

## Implementation of data integration toolkit for ITER Physics Data Model

Xiaojuan LIU , Zhi YU, Nong XIANG

Institute of Plasma Physics , Chinese Academy of Sciences

Institute of Plasma Physics, Chinese Academy of Sciences

## Outline

- Background and Motivation
  - -Why do we need a data integration toolkit in the fusion research ?
- What is the **data integration ?** 
  - Design concept and framework
- Database mapping
  - Unified data access interface for experiment database.
- Lightweight data description
  - Unified data exchange interface between physical programs.
- Summary



## Motivation: Fusion research needs data integration tools.

- The magnetic confinement fusion is a complex large-scale scientific project with different time scales and space scales. Fusion research activities are mainly carried out around the experimental data and simulation data.
- ITER defined the Physical Data Model 1 (PDM) which aims at being the main gate and general standard to fusion data.
  - IMAS<sup>2</sup> being a particular case of application, choose the PDM as the backbone to realize the data interaction between physical modules and experimental by UDA and UAL
    - Integrate the data of different systems into the same system to realize the integrated analysis of the entire device. This process is the so-called

#### data integration process.

- The European Fusion Community recently proposed the FAIR4FUSION<sup>3</sup> project.
  - One of the purposes is to improve the sharing and accessibility of fusion data and to achieve data integration and standardization.
- We develope a data integration toolkit to provide a platform environment for data modeling and analysis in fusion research.



1.Imbeaux, Frederic A generic data structure for integrated modelling of tokamak physics and subsystems 2010 Computer Physics Communications , Vol. 181, No. 6 Elsevier p. 987-998 2.Imbeaux, Frédéric Design and first applications of the ITER integrated modelling & analysis suite 2015 Nuclear Fusion , Vol. 55, No. 12 IOP Publishing p. 123006 3.https://www.fair4fusion.eu/



codeparam 🕂

equilibrium [

Description of a 2D,

equilibrium code.

axi-symmetric, tokamak equilibrium: result of an

Time-dependent CPO

## Background: The de facto data standardization format and scheme

#### • Proprietary, non-universal data format:

- SWIM<sup>1</sup> plasma state file (which is different from the ONETWO plasma state file)
- Consistent Physical Object<sup>2</sup> is used by the ITM-TF,
- BPSD which is a "in-house" standard data structure is used in TASK<sup>3</sup>.
- OMFIT<sup>4</sup> treats files, data and scripts as a uniform collection of objects organized into a single, self-descriptive, hierarchical structure (the OMFIT tree structure).

#### • Universal, non-specific field data format:

- OMAS<sup>5</sup>: Ordered Multidimensional Array Structure
  - It provides a convenient Python API
  - it stores data with different file/database formats, mapping the physics codes I/O with OMFIT framework to the IMAS data model.
- HDC<sup>6</sup>: Hierarchical Data Containers
  - HDC is a tiny library for exchanging hierarchical data (arrays of structures) in shared memory between multiple programming languages, currently supporting C, C++, Python, Fortran and MATLAB.
  - The hierarchical structure is organized as a tree.



<sup>1.</sup>W.R. Elwasif et al., in Parallel, Distributed and Network\_x0002\_Based Processing (PDP), 2010 18th Euromicro Interna\_x0002\_tional Conf. on, Institute of Electrical and Electronics En\_x0002\_gineers (IEEE), pages 419427 (2010)

2..A. Fukuyama et al., in Proc. 20th Fusion Energy Conf. Villamoura, Portugal (2004).

3.B. Guillerminet et al., Fusion Eng. Des. 83, 442 (2008).

4.Meneghini, O., Smith, S. P., Lao, L. L., Izacard, O., Ren, Q., Park, J. M., ... Staebler, G. M. (2015). Integrated modeling applications for tokamak experiments with OMFIT. *Nuclear Fusion*, *55*(8), 083008. 5.Meneghini, Orso-MariaOMAS: A Python Library to Interface with the ITER Integrated Modeling and Analysis Suite (IMAS) 2020 APS Division of Plasma Physics Meeting Abstracts, Vol. 2020 6.JM10-0092.https://bitbucket.org/compass-tokamak/hdc/src/master/

## Fundamental design concept —— "Database is a collection of documents"

- Purpose :The data interfaces of different data sources are unified into the ITER physical data model.
- Data Intergration: Integrate different data sources into a **unified namespace**, which can be accessed through **a URI**.
- Namespace is a container that organizes and uniquely restricts content based on PDM
  - The top container is a collection of documents.
  - The document is a unified hierarchical tree structure that integrates different data objects.
- Data objects in Fusion can be divided into three categories :



- Configuration data (readonly, identified by 'device')
- Experimental data (readonly, identified by 'shot)
- Analysis data readonly, identified by 'run')

### Framwork



#### • A uniform access entry:

- a uniform URI determined by a variable key-value:
- a global unified request path which is defined by a hierarchical tree structure in the document

#### A virtual mediated conversion layer:

- virtual: not store any data
- **matching** the key-value in URI to determine the source of the data
- **mapping** the schema in request to the target schema in mapping file
- query reformulation is to translate the user queries into a set of subqueries based on source schemas.
- Different heterogeneous data sources:
  - device database
  - physical modules IO
  - IMAS, OMAS



## Data model and Data representation



- ASIPP
- Data format conversion and data schema mapping are separated. When the data model is changed, the content of the mapping file or configfile is modified without changing the API.

### Database mapping

- Unified data access interface for experiment database.
  - Original Request Configuratio Data-Data (static) Searching in Mapping Files Experimenta **TDI Request** Data Data or Request (dynamic) is Request? Yes No **MDSplus** ID **MDSplus** client Data Database Data Return result
- The mapping between data sources is the conversion of paths.
- Request is a PATH of data in a hierarchical tree-like structure;
- A configurable mapping rules are defined in mapping files.
  - Support the **merging** of multiple data sources:static data and dynamic data
  - The position of the data element in a **hierarchical tree** is located by the **path**.
    - static data : return result
    - dynamic data :return TDI request
- Data mapping is the conversion of paths between different data schema.



## Unified access to experiment databases

- A uniform URI with different key-value to identify two different databases.
- The request to access data from the two different databases is unified .
- The corresponding resultes are shown in the figure respectively.





## Lightweight data description mechanism

- Unified data exchange interface between physical programs.

-mapping -> equilibrium

-read core profiles

EAST

Database

IMAS

Database

- Requirement for physics codes:
  - The input parameters are dynamically changing and multi-source gathered from experimental data,

device data ,config file and the output of other physical modules.

- Data format and semantics of input/output parameters are specified by code themselves
- The process of preparing input parameters is actually to gathering data from different sources into the same semantics and expression .
- It is essentially a data integration process.
- Implementation:
  - One-to-one corresponding lightweight configuration description file :
    - The file defines the data semantics and data format of in/out parameters. when the module is called, the tool converts data into the format required by the code based on these descriptions.
- The data exchange between physics modules which are called through a



unified API is achieved in the data semantic layer centered on PDM.



## Lightweight data integration: Genray+Cql3d



## Summary

- A data integration toolkit was developed for experiments, modeling and simulation in fusion research.
- Using the ITER Physical Data Model as the global unified schema.
- The core design concept :
  - Hierarchically structured data can be organized into collections of documents.
  - Data is defined by data model and data representation, which should be separated in implementation.
- The Framework is divided into three layers:
  - A uniform access entry,
  - A virtual mediated conversion layer
  - Plugins of different heterogeneous data sources
- A **lightweight data description** is defined to unify the data exchange interface between the physics module.





# Thank you for attention



Institute of Plasma Physics, Chinese Academy of Sciences