced in Johnson IAEA talk | johnsonm | 2015-03-11 14:25:28.355386 |
| **smitha's EFIT runs**<br>UID:<br>d18e5096-d4d0-4e40-a5e1-e52d536e | EFIT runs and snap files of interest from 6/2014 | smitha | 2015-04-07 13:53:37.678777 |
| **Collection of snap files**<br>UID:<br>fbe6a178-a62c-4896-9043-9dd23e38 | EFIT snap files used for 2014 MSE runs | smitha | 2015-03-11 14:25:22.843804 |

**Select to view details**

GENERAL ATOMICS

# Interactive UI Page Example: Collection Details

# Current Status

- **The MPO System V1.0 is released**
  - http://mpo.psfc.mit.edu provides detailed information

- **Python, IDL, shell API clients are provided**
  - Used to instrument MPO calls

- **Integration with multiple workflows**
  - DIII-D between-pulse EFIT
  - SWIM (Simulation of RF Wave Interactions with Magnetohydrodynamics)
  - GYRO (Nonlinear tokamak microturbulence software package )
  - AToM (Advanced Tokamak Modeling)
- **Planned Integration with a Climate Modeling Project–Calibrated and Systematic Characterization, Attribution and Detection of Extremes (CASCADE)**

# Future Work

- **Expand the reach of MPO framework**
  - Harden the system - ease of adoption, robustness, scalability
  - Reach out to more science domains – including non-fusion

- **Provide data exchange capability between MPO and W3C standard based software (e.g. PROV)**
  - W3C PROV
  - Import and export data

- **Improve user interface and analysis**
  - How to provide better/faster graphical navigation?
  - Additional visualizations and analysis

# Summary

- **MPO System is a software for documenting scientific workflows and data**
  - A new type of logbook with automation and analysis capabilities built-in

- **Production workflows have been MPO instrumented**
  - Proven useful
  - Approach is valid and general

- **MPO team seeks partners**
  - Test, deploy and feedback
  - Contribute
  - Contact email: mpo-info@fusion.gat.com

**GENERAL ATOMICS**

# Acknowledgments

- **MPO team members**
  - Gheni Abla, Liz Coviello, Sean Flanagan, Xia Lee, David Schissel– ***GA/DIII-D***
  - Alex Romosan, Arie Shoshani, John Wu – ***LBNL***
  - Martin Greenwald, Josh Stillerman, John Wright – ***MIT/PSFC***

- **Our colleagues**
  - Members of SWIM, AToM teams
  - Staff at GA, LBNL, MIT

- **This work is supported by US Department of Energy**
  - Office of Advanced Scientific Computing Research
  - Office of Fusion Energy Sciences