

The IAEA Fusion Data Lake Project

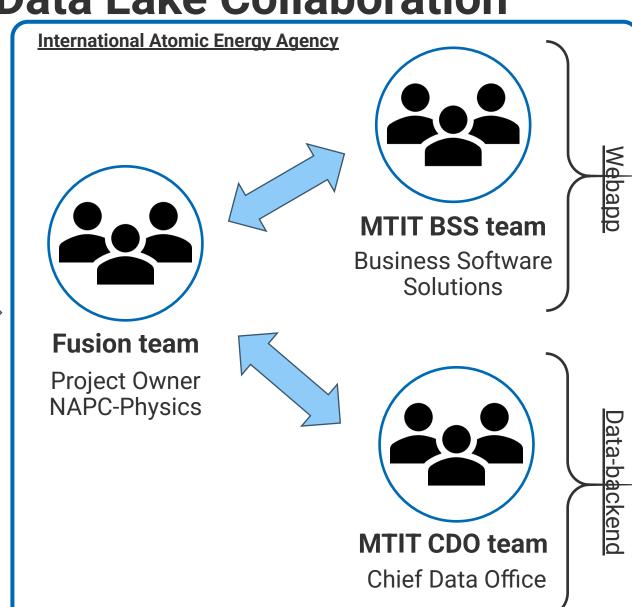
Accelerating AI and Big Data Applications through Open Science and FAIR Data

D.S. Gahle, M. Barbarino, et al, IAEA 6th IAEA Technical Meeting on Fusion Data Processing, Validation and Analysis September 9th-12th, Fudan University, Shanghai

Al for Fusion and Fusion Data Lake Collaboration

**IAEA AI for Fusion CRP Partners

The Australian National University (Australia), Institute of Plasma Physics, Chinese Academy of Sciences (China), Shanghai Jiao Tong University (China), Southwestern Institute of Physics (China), Huazhong University of Science and Technology (China), Southwestern Institute of Physics (China), Institute for Plasma Research (India), Eni S.p.A. (Italy), University of Cagliari (Italy), Osaka University (Japan), National Institute for Fusion Science (Japan), Korea Institute of Fusion Energy (Korea), Chung-Ang University (Korea), Centro de Investigaciones Energéticas, Medioambientales y Tecnológicas (Spain), Chalmers University of Technology (Sweden), Swiss Federal Institute of Technology Lausanne (Switzerland), Imperial College London (UK), Culham Centre for Fusion Energy – UK Atomic Energy Agency Culham Science Centre (UK), First Light Fusion Ltd. (UK), Massachusetts Institute of Technology (USA) (USA), Princeton Plasma Physics Laboratory (USA), University of Notre Dame (USA), General Atomics (USA), University of Wisconsin-Madison (USA), and College of William & Mary (USA).



Al for Fusion CRP and the Fusion Data Lake Project

"Al for Fusion: accelerating fusion R&D with Al, through the creation of a platform and cross-community network for innovation and partnership" - Al for Fusion (Homepage)

Who and what do we need for AI accelerated development of fusion energy?

The who:

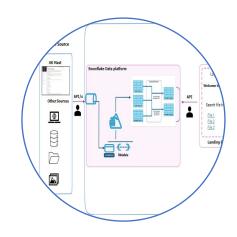
- Physicists
- Software developers
- Statisticians

The what:

- Surrogate models
- Digital twins
- Large scale data access!

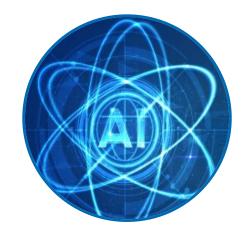
The Fusion Data Lake: Accelerating AI development through Open Science and FAIR Data

Today's Topics



Platform Design

Architecture design, data and metadata models, and technology and data transfers



Proof of Concept

What has been developed, and what is coming next



Data Governance

Data access and terms of service

Platform Design: Goal

The Fusion Data Lake has three core deliverables as its goals:

- 1. Catalogue Query-able, central catalogue of experimental fusion experiments (devices, shots, and diagnostic signals)
- **2. Data Federation** Decentralised network of experimental repositories accessible through a single point
- 3. Medium Term Storage Centralised storage for data projects

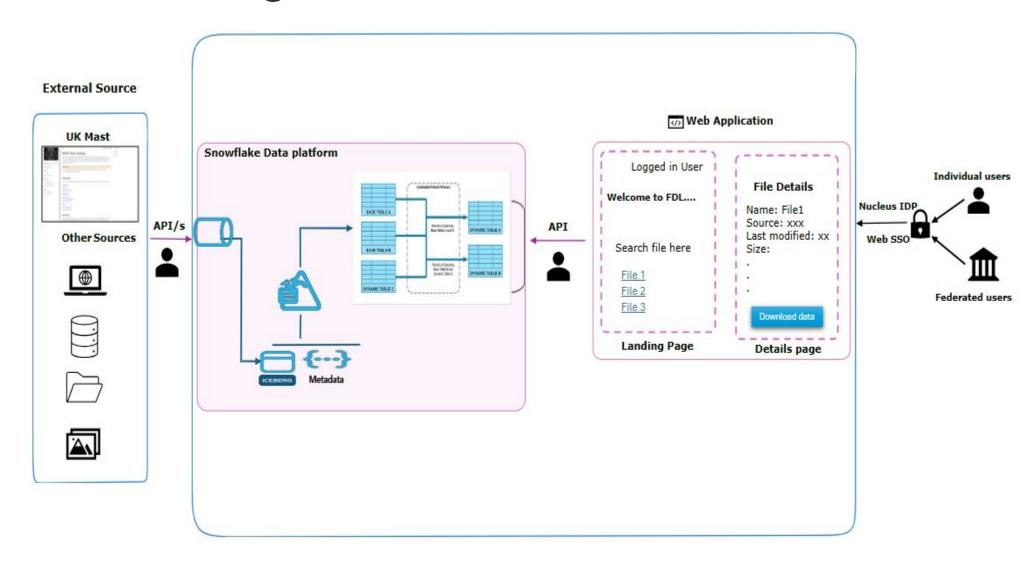
Platform Design: Goal

The Fusion Data Lake has three core deliverables as its goals:

- **1.** Catalogue Query-able, central catalogue of experimental fusion experiments (devices, shots, and diagnostic signals)
- **2. Data Federation** Decentralised network of experimental repositories accessible through a single point
- 3. Medium Term Storage Centralised storage for data projects

Building an inclusive fusion data platform for everyone!

Platform Design: Architecture



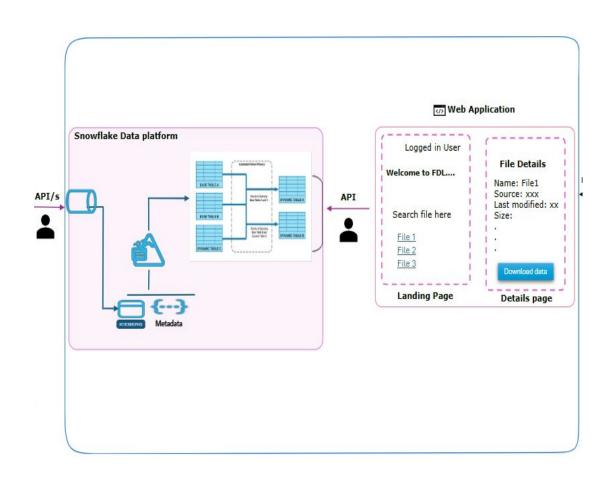
Platform Design: Architecture

Tech stack:

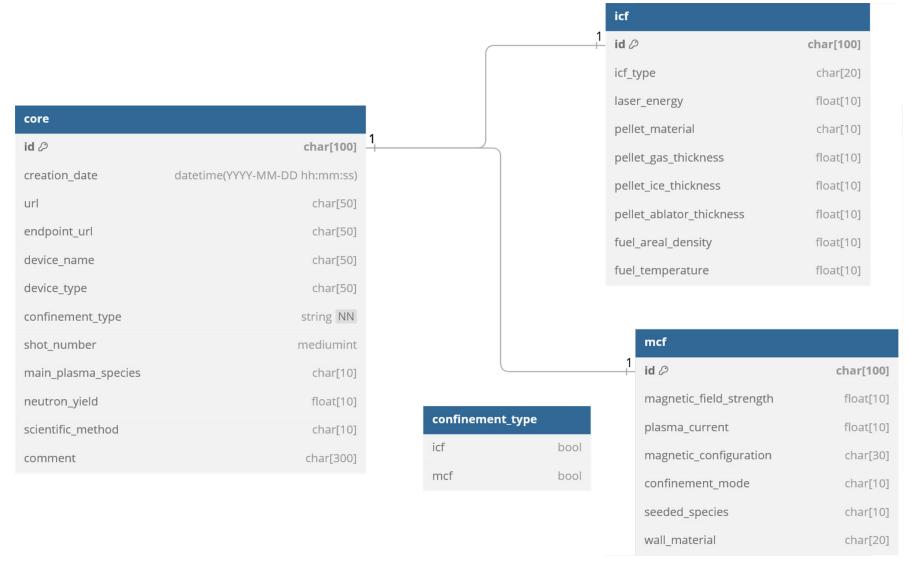
- Compute Snowflake
- Storage Azure
- Webapp C#/.Net

Data Engineering:

- Python (Snowflake native packages)
- Metadata driven ETL pipelines
 - O Ingestions configuration
 - Transformation configuration
- Medallion structure



Platform Design: Data and Metadata Model





View interactively here: <u>Fusion Data Lake - Minimal Metadata Model - dbdiagram.io</u>

Platform Design: Data & Metadata Model

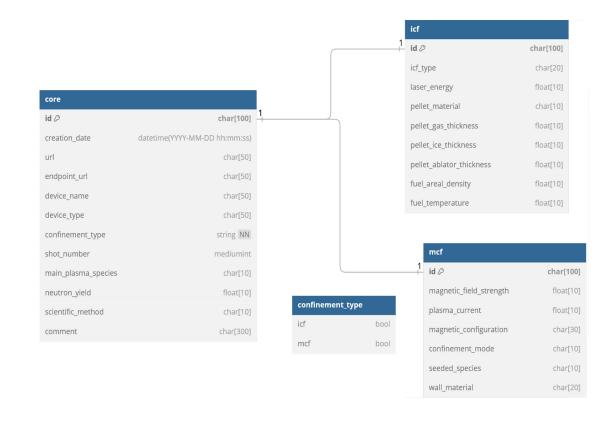
For a catalogue:

Minimum data set for querying (IFE/MFE)

In general:

- Ontologies from the "ITER Data Dictionary" and global standards such as "Dublin Core"
- Units and data scales
- Data types and limits

Storage: Dynamic allocation (machine and priendly)

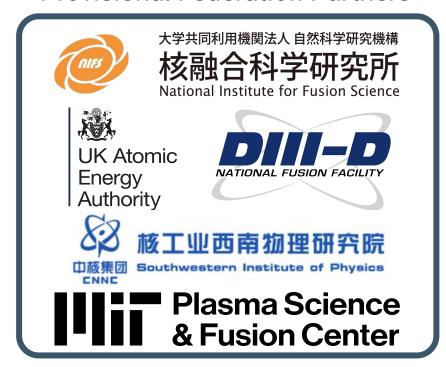


Platform Design: Technology and Data Transfers

Provisional Development Partners



Provisional Federation Partners

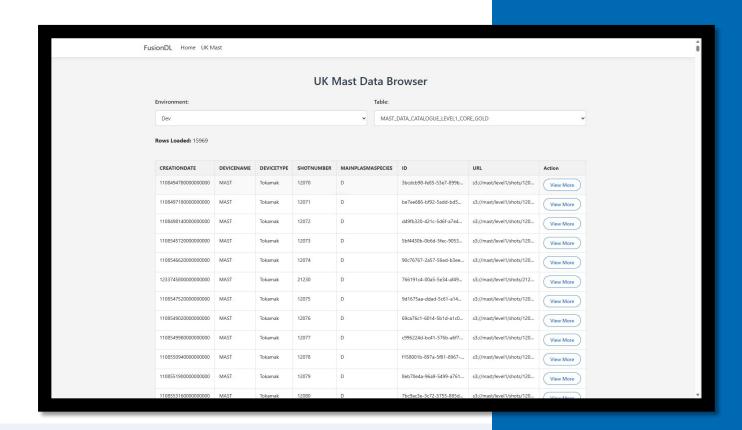


Do you want to share data and technology with the Fusion Data Lake project?

Proof of Concept: Phase 1 – Complete

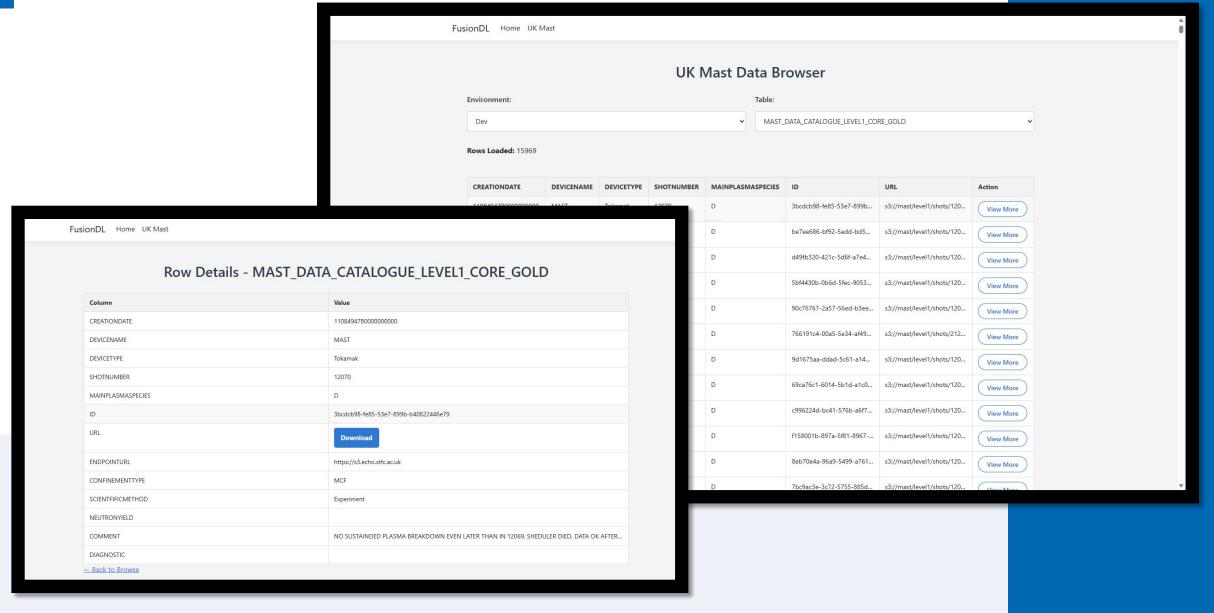
What was done:

- Data pipeline
- MAST Data Catalog connection
- Medallion data storage
- Webapp catalogue



Demonstrated cataloguing and federation capacity

Proof of Concept: Phase 1 – Complete



Proof of Concept: Phase 2 – In Progress

Goals: Expand the Fusion Data Lake with the LHD and Alcator C-Mod logs

<u>Timeline</u>: September 2025 to January 2026



Networking and Dataset ingestion



Data model mapping



E2E Testing

To demonstrate scalability!

Regulation: Terms of Service (ToS)



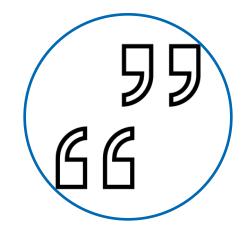
Uploads/Ingress

- Users who upload data must have the authority to share the data to the platform,
- Data must be licenced, either linked or in the metadata.



Downloads/Egress

- Data is only used and shared in line with ToS & data licensing,
- Data licenses will be transparent and accessible to enable compliance.



Attribution/Citation

- Data must be acknowledged in line with the data license,
- The platform must be referenced.

Regulation: Data Access Levels

Levels of data restriction:

- Public Open data. No login credentials required.
- Internal Open to anyone with login credentials to the platform.
- 3. <u>Restricted</u> Only open to individuals with login credentials from institutions approved by the data owner.
- 4. <u>Closed</u> Only be accessible to individuals approved by the data owner, with login credentials to the platform.

Institutional credential authentication using OpenAthens

You choose the level of access to your data!

Requires a

NUCLEUS (IAEA)

account

Conclusions

The IAEA's Fusion Data Lake is being built to accelerate AI development in fusion research through Open Science and FAIR Data.

Progress:

- Scale able engineering design within standard fusion ontologies
- Proof of Concept Phase 1 to Phase 2
- Developing regulation for an inclusive and safe platform







Come support the development and adoption of the Fusion Data Lake!



Thank you for listening!

Any Questions?

Daljeet Singh Gahle *et al*, IAEA 6th IAEA Technical Meeting on Fusion Data Processing, Validation and Analysis September 9th-12th, Fudan University, Shanghai

Conclusions

The IAEA's Fusion Data Lake is being built to accelerate AI development in fusion research through Open Science and FAIR Data.

Progress:

- Scale able engineering design within standard fusion ontologies
- Proof of Concept Phase 1 to Phase 2
- Developing regulation for an inclusive and safe platform







Come support the development and adoption of the Fusion Data Lake!