

Enhancing Fusion Research with TokSearch: Updates and Integration into the Fusion Data Platform

presented by **Brian Sammuli, General Atomics**

Brian Sammuli¹, Erik Olofsson¹, Tom Neiser¹, Mitchell Clark¹, David Orozco¹, Cihan Akcay¹, Javier Hernandez Nicolau², Fabio Andrijauskas², Annmary Justine Koomthanam³, Aalap Tripathy³, Rishabh Sharma³, Matthew Waller⁴, Ruqi Pei⁴, Zeyu Li¹, Amitava Majumdar², Rose Yu², Sicun Gao², Frank Wuerthwein², Raffi Nazikian¹, Martin Foltin³, Craig Michoski⁴, David Schissel¹

¹General Atomics

²University of California, San Diego

³Hewlett Packard Enterprise

⁴Sapientai

14th IAEA Technical Meeting on Control Systems, Data Acquisition, Data Management and Remote Participation in Fusion Research

July 17, 2024

Work supported by the U.S. Department of Energy under Award No. DE-FC02-04ER54698 and Award No. DE-SC0024426

Overview

- **The Fusion Data Platform project**
- **TokSearch**
- **Project Timeline**

Motivation: Why do we need a Fusion Data Platform?

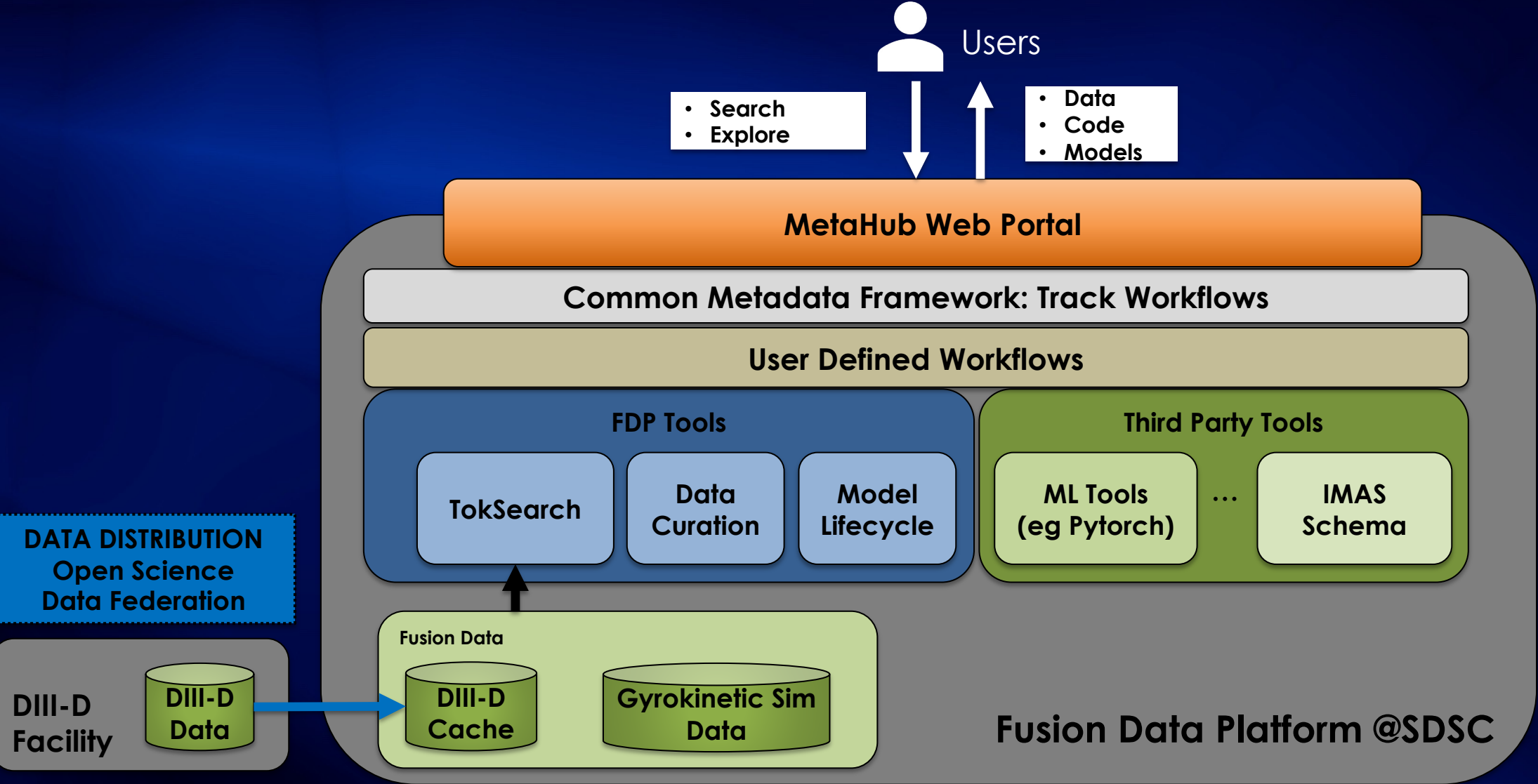
	Need	FDP Solution
Data Distribution	No community-wide broad distribution of fusion data	Use Open Science Data Federation software stack (HEP tested)
Data Curation	No standard tools for curating fusion data	Equip FDP with fusion-specific data curation tools
Standardization	MDSplus commonly used, but no unified schema	Use IMAS schema for curated data
Reproducibility and Provenance	ML models often published with insufficient info to reproduce results	Use HPE's Common Metadata Framework for workflows, with complete provenance tracking of data, models, and metadata
Discoverability	No standard repository of curated data, no ability to search existing models/data	Provide metadata database with web portal and search interfaces

The Fusion Data Platform (FDP) project aims to address these challenges

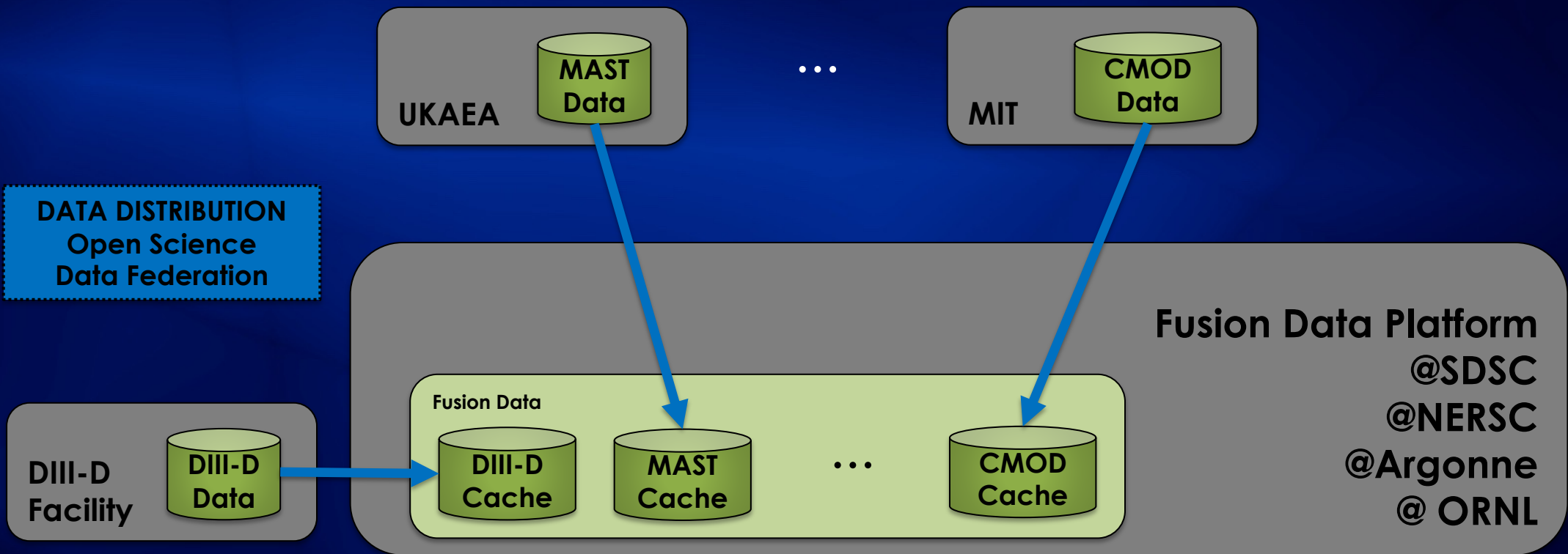
- Development environment for fusion data-driven applications with fusion-specific data processing and curation tools
- Uses distributed version control semantics for code, metadata, generated artifacts
- Provides federated access to data
- Team:
 - General Atomics (PI: Brian Sammuli)
 - Hewlett Packard Enterprise (PI: Martin Foltin)
 - UCSD/SDSC (PI: Frank Wuerthwein)
 - Sapien^α (PI: Craig Michoski)
- **3 year, \$7.4M project, FES funded**



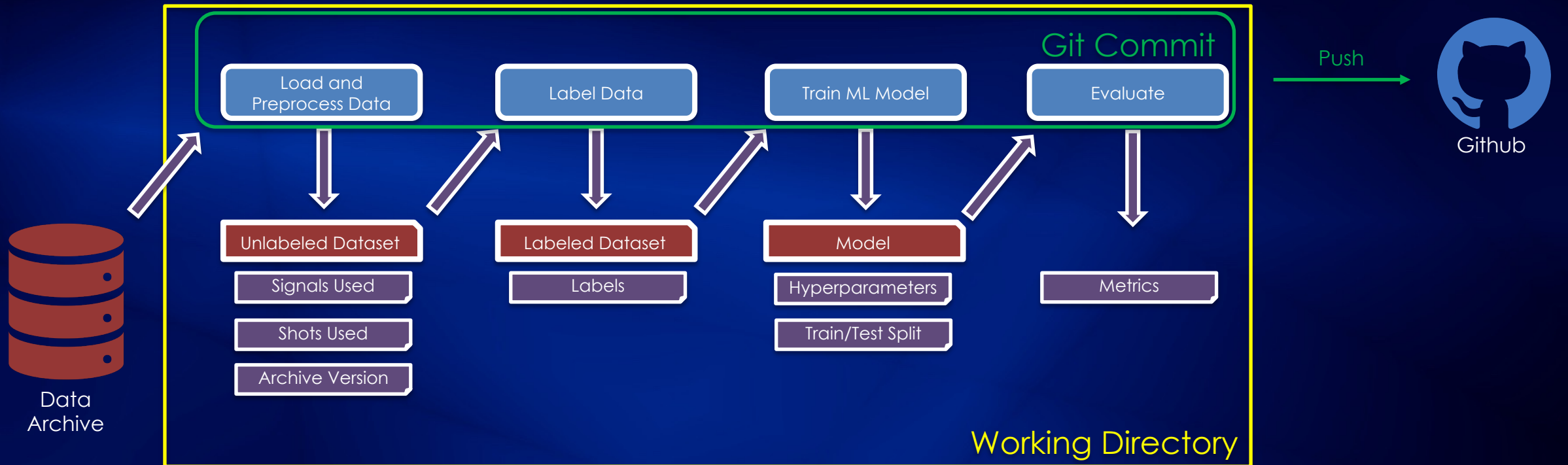
Fusion Data Platform will provide comprehensive environment for community-driven AI/ML development, supporting both experimental and simulation data



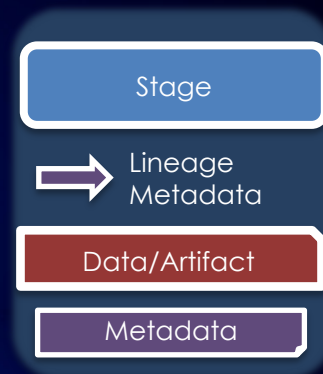
Aim is to expand platform to include additional devices, and deploy tools across multiple computing sites



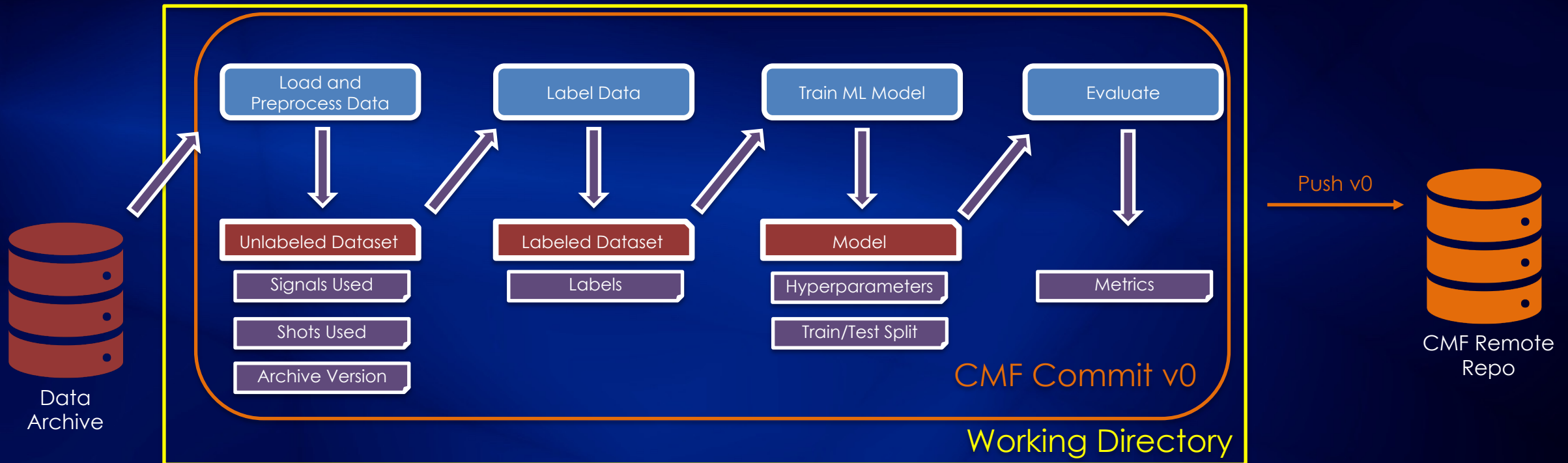
Workflows can be represented as directed graphs that capture the relationship between code, data, generated artifacts, and metadata



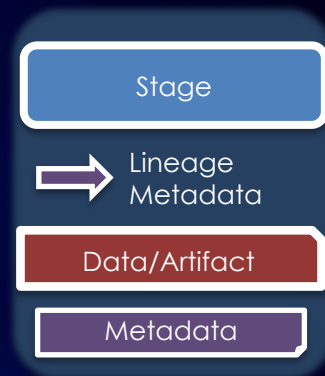
- **Typical Usage Pattern – Just version the code**
 - Code and model can diverge



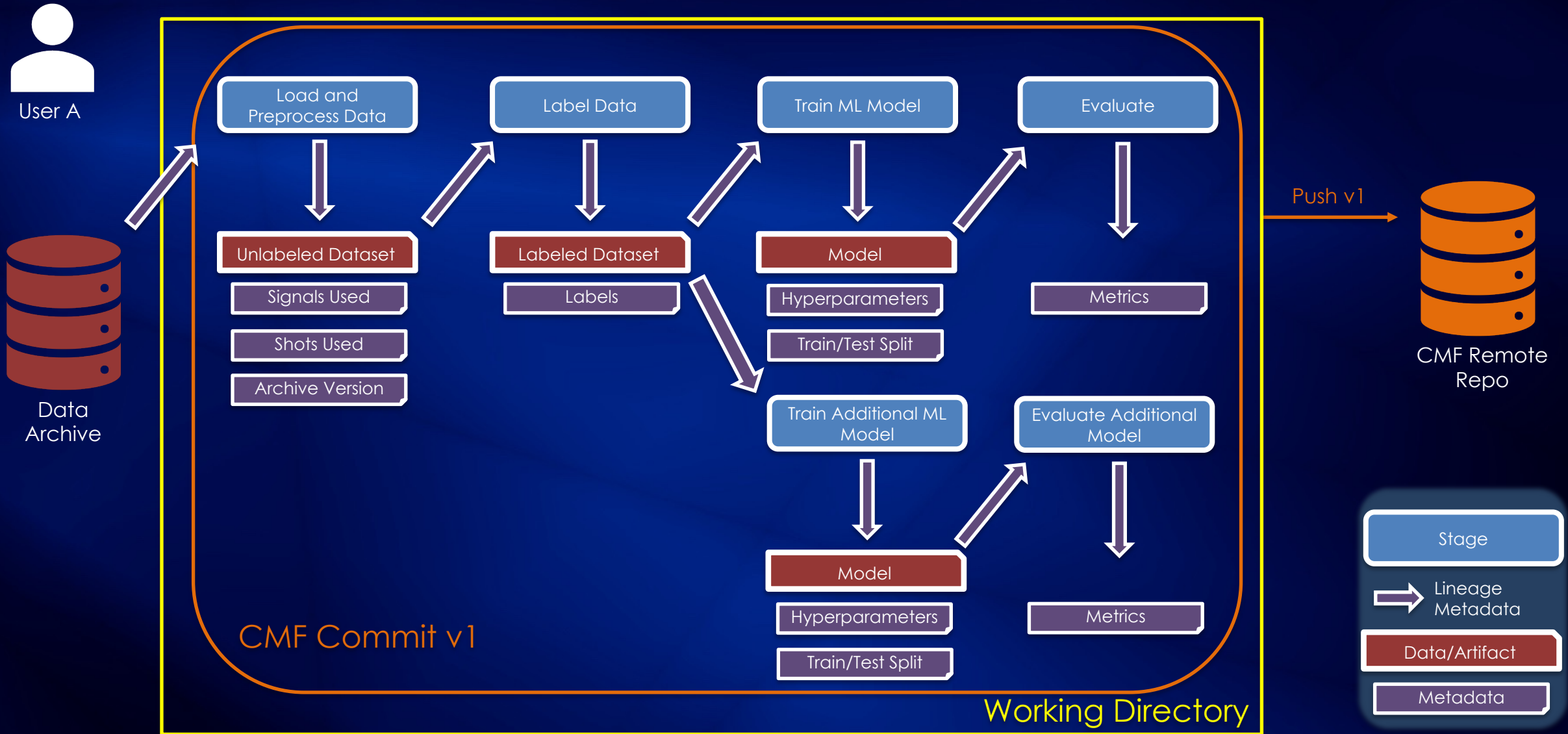
Workflows can be represented as directed graphs that capture the relationship between code, data, generated artifacts, and metadata



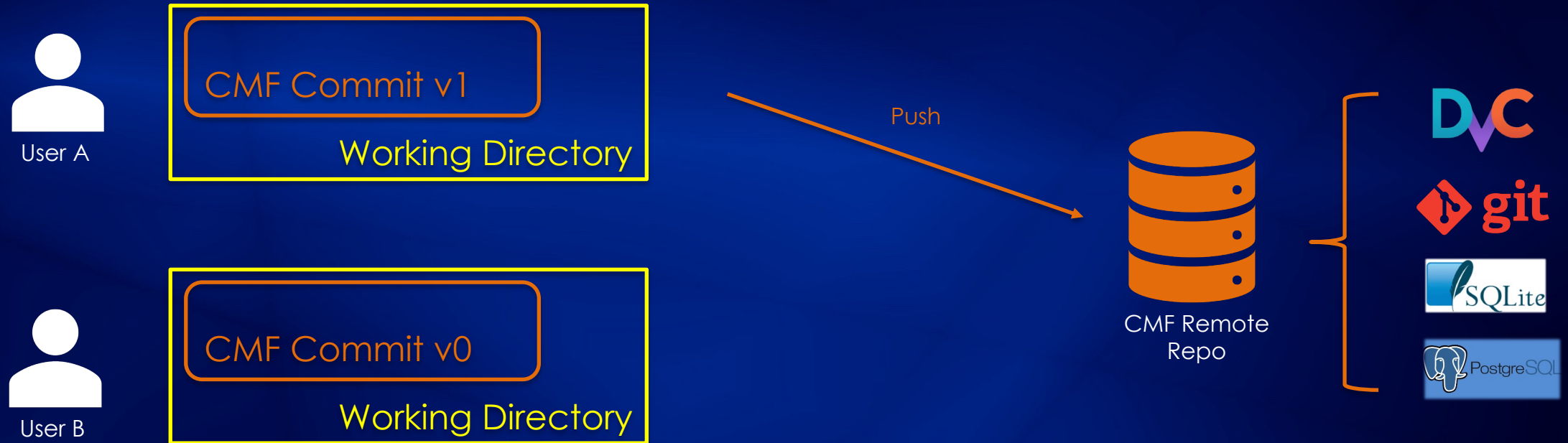
- **HPE's Common Metadata Framework (CMF) bundles entire workflow in a single atomic unit, allowing for consistent snapshots**
 - Enables provenance tracking, reproducibility, and sharing



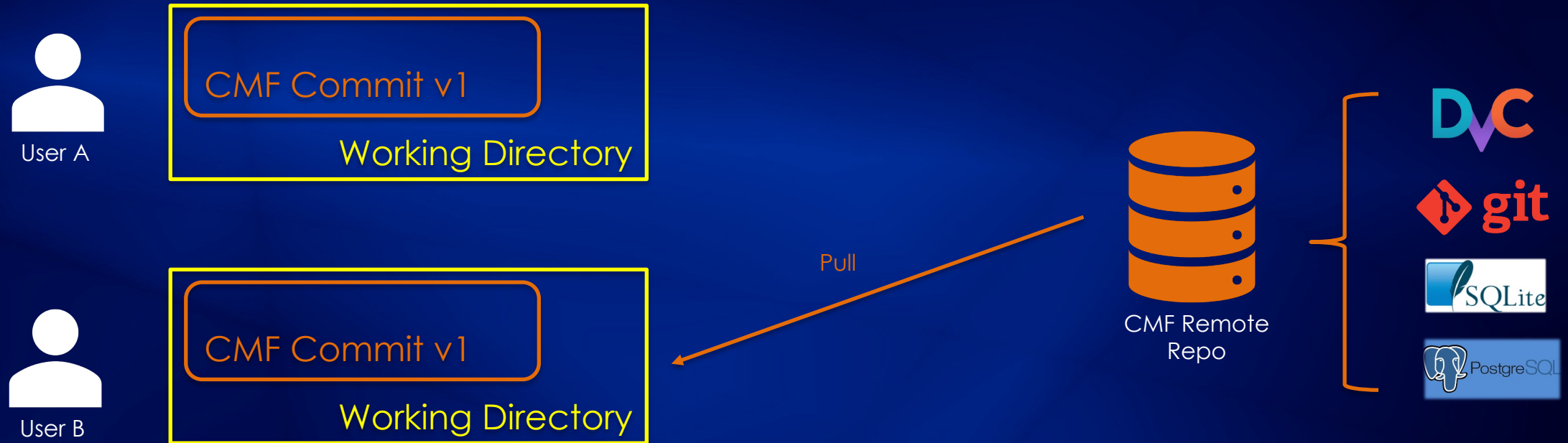
The FDP uses distributed version control semantics, allowing for collaborative development and workflow sharing



The FDP uses distributed version control semantics, allowing for collaborative development and workflow sharing



The FDP uses distributed version control semantics, allowing for collaborative development and workflow sharing



Hewlett Packard Enterprise's Common Metadata Framework (CMF) orchestrates coordinated version control of code, metadata, generated artifacts

- CMF supports definition of arbitrarily complex workflows and integrates with popular AI/ML tools

Simplified active learning workflow definition:

```

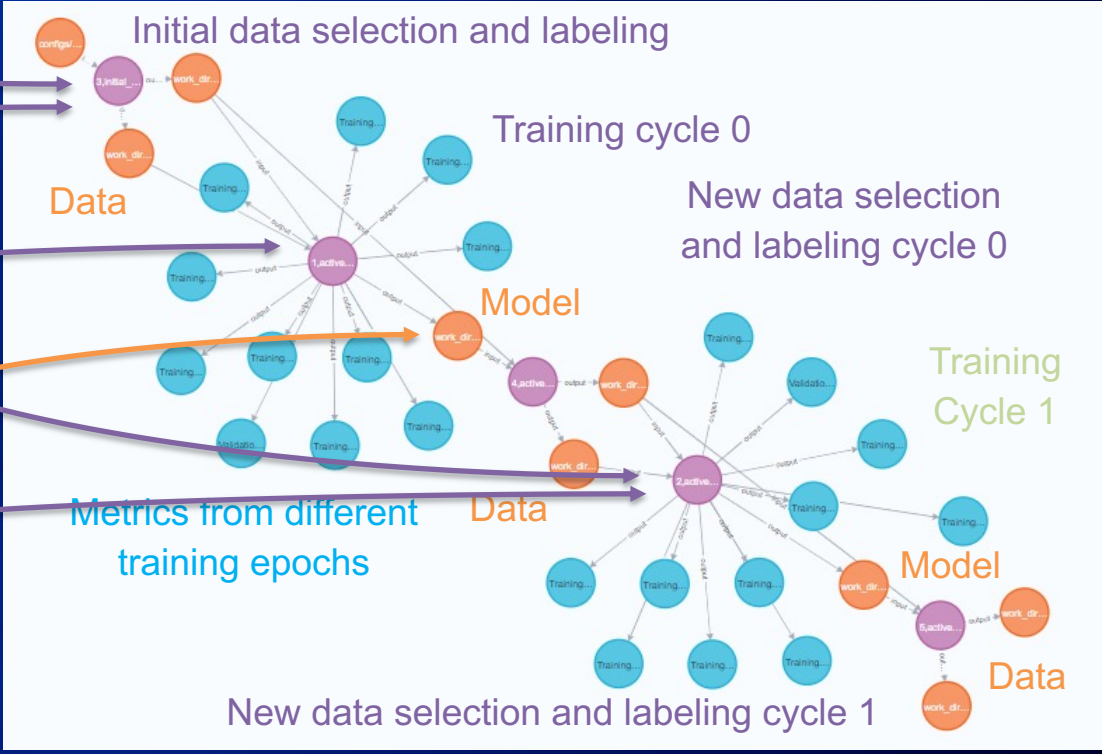
initial_data_selection:
  cmd: python initial_select.py input_data selected_data_0 ...
  deps: - input_data
  outs: - selected_data_0

labeling:
  foreach: - cycle: 0 - cycle: 1
  do:
    cmd: python label.py selected_data_${item.cycle} labeled_data_${item.cycle} ...
    deps: - selected_data_${item.cycle}
    outs: - labeled_data_${item.cycle}

training:
  foreach: - cycle: 0 - cycle: 1
  do:
    cmd: python train.py labeled_data_${item.cycle} ai_model_${item.cycle} ...
    deps: - labeled_data_${item.cycle}
    outs: - ai_model_${item.cycle}

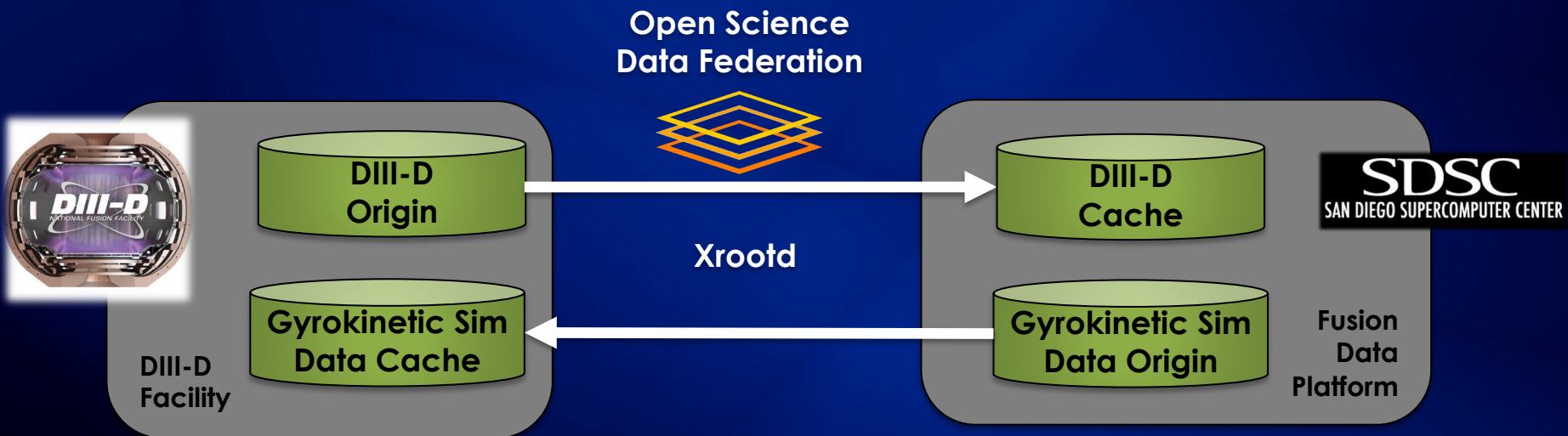
new_data_selection:
  foreach: - cycle: 0 next_cycle: 1 - cycle: 1 next_cycle: 2
  do:
    cmd: python select.py input_data selected_data_${item.cycle}
        ai_model_${item.cycle} selected_data_${item.next_cycle} ...
    deps: - input_data, selected_data_${item.cycle}, ai_model_${item.cycle}
    outs: - selected_data_${item.next_cycle}
  
```

Corresponding CMF lineage
(Labeling & Selection merged for simplicity):

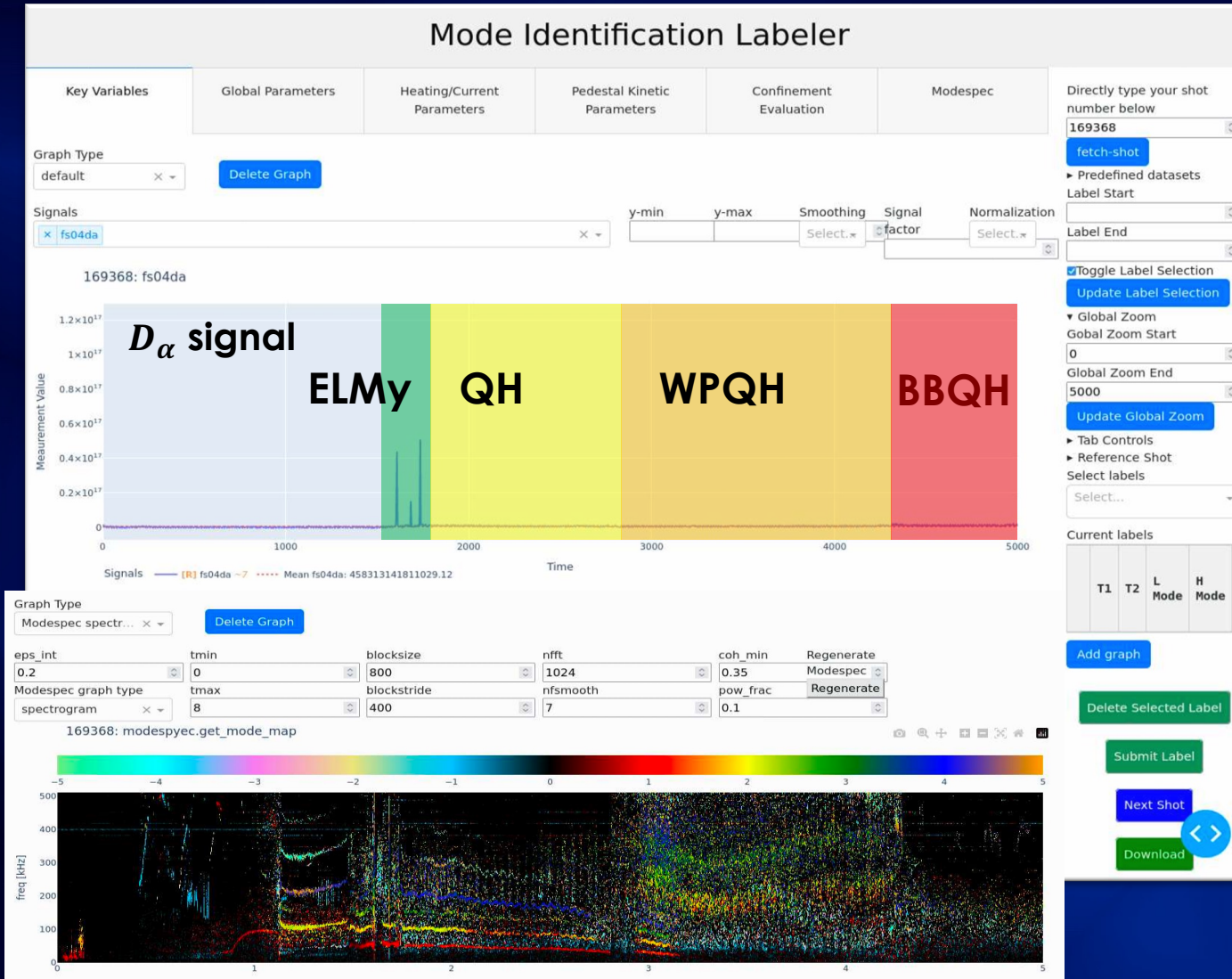


Federated access to data achieved using Open Science Data Federation (OSDF)

- OSDF uses origin/cache scheme
 - Origin: Original source of data
 - Cache: Access data at remote site on fast storage
- DIII-D origin/cache operational
- Gyrokinetic sim. origin/cache in progress



Sapientai has created a new visual labeling tool: First applied to edge plasma regimes at DIII-D



Labels:

- Edge-localized-mode (ELMy)
- Standard quiescent-H mode (QH)
- Wide pedestal quiescent-H mode (WPQH)
- Broadband quiescent-H mode (BBQH)

	shots	Time Spent	Speed up
Previous database	145	~2-3 weeks	-
New Database	400	1 week	~5-8x

Sapientai is developing fusion-specific data exploration tools and ML-assisted labeling that will be integrated with visual labeling capability

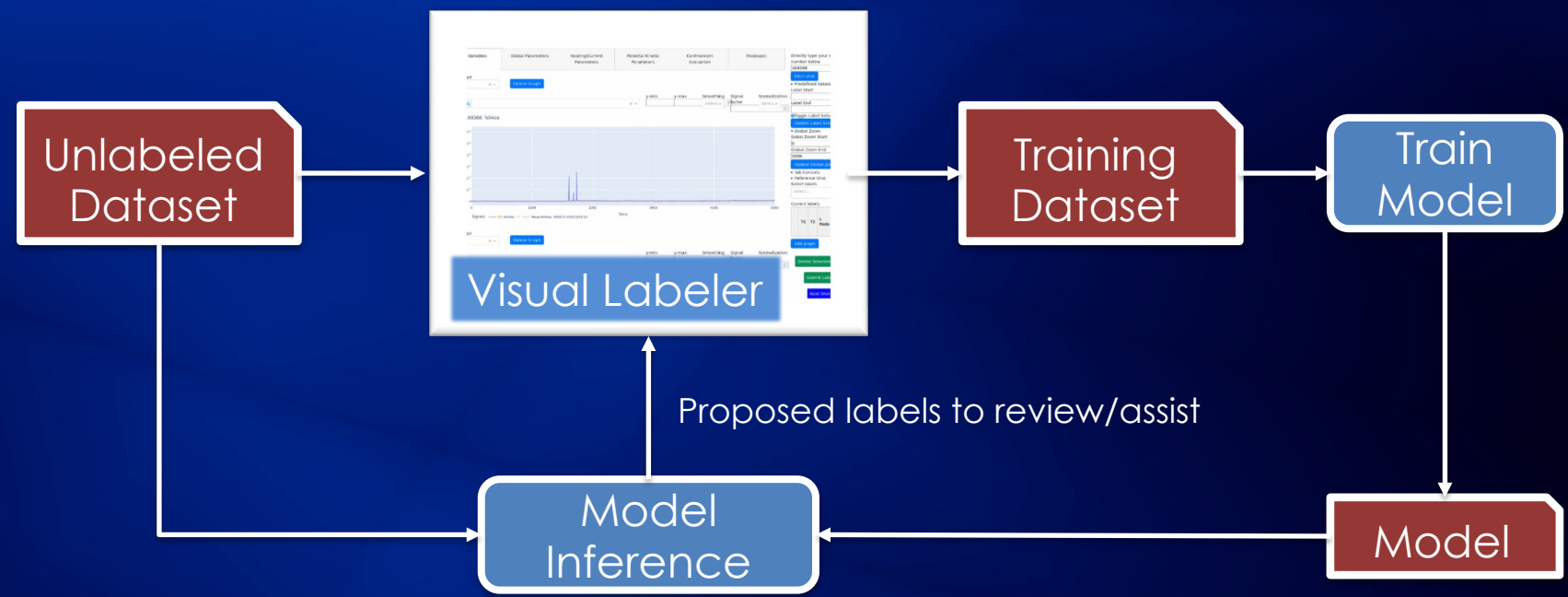
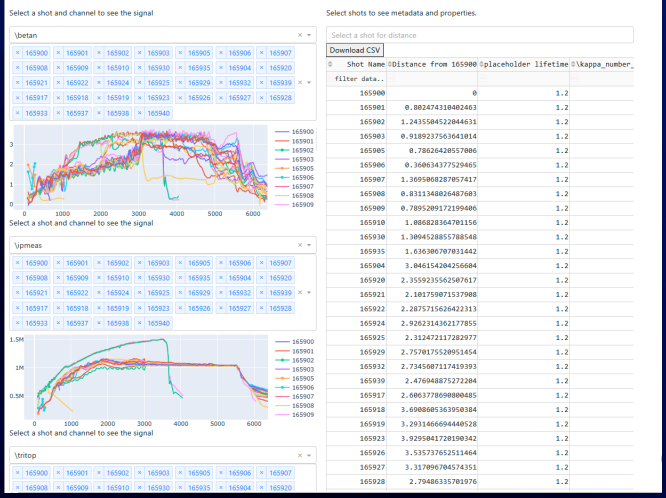


Unsupervised clustering helps inform labeling



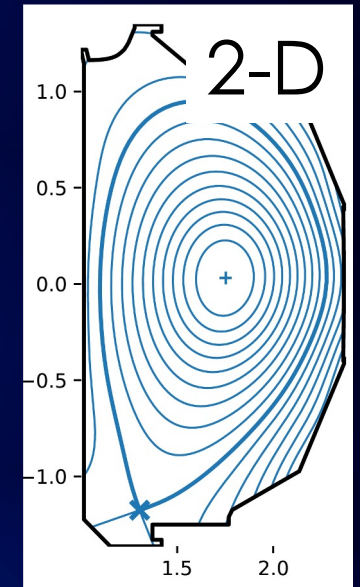
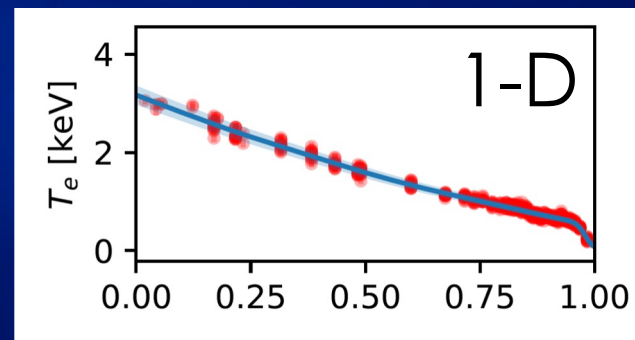
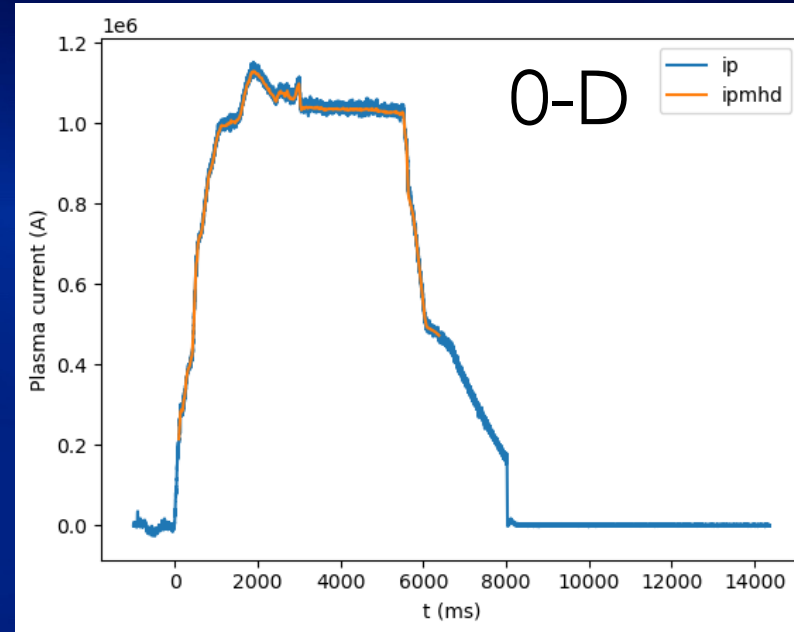
ML-assisted labeling shows proposed labels and allows user to not start from scratch

Historical context of many shots guides curation process



Constructing AI/ML datasets from fusion data is a challenge!

- Large volume of data
- Large variability of data types
 - Sample rates span orders of magnitude: $O(1 \text{ Hz}) - O(100 \text{ Mz})$
- Variety of dimensionalities
 - 0-D scalar time series
 - e.g. magnetics, currents
 - 1-D profile time series
 - e.g. temperature profile
 - 2-D grid data time series
 - e.g. equilibrium reconstructions
 - Image time series
 - e.g. infrared camera data



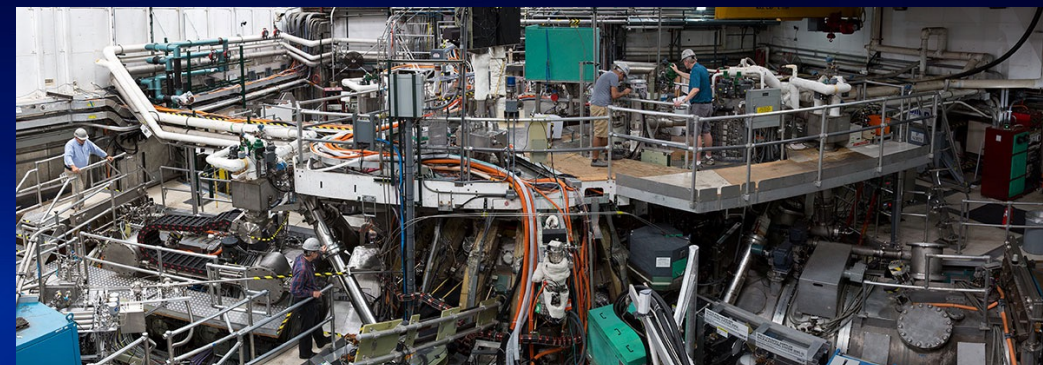
DIII-D data is stored as in multiple formats, adding to data curation difficulty

Total data volume: 700 TB

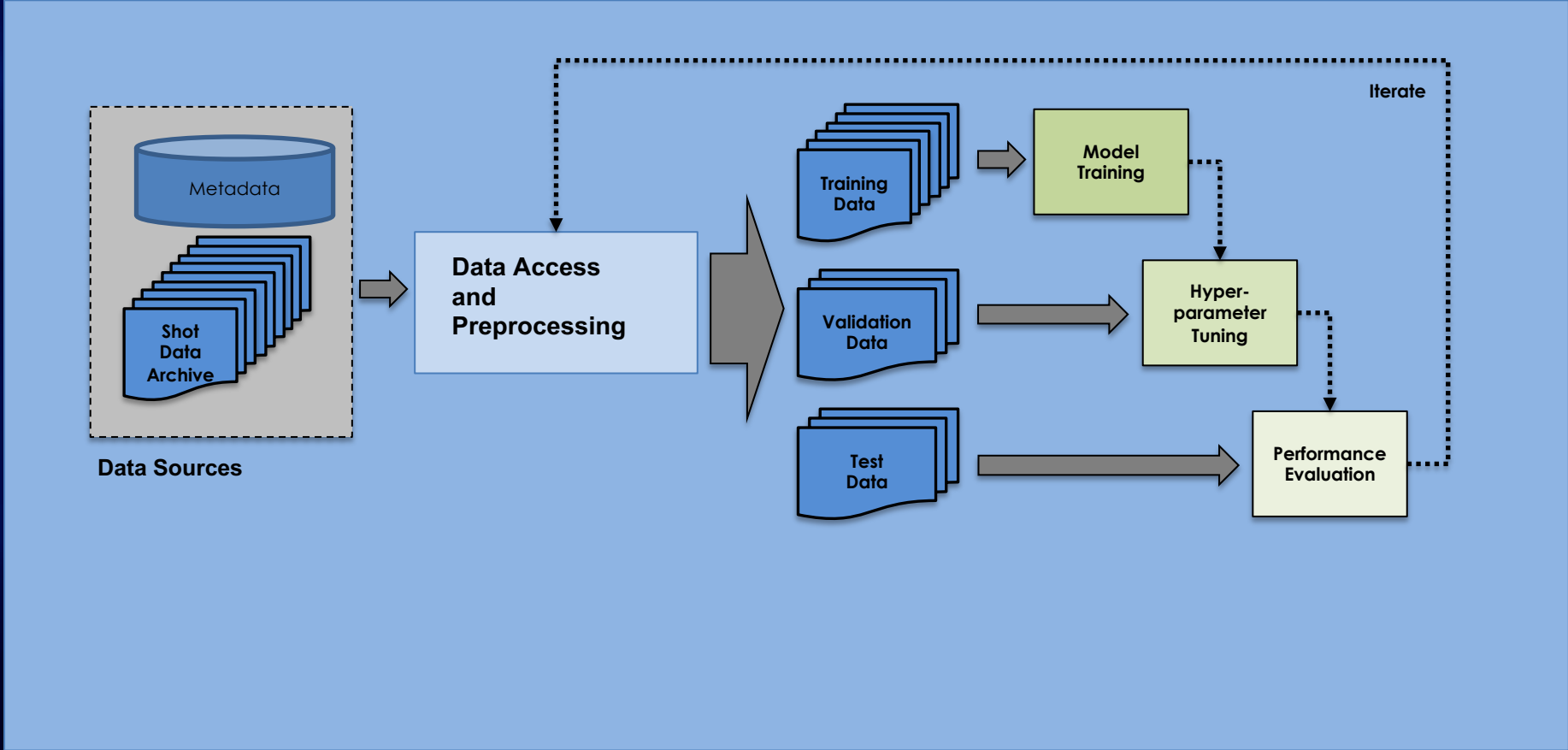
Data Archive Types:

- **PTDATA**
 - Raw data from data acquisition systems (ie digitizers)
 - ~90% of total data
- **MDSplus**
 - Store output of analysis code (e.g. plasma equilibria)
 - ~10% of total data
- **Relational Database**
 - Shot metadata (e.g. date, operator notes)
 - \ll 1% of total data

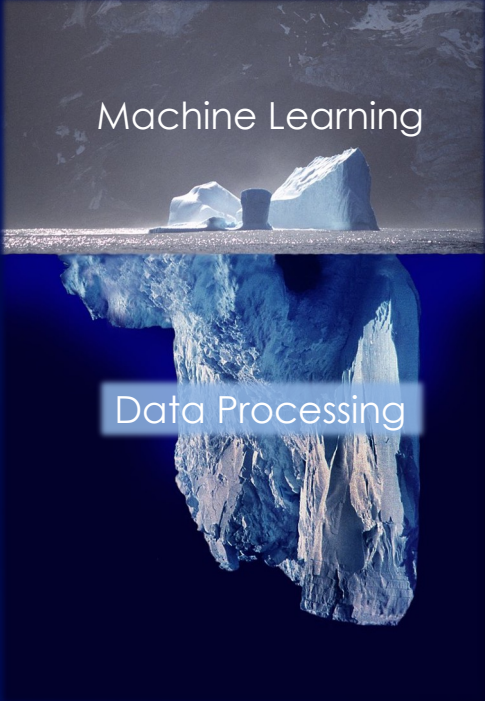
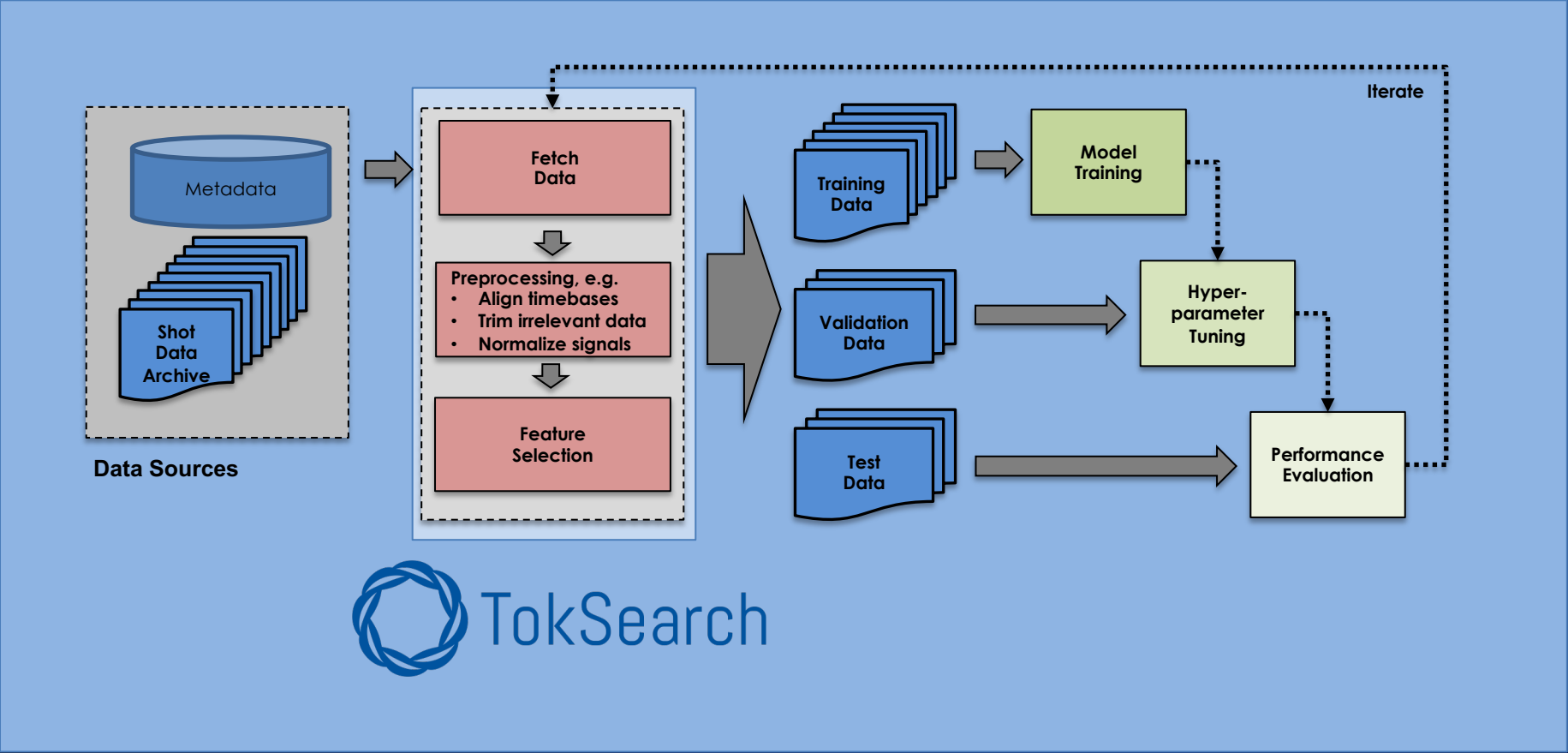
→ Typically use all three when doing AI/ML



Creating real world datasets involves significant data retrieval and preprocessing



Creating real world datasets involves significant data retrieval and preprocessing



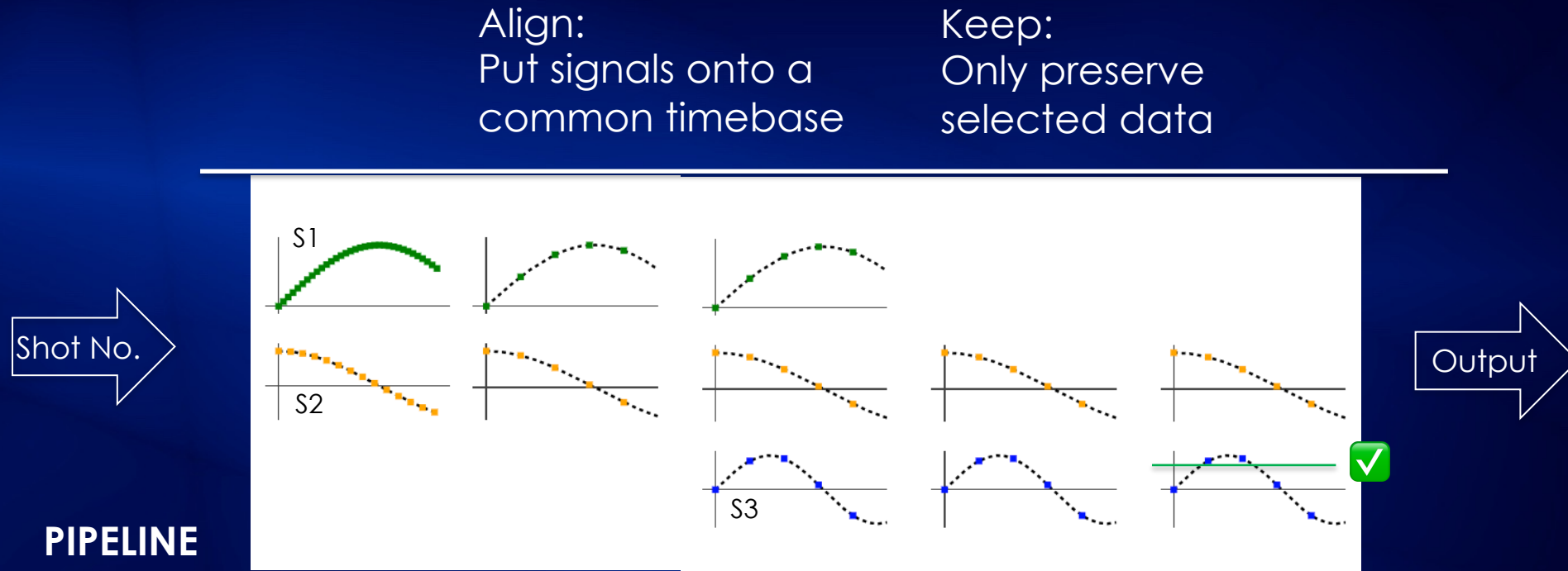
Quick poll: How many of you have written something like this?

```
results = []  
for shot in shots:  
    data = read_data(shot)  
    processed_data = process_data(data)  
    if meets_some_criteria(processed_data):  
        results.append(processed_data)
```



TokSearch accelerates this type of processing pipeline!

TokSearch composes built-in and user-defined functions into a pipeline applied to each shot



Align:
Put signals onto a
common timebase

Keep:
Only preserve
selected data

PIPELINE

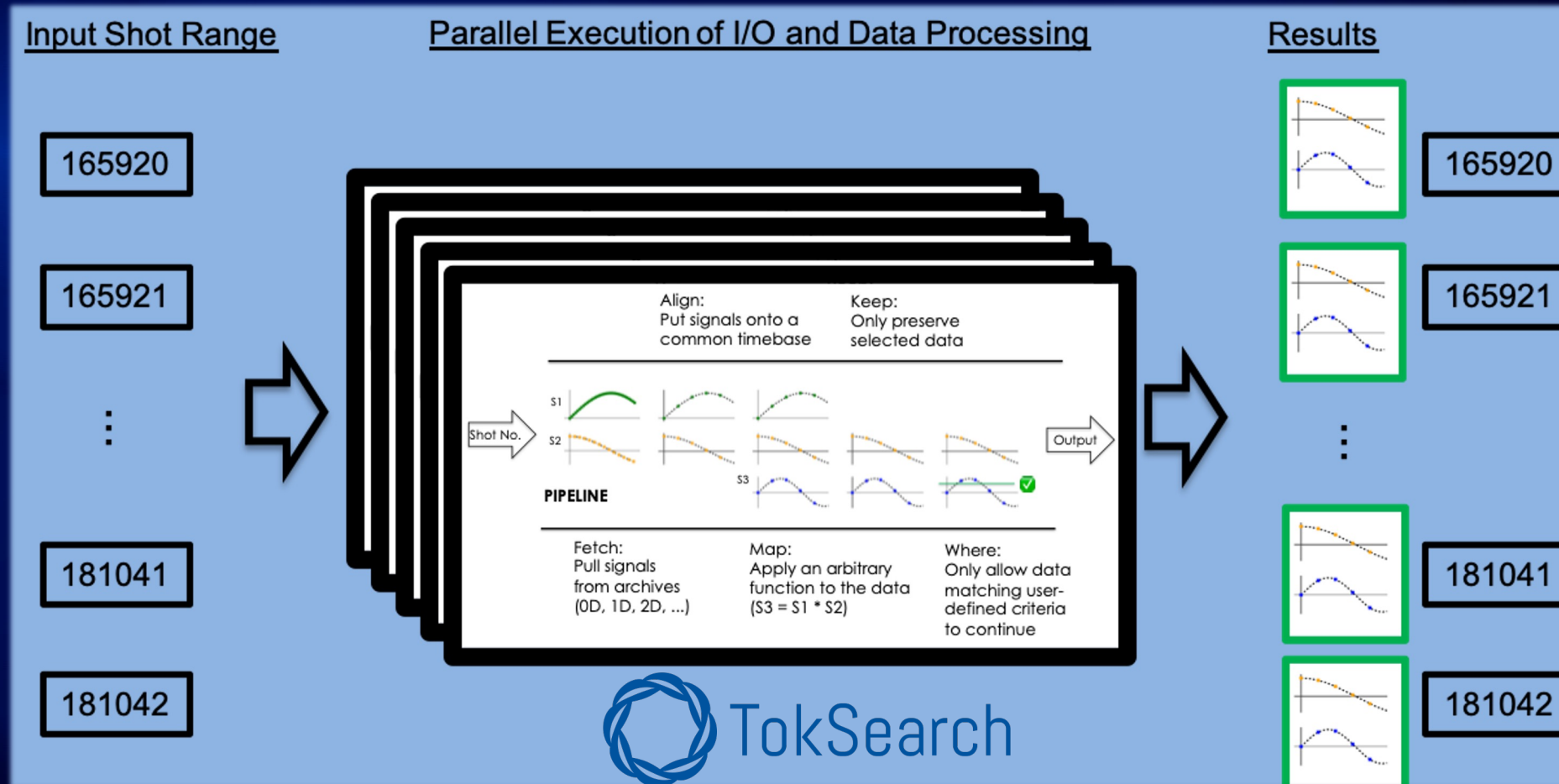
Fetch:
Pull signals
from archives
(0D, 1D, 2D, ...)

Map:
Apply an arbitrary
function to the data
($S3 = S1 * S2$)

Where:
Only allow data
matching user-
defined criteria
to continue



TokSearch pipelines can be executed in parallel for high throughput data processing



- Access to ~1PB of DIII-D data at high throughput

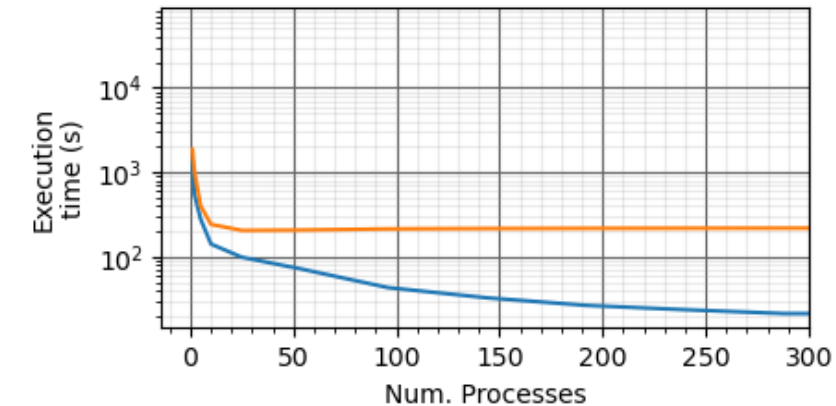
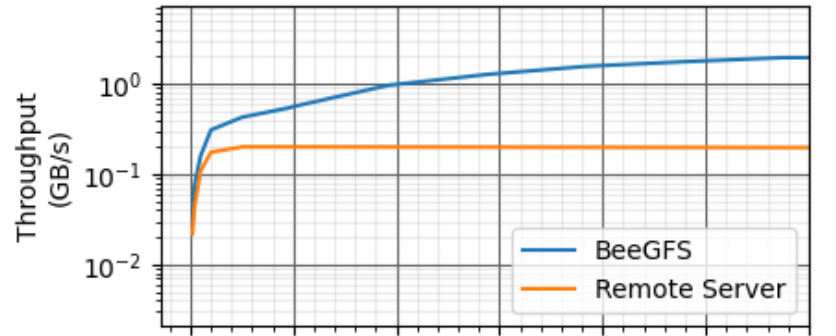
TokSearch can exploit distributed computing resources

- **Supported methods of parallelization: Ray, Apache Spark, Python Multiprocessing**
 - Others possible (MPI, Dask,...). Very easy to extend.
 - Runnable on HPC systems via SLURM
- **Full copy of DIID-D archives (0.7 PB) available on BeeGFS fast file system**
 - Currently available on Saga cluster hosted at GA
 - Will be available via OSDF on FDP

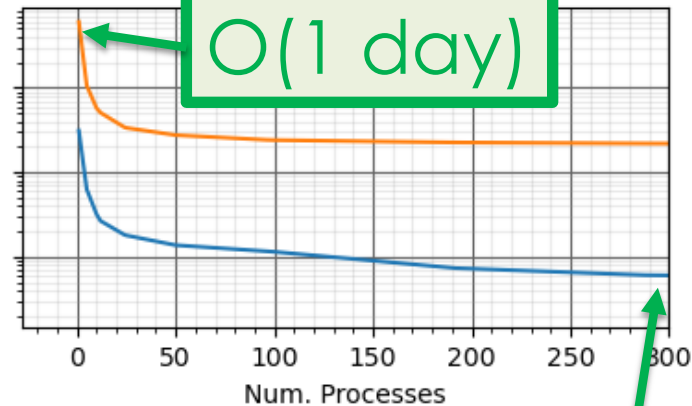
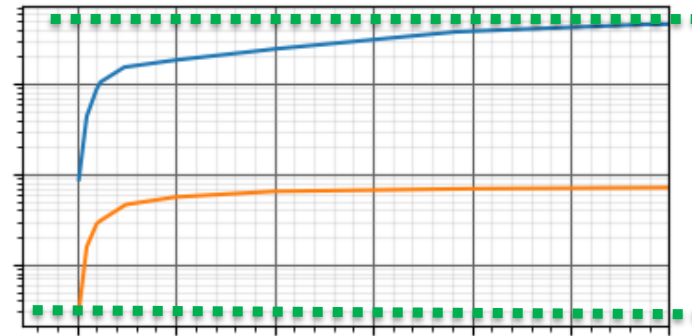


TokSearch Performance Benchmarks: Orders of magnitude speedup via parallelization

MDSplus
Reading 65x65 Grid Data



PTDATA
Reading 10x 20 kHz scalar signals



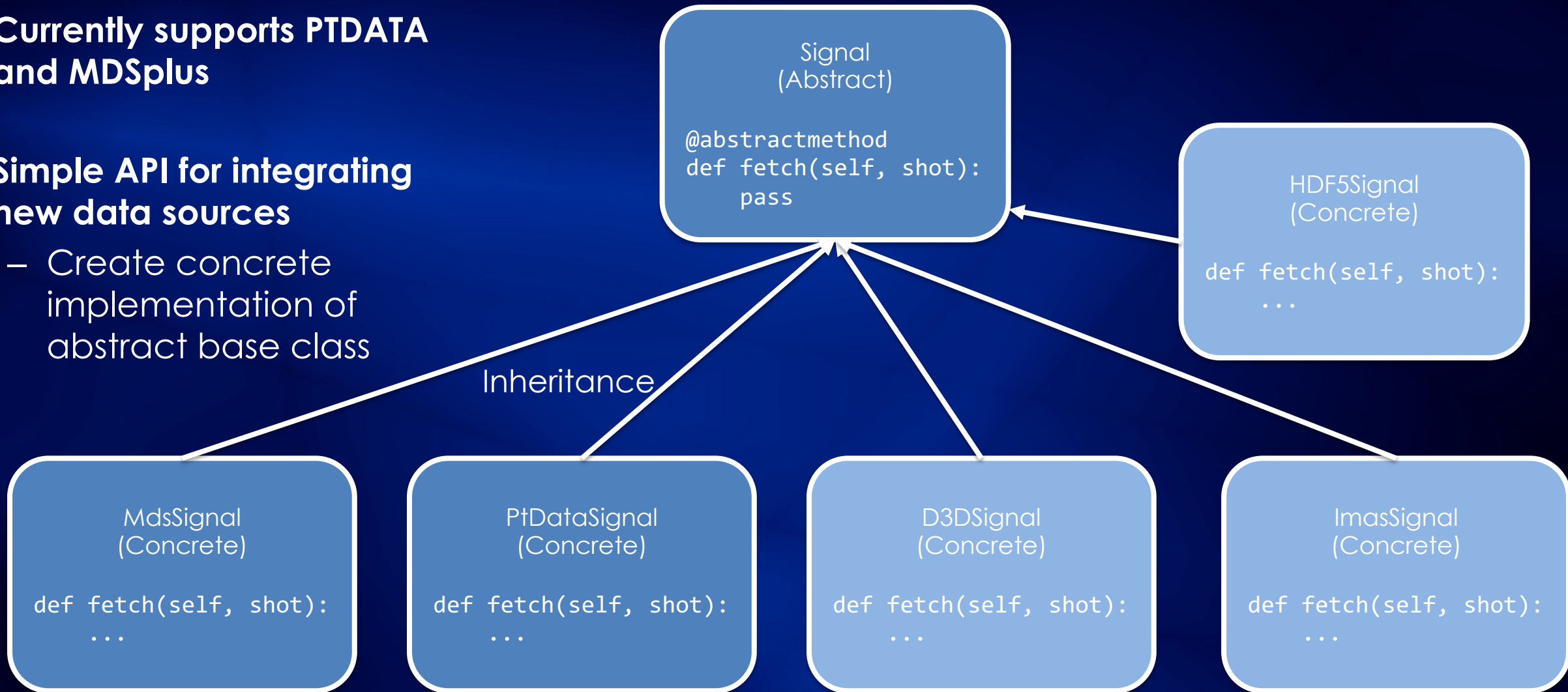
1700x Speedup

- Load and process 10,000 shots in minutes instead of days

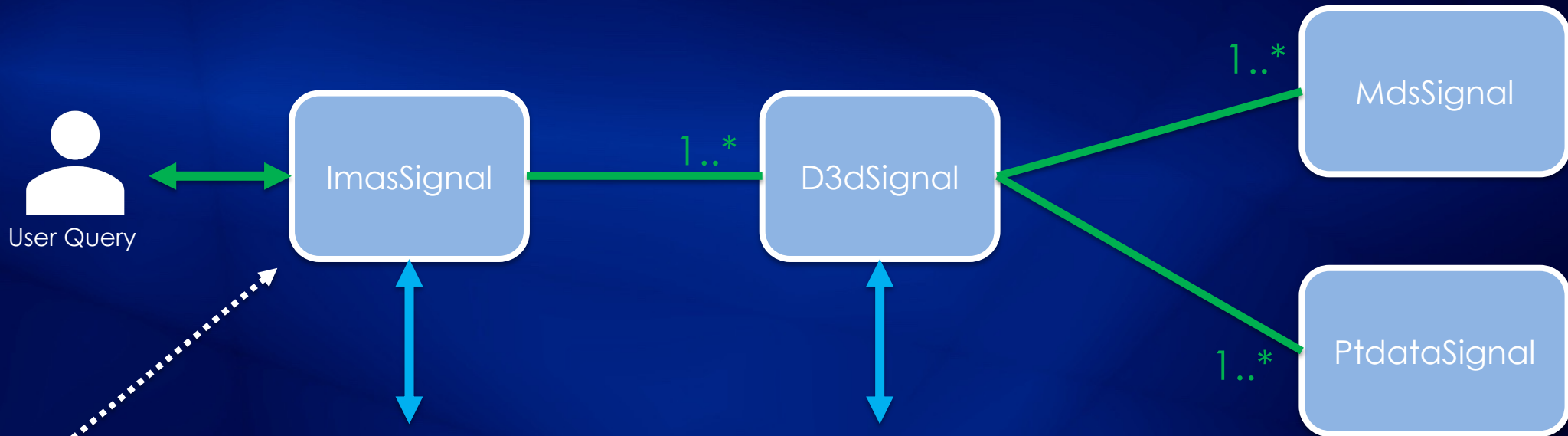
O(1 minute)

TokSearch data access is designed to be extensible

- **Currently supports PTDATA and MDSplus**
- **Simple API for integrating new data sources**
 - Create concrete implementation of abstract base class



Mapping of TokSearch to IMAS Schema is Underway



- Get metadata needed to resolve IMAS entry
- Assemble compound entries
- Perform COCOS/UNIT conversions

TokSearch recently open-sourced under Apache 2.0 license

- Source available on GitHub (<https://github.com/GA-FDP/toksearch>)
- Docs site has tutorials, extensive API documentation

TokSearch 2.0.X

GA-FDP/toksearch
☆0 🗨0

TokSearch

Tutorials

- The Basics
- Working with signals
- Xarray and signal alignment
- Parallelization
- Distributed computing
- Creating pipelines from SQL
- Combining data after pipeline computation

API

- Pipeline
- Record
- MDSplus Data Access
- Backend
- Abstract Interfaces

TokSearch

Welcome to TokSearch

TokSearch is a Python package for parallel retrieving, processing, and filtering of arbitrary-dimension fusion experimental data. TokSearch provides a high level API for extracting information from many shots, along with useful classes for low level data retrieval and manipulation.

The fundamental class in TokSearch is the `Pipeline`. A `Pipeline` object takes a list of shots and, for each shot in the list, creates a dict-like object called a `Record`. The `Pipeline` object then provides methods for defining a sequence of processing steps to apply to each record. These processing steps include:

- Passing user-defined functions to the pipeline via the `map` method.

or...

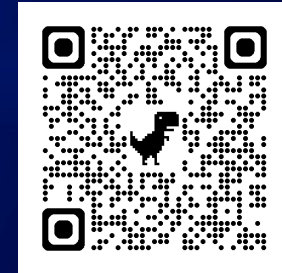
- Using a set of built-in methods, such as `fetch`, `fetch_dataset`, `align`, `keep`, or `discard`.

Table of contents

Installation

- Installation with Conda in an existing environment
- Installation with Conda in a new environment
- Installation from Source

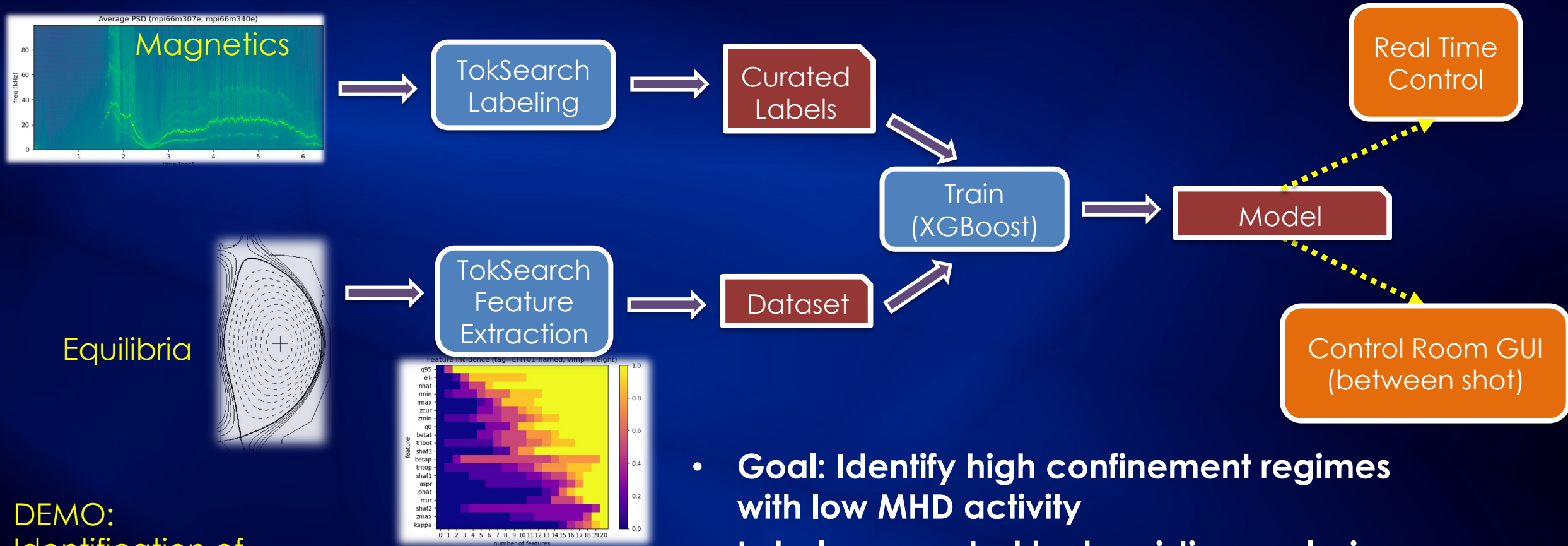
Source



Docs



A workflow that applies hazard analysis to neoclassical tearing modes has been completed using magnetic + MSE equilibria, will be extended to kinetic equilibria



DEMO:
Identification of
safe operating
regimes

- Goal: Identify high confinement regimes with low MHD activity
- Labels generated by heuristics analyzing locked and rotating modes in TokSearch

Hazard Analysis using FDP tools near ready for control room deployment

- Control room tool allows experimentalists to examine causes of tearing mode precursors from previous shot, adjust subsequent control setup
- Developed collaboratively with LLNL colleagues (Holcomb, Victor)



Project Timeline

	Year 1 (FY 2024)	Year 2 (FY 2025)	Year 3 (FY 2026)
Platform Infrastructure and Tools	<p>Initial FDP deployed and ready for use:</p> <ul style="list-style-type: none"> • Core infrastructure • Initial curation tools (Visual labeling, TokSearch) • Documentation • DIII-D + gyrokinetic data 	<p>MetaHub portal deployed</p> <p>Model lifecycle management tools</p> <ul style="list-style-type: none"> • Continual learning • Drift detection <p>TokSearch LLM integration</p>	<p>Metadata indexing capability added to MetaHub</p> <p>Data discovery and workflow recommendation tools</p>
Community engagement	<p>Announcement of alpha release at APS-DPP</p>	<p>Broad beta release</p> <p>Web portal</p> <p>DIII-D team engagement</p> <p>Student engagement</p> <p>Training, outreach</p>	<p>Full public release</p> <p>Online tutorials</p> <p>User community board</p>
Demonstrations	<p>Demos initiated, documented</p> <p>Alpha users engaged</p>	<p>Migrate development effort to platform at SDSC</p> <p>Publish Demos and use cases on platform</p>	<p>Integrate model lifecycle tools, incorporate in CI/CD pipelines</p> <p>Publish workflows at project completion</p>

Wrapping up

- **FDP project is underway, and is being demonstrated with DIII-D data**
- **Provides flexible environment for developing data-driven workflows**
 - Modeling
 - Data curation (labeling, tagging, cleaning)
 - Simulation
- **Distributed version control semantics and federated data access allows easy deployment at multiple computing centers (e.g. DOE, cloud)**
- **TokSearch is a key element of the FDP, allowing for accelerated data processing**

Thank you!

- **Let me know if you are interested in helping using/developing these tools**
 - Use cases welcome!
- **Also, I'm happy to provide an invite for the initial release (Oct. 2024)**
- **Contact info:**
 - Brian Sammuli
 - sammuli@fusion.gat.com