

IMAS Data Model and I/O library: Status and needs

O. Hoenen¹, L. Fleury²

¹ ITER Organization, Route de Vinon-sur-Verdon, CS 90 046, 13067 St. Paul Lez Durance Cedex, France

² CEA, IRFM, 13108 St Paul-lez-Durance, France

Disclaimer: The views and opinions expressed herein do not necessarily reflect those of the ITER Organization

Integrated Modelling & Analysis Suite

IMAS is the collection of physics software that will be used to support ITER operations and research as defined in the ITER Integrated Modelling Programme.



Applications

- Independent physics codes
- Complex workflows
- Experimental data processing pipelines

Generic Tools

- **Data access, storage**, discovery, manipulation, visualization
- Support for simulation, datasets management and exploitation

Data Model

- Machine independent **data structures**
- Can serve as code coupling interface

Data Model

Covers both **simulation** and **experimental** data

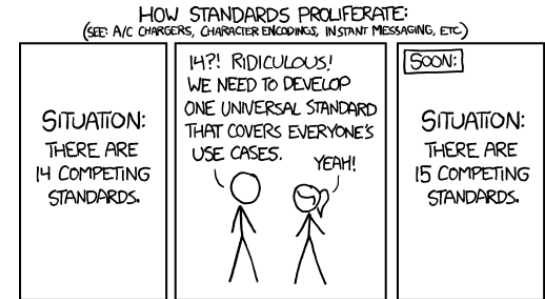
- Organized as a set of **Interface Data Structure (IDS)**

Created for **ITER**, but **generic by design**

- Used directly in experiments (ITER, WEST)
- Can map existing experiments data formats (AUG, DIII-D, EAST, HL-2A, JET, KSTAR, MAST, TCV)

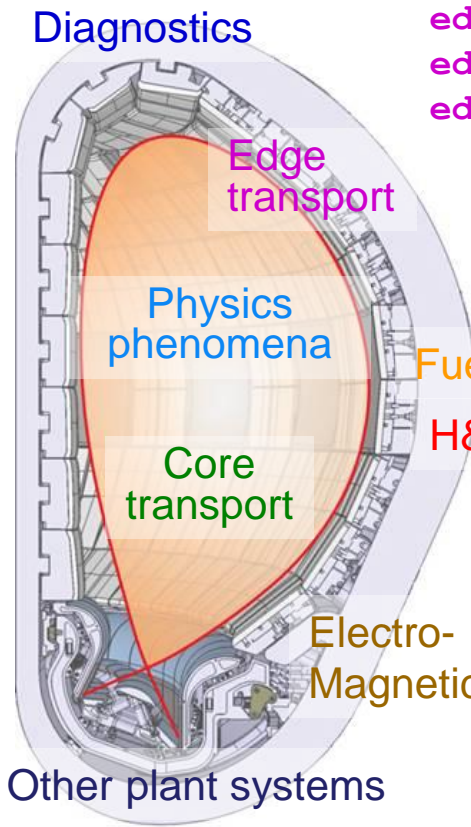
Can serve as a **community data standard**

- Ease **comparison**, **benchmarking**, **scaling**
- Towards using **data science** for research
 - **Findable**, **Accessible**, **Interoperable**, **Reusable**
 - **Machine-Learning**, **Data-Mining**, etc...



Current IDS coverage

barometry bolometer
 bremsstrahlung_visible
 calorimetry camera_ir
 camera_visible
 camera_x_rays disruption
 charge_exchange gyrokinetics
 ece hard_x_rays mhd
 interferometer mhd_linear
 langmuir_probes ntms
magnetics mse radiation
 neutron_diagnostic sawteeth
 polarimeter turbulence
 reflectometer_profile
 refractometer
 soft_x_rays cryostat
 spectrometer_mass divertors
 spectrometer_uv wall
 spectrometer_visible
 spectrometer_x_ray_crystal
 thomson_scattering



edge_profiles amns_data
edge_sources controllers
edge_transport dataset_description
 core_instant_changes dataset_fair
 core_profiles numerics
 core_sources pulse_schedule
 core_transport real_time_data
 summary
 temporary
 transport_solver_numerics
 workflow
 coils_non_axisymmetric
 em_coupling distribution_sources
equilibrium distributions
 iron_core ec_launchers
 pf_active ic_antennas
 pf_passive lh_antennas
 tf gas_injection nbi
 gas_pumping waves
 pellets

An evolving standard

New IDS added

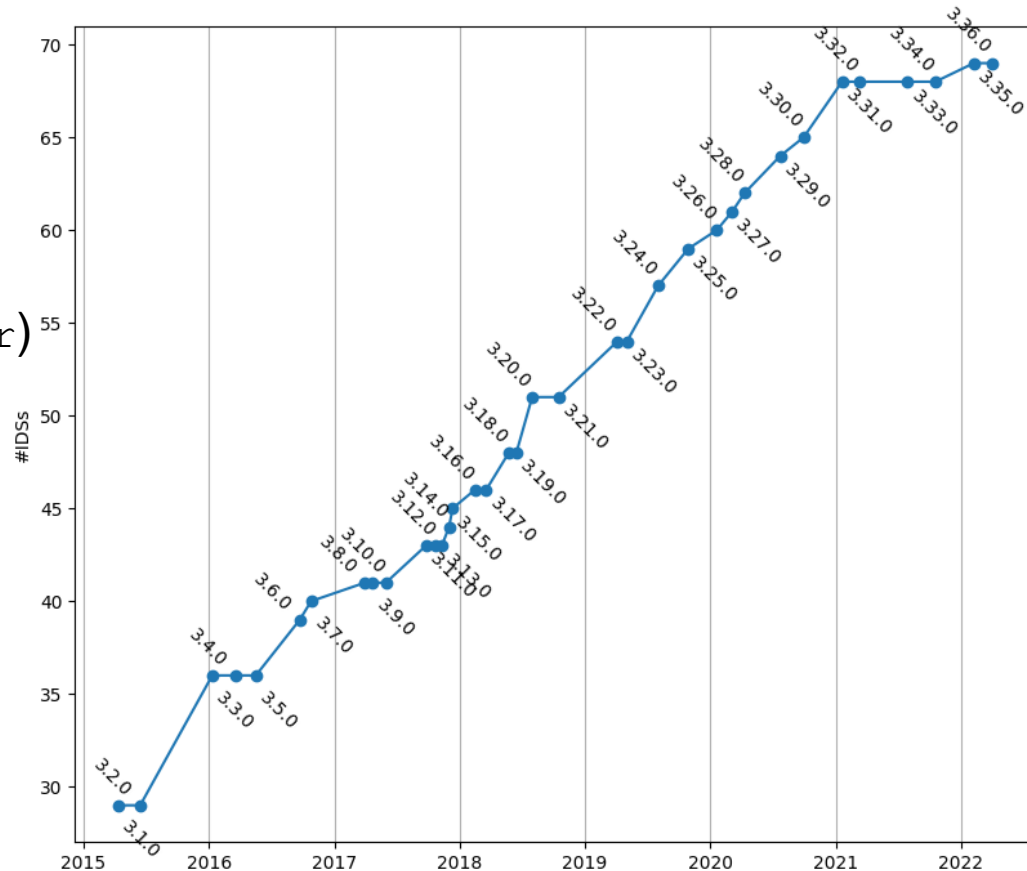
- Diagnostics
- Actuators and Plant Systems
- Physics
- Meta-Data (e.g. `dataset_fair`)

Corrections

- Fixing inconsistencies
- Clarifying definitions
- Missing quantities

Maturity level

- alpha → active
- Stricter rules for compatibility

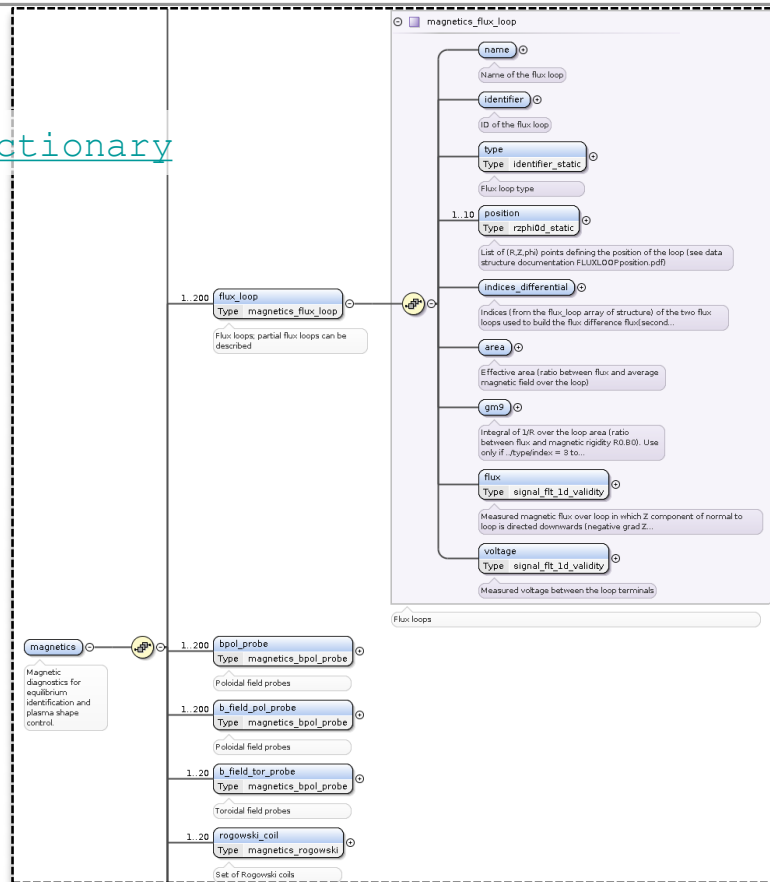


Data Dictionary

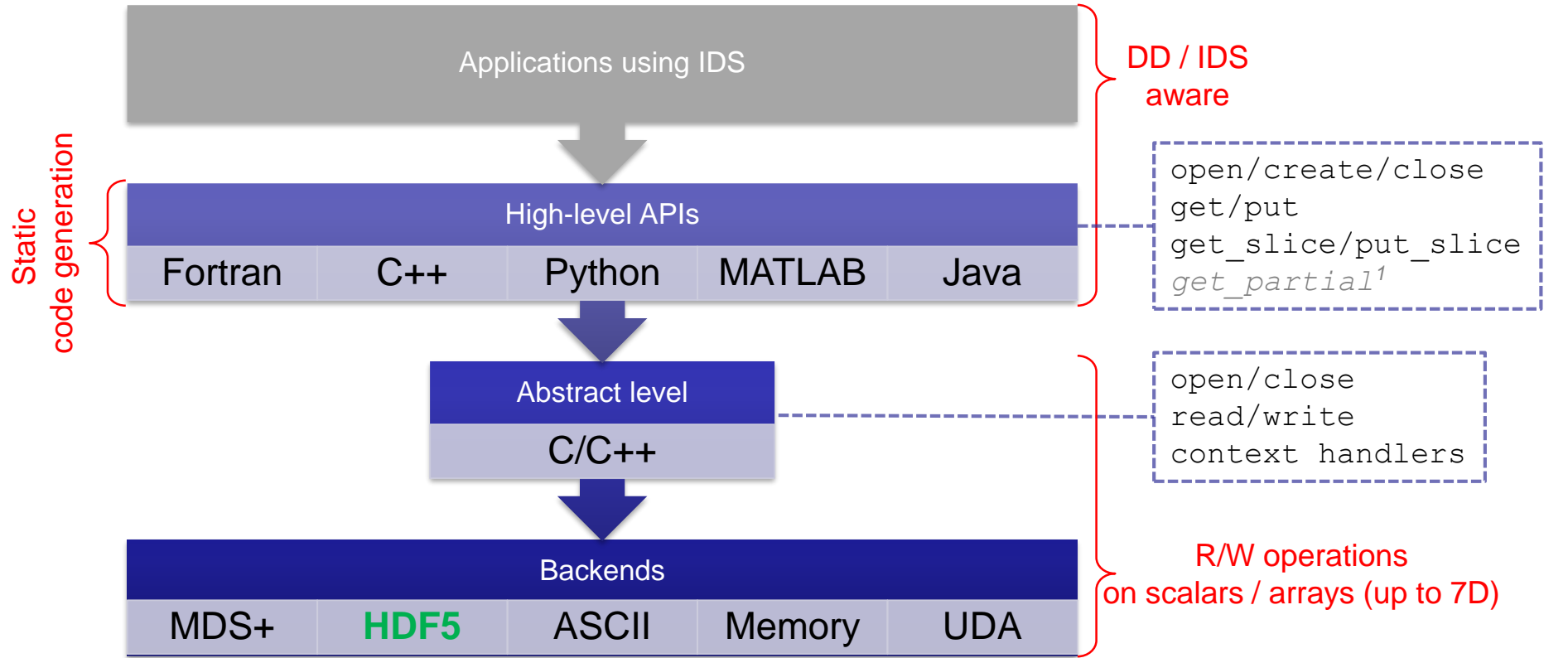
Implementation of the Data Model

<https://git.iter.org/projects/IMAS/repos/data-dictionary>

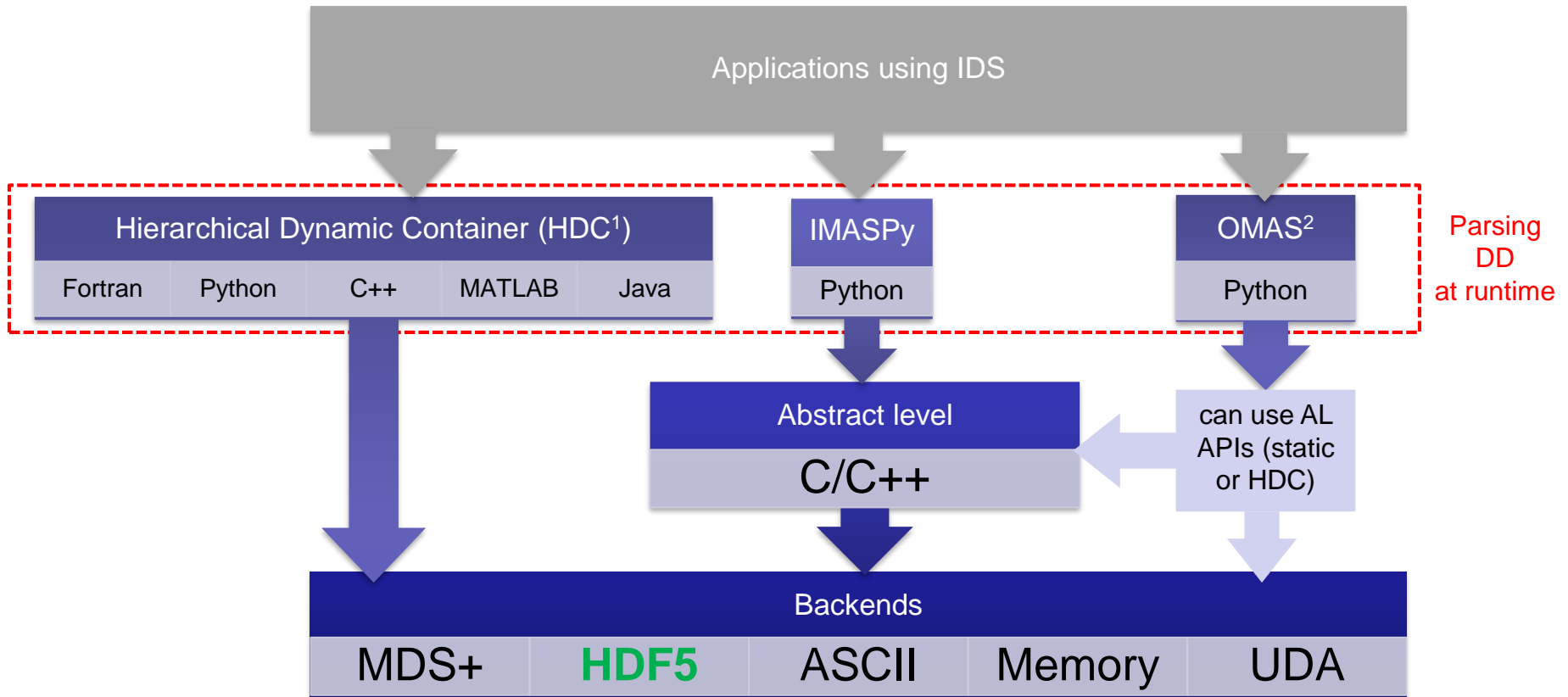
- XML description of IDS fields
 - names, types, units, coordinates, meta-data, documentation, etc...
- Fields
 - Physics quantities (with error bars)
 - Numerics (grids, ...)
 - Meta-data (provenance, identifiers, ...)
- Used for **static code generation** or **parsing at runtime** → **Access-Layer**



Access-Layer library (v4)



Alternate solutions for AL



Parsing DD at runtime

¹ <https://bitbucket.org/compass-tokamak/hdc/>

² <https://github.com/gafusion/omas/>

HDF5 backend: datasets layout

IDS are structures composed of scalars, multidimensional arrays and 1D arrays of structure (AoS)

$A[i].B[j].f_n$

Examples:

`process(i1)/reactants(i2)/metastable(:)` (from `amns_data` IDS).

`profiles_1d(itime)/ion(i1)/z_ion_square_1d(:)` (from `core_profiles` IDS)

Representation of f_n

- Hierarchical approach
 - $N_i \times N_j$ groups with one n -rank dataset
 - poor metadata / raw data ratio (`core_profiles` ~1.28)
- Tensors approach
 - a single $(n+2)$ -rank dataset
 - improved metadata / raw data ratio (`core_profiles` ~0.0053)

HDF5 backend: files organization

« master » pulse file

```
HDF5 "/home/ITER/fleuryl/public/imasdb/test/3/9998/9998/master.h5" {
GROUP "/" {
  ATTRIBUTE "HDF5_BACKEND_VERSION" {
    DATATYPE H5T_STRING {
      STRSIZE 10;
      STRPAD H5T_STR_NULLTERM;
      CSET H5T_CSET_UTF8;
      CTYPE H5T_C_S1;
    }
  }
  DATASPACE SCALAR
  DATA {
    (0): "1.0"
  }
}
ATTRIBUTE "RUN" {
  DATATYPE H5T_STD_I32LE
  DATASPACE SCALAR
  DATA {
    (0): 9998
  }
}
ATTRIBUTE "SHOT" {
  DATATYPE H5T_STD_I32LE
  DATASPACE SCALAR
  DATA {
    (0): 9998
  }
}
EXTERNAL_LINK "edge_transport_1" {
  TARGETFILE "./edge_transport_1.h5"
  ...
}
EXTERNAL_LINK "equilibrium" {
  TARGETFILE "./equilibrium.h5"
  ...
}
}
```

Backend version

Run number

Shot number

Link to edge_transport IDs
(occurrence 1)

Link to equilibrium IDs
(default occurrence)

- One file per IDS and per occurrence
- One master file with reference to each IDS

```
[fleuryl@sdcc-login03]$ ls -alh ~/public/imasdb/test/3/9998/9998
total 60M
drwxrwsr-x. 2 fleuryl fleuryl 4.0K Jan 25 14:15 .
drwxrwsr-x. 3 fleuryl fleuryl 4.0K Jan 25 13:58 ..
-rw-rw-r--. 1 fleuryl fleuryl 6.4M Jan 25 14:15 edge_profiles.h5
-rw-rw-r--. 1 fleuryl fleuryl 17M Jan 25 14:15 edge_transport_1.h5
-rw-rw-r--. 1 fleuryl fleuryl 34M Jan 25 14:15 edge_transport.h5
-rw-rw-r--. 1 fleuryl fleuryl 2.4M Jan 25 14:15 equilibrium.h5
-rw-rw-r--. 1 fleuryl fleuryl 2.3K Jan 25 14:15 master.h5
```

edge_transport_1.h5
(occurrence 1)

```
HDF5 "/home/ITER/fleuryl/public/imasdb/test/3/9998/9998/edge_transport_1.h5" {
GROUP "/" {
  ...
  GROUP "edge_transport_1" {
    ...
    DATASET "grid_ggd[&grid_subset[&base[&jacobian" {
      DATATYPE H5T_IEEE_F64LE
      DATASPACE SIMPLE { ( 3, 3, 3, 3 ) / ( 3, 3, 3, 3 ) }
      DATA {
        (0,0,0,0): 320.188, 352.19, 687.804,
        (0,0,1,0): 548.362, 99.6258, 361.164,
        (0,0,2,0): 685.705, 298.591, 708.262,
        ...
      }
    }
    ...
  }
  ...
}
```

Storage

- MDS+ backend (historical default format)
 - Requires 80M models (.tree and .characteristics files)
 - Increases with DD complexity → self-imposed size limitations (occurrences, AoS)
- HDF5 backend
 - Only a single lightweight master.h5 (~1-32K)
 - No extra costs for empty IDS/AoS → remove previous limitations

Scenario	# IDS	# time slices	MDS+ ¹	HDF5 ²
DINA (105013/1)	11	2352	396M	172M
METIS (130011/4)	13	108	2.2G	345M
JINTRAC-DINA (134174/117)	8	106 – 650	656M	387M
RANDOM (9998/9998)	69	3	788M	134M

¹ do not account for models (~80M) ² with compression enabled

R/W access time

IDS	get			get_slice			put			
core_profiles	4.71	10.02	0.71	0.04	0.40	0.03	1.77	5.50	0.34	mds+
	4.22	9.1	0.66	0.13	1.38	0.19	1.95	6.85	0.59	hdf5
core_sources	2.38	1.59	1.66	0.04	0.11	0.09	1.15	0.86	0.78	
	2.41	1.49	1.66	0.04	0.15	0.13	1.47	1.02	1.26	
core_transport	2.48	0.21	1.08	0.02	0.02	0.05	1.03	0.11	0.54	
	2.29	0.19	1.06	0.04	0.06	0.23	1.03	0.14	0.89	
equilibrium	8.07	2.23	0.50	0.39	0.10	0.02	3.80	1.36	0.33	
	6.72	3.10	0.72	0.22	0.69	0.18	8.36	18.2	1.53	
summary	0.07	0.11	0.08	0.24	-	0.20	0.09	0.20	0.11	
	0.02	0.05	0.21	0.07	-	0.07	0.02	0.05	0.37	
	DINA	METIS	JINTRAC-DINA							

Open points

- Observations

- Performance is related to IDS inner structure (but not only, e.g. equilibrium)
 - Gather do's and don'ts when defining DD extensions
- Different backends have different non-overlapping performance sweet-spots
 - hierarchical vs tensor representation, full read vs slicing, simulation vs acquisition data
 - also depends on the backend parameters (chunks size, compression, etc...)

- Needs

- Moderate data size (scenario simulation) → tests at scale experimental data (e.g. WEST) or synthetic diagnostic data
- Continued improvement (performance, reliability) of the storage backend(s) → balance between fine tuning and general purpose
- Parallel I/O
- Recent interest in storing data + functions → towards a more compact Data Model