# ILDG Lunch



Hubert Simma
on behalf of the
ILDG Working Groups

10 August 2022

## Overview

# Naive Data Consumer

**Collaboration internal:**

- everything is known about our configs (location, tracking, reliability, . . . )
- we have a clear data managment plan
- data stewards take care of our (meta)data
- usage rules are well defined and known

**Community wide:**

- we somehow know about existence of useful data
- get the data at no cost (human and CPU)
- use data freely to do high quality research
- be nice and acknowledge source of data

# Naive Data Provider

**Collaboration internal:**

- follow a well-defined and smooth workflow
- public and internal data can be handled in the same way (no extra efforts at end of embargo times)
- public data becomes (easily) published (citable)
- efforts are rewarded by funding agencies

**Community wide:**

- dump precious (meta) data to some storage at no cost (human and storage resources)
- declare it public
- get it used by others
- receive credits/citations

# Naive Implementation

Have one machine with a big disk for all (collaboration internal or community shared) configs?

Problems / Challenges:

- resources
- scalabilty
- usability
- access control
- bookkeeping
- credits and citation

FAIR data?

No free lunch! But ILDG can guide . . .

# FAIR Principles

**F**indable

       **A**ccessible

             **I**nteroperable

                    **R**eusable

force11.org
⋮
Wilkinson 2016
⋮
go-fair.org

- It is becoming a mandatory requirement by funding agencies
  "The [European] Commission will work with global policy and research partners to foster cooperation and to create a level playing field in scientific data sharing and data-driven science."

  EU Commission, COM(2016)178

- provides guiding principles, not an implementation
- conceptually refers to three types of entities:
  - data = any digital object
  - metadata (MD) = information about digital object
  - infrastructure
- requires machine actionable (meta)data

# What does "findable" mean?

**Findable**

F1 globally unique and persistent ID assigned to (M)D

F2 data described with rich MD

F3 MD includes data ID of data

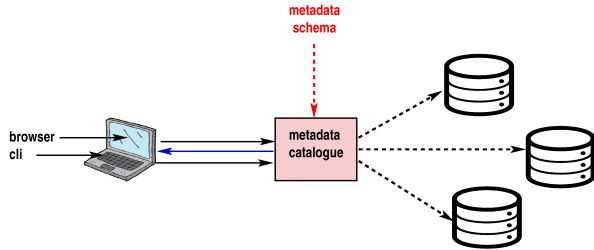F4 (M)D registered or indexed in a searchable resource

Metadata includes information on

- content (general and domain-specific vocabulary)
- provenance (who, when, where, how?)
- access (format, path, license, . . . )
- . . .

Metadata

- follows a well-defined and rich schema
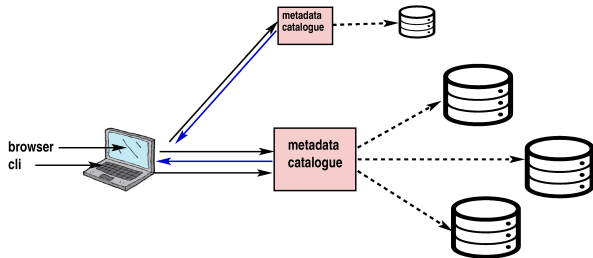- stored separately from data (big)
- searchable in central catalog of each RG

Metadata

- follows a well-defined and rich schema
- stored separately from data (big)
- searchable in central catalog of each RG

# Unique identifiers
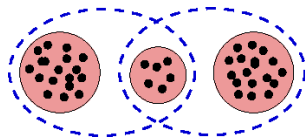
- **Ensembles:** have only MD (content, access permissions, ...)
  $$\texttt{mc}://\langle rg\rangle/\langle collab\rangle/\langle proj\rangle/\ldots$$
- Configurations: MD (related ensemble, provenance info) **and** actual data
  $$\texttt{lfn}://\langle rg\rangle/\langle collab\rangle/\langle proj\rangle/\ldots$$



| ID | entity | relation | content | data storage | access control |
|---|---|---|---|---|---|
| lfn | config | mc | yes | yes | no |
|  | ↓ |  |  |  | ↑ |
| mc | ensemble | — | yes | no | yes |
|  | ↑↑↑ |  |  |  |  |
| ∗) | publication | set of mc | yes | no | no |

∗) ILDG 1.0 has no official registration of IDs or publication metadata yet!

# DOI and Data Publishing

## Data Publishing

- Registration of persistent identifier (DOI)
- Metadata for registration (DataCite)
- Landing Page (hosting and automatic generation)
- Harvesting of metadata

Exploratory setups by JLDG and USQCD
- using national registration authorities (JaLC, OSTI)
- workflow and metadata for registration and generation of landing pages
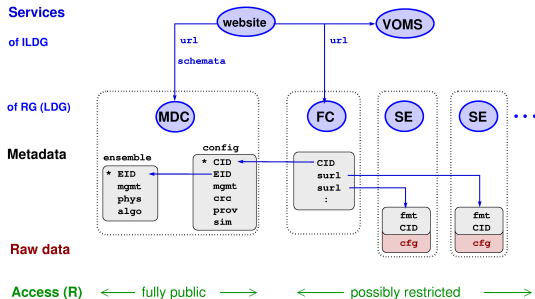
Possible directions in ILDG 2.0
- establish workflow for registration, generation and hosting of landing pages (e.g. Zenodo)
- extended metadata support
- dedicated metadata harvesting (e.g. by INSPIRE)
- common registration authority

# What does "accessible" mean?

**Accessible**

| | |
|---|---|
| A1 | (M)D retrievable by ID using standardized protocols |
| A1.1 | protocol is open, free, and universally implementable |
| A1.2 | protocol allows authentication/authorization procedure where necessary |
| A2 | MD accessible even if data is no longer available |

- A1 can be achieved e.g. by a File Catalog: ID $\mapsto$ storage location(s)
- Accessible does not imply (unrealistic) public access without authentication
- MD is precious even without the associated data

# How does ILDG address "accessible"?



- all metadata is publicly accessible (from MDC)
- well-defined community-wide metadata schema
- metadata available in a standard markup language
- standardized protocols and API of services for access to data and metadata

# What does "interoperable" mean?

**Interoperable**

I1 (M)D use a formal, accessible, shared, and broadly applicable language

I2 (M)D use vocabularies that follow FAIR principles

I3 (M)D include qualified references to other (M)D

- ability of data (or tools) from non-cooperating resources to integrate (or work together) with minimal effort

## Common standards for

- Metadata schema
- Data format
- API and URL for web services of regional grids

New directions:

- Extend ILDG format to include support for HDF5
  - definition of ILDG packing rules
  - convenient tools for packing and conversion
- Token-based authentication
- REST API

# What does "reusable" mean?

## Reusable

**R1** (M)D richly described with plurality of accurate and relevant attributes

**R1.1** (M)D released with clear and accessible data usage license

**R1.2** (M)D associated with detailed provenance
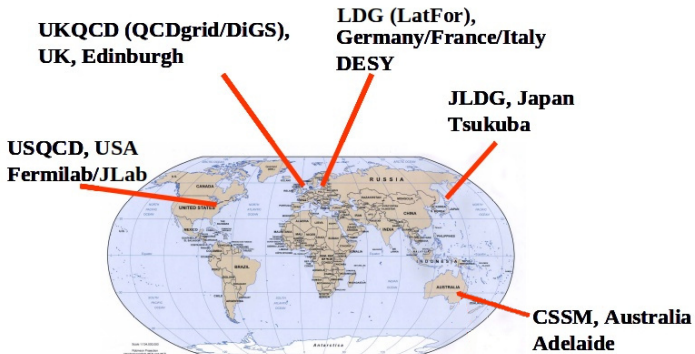
**R1.3** (M)D meet domain-relevant community standards

- reference to a paper may not be sufficient
- good scientific practice ↔ FAIR
- also related to verifiable invariance of results     (see presentation by Ed Bennett)
  - reproducibility: same data + same analysis
  - replicability: new data + same analysis
  - robustness: same data + new analysis

# How does ILDG address "reusable"?

**Ensemble MD**

- Physics
- Algorithm
- Management

**Config MD**

- Markov step
- Implementation (machine, code)
- Management (creator, date, checksum)

But no license aspects! (cf. R1.1)

# Global Structure of ILDG

## ILDG

- Federation of autonomous Regional Grids (RG)
- Virtual Organisation (VO)
- Agreed standards



UKQCD (QCDgrid/DiGS), UK, Edinburgh

LDG (LatFor), Germany/France/Italy DESY

JLDG, Japan Tsukuba
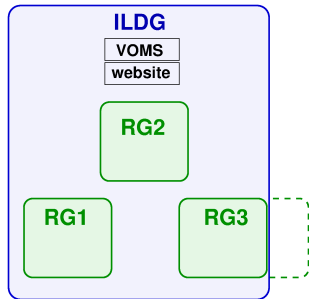
USQCD, USA Fermilab/JLab

CSSM, Australia Adelaide

# Services and Organization of ILDG

ILDG operates only 2 global services

- VO registration (VOMS)

  registry of ILDG users (groups and roles)

  used for authentication to storage elements

- Website (temporary mirror)

  specification of standards and conventions

  URLs of services of each regional grid (`Services.xml`)

Organization

- Board
- Metadata Working Group (MDWG)
- Middleware Working Group (MWWG)

# Regional Grids

## Services operated by each Regional Grid

- Metadata Catalog (MDC)
- File Catalog (FC)
- Storage Elements (SE)
- Website with RG-specific information

Regional Grids: CSSM, JLDG, LDG, UKQCD, USQCD

- are implemented with different architectures and technologies
- operate in an autonomous way with individual policies

Examples

- JDLG: single SE, no specific access control
- LDG: multiple SE, fine grained access control

# ILDG use cases

Consumer (collaboration internal)

- `lfind`: search in metadata catalog
- `lget`: download data and metadata

↓

Consumer (community wide)

- optionally also use common search engines
- cite DOIs for published data used

Provider (collaboration internal)

- `lpack`: generate markup*) and pack data
- `linit`: register ensemble metadata
- `lput`: upload config data and metadata

↓

Provider (community wide)

- optionally register DOI and generate landing page
- change access control flag
- have data citation record

∗) trivial if information is already collected during production!

# Plans for ILDG 2.0

Keep basic concepts of ILDG 1.0 (well defined metadata schema, interoperable services, . . . ), but make implementation fully compliant with FAIR principles and upgrade to modern technologies

Restore and improve usability:

- Support for **DOI registration and data publishing**
  Then (sets of) ensembles in ILDG become also findable e.g. by
  INSPIRE or other search engines $\xrightarrow{DOI}$ landing pages $\xrightarrow{MDC}$ ensemble IDs
  and properly citable in journal papers
- Adjustment of **metadata schemas** and data format
- User **tools and documentation**!

Upgrade of technologies:

- Token-based authentication
  (commonly used for cloud services and replacing cumbersome Grid Certificates)
- REST APIs for services

Active contributions (e.g. in MDWG and MWWG)
from all of you are needed and welcome!

# Backup Slides

## Tentative Summary of the Parallel Session

|  | Collab | Public | ILDG | #ens | #cfg | TB |
|---|---|---|---|---|---|---|
| Steve | MILC | 1 | 0 | >25 | 75k | 1000 |
| Peng | CLQCD $T = 0$ | 1 | 1 | 10 | 5k | 14 |
|  | CLQCD $T > 0$ | 1 | 1 | 28 | 150k | 150 |
| Yoshinobu | PACS | 1 | 2/3 | 3 | 100 | 60 |
| Ryan | FASTSUM | 1 | 1 | 25 | 22k | 40 |
| Anthony | OpenLAT | 1 | 2 | 10 | 6k | 13 |
| Rajan | JLab/W&M/LANL/MIT | 0 | 0 | 13 | 90k | 2000 |
| Issaku | JLQCD | 1 | 2/3 | 230 | 60k | 20 |
| Andrey | TMFT | 1 | 1 | 60 | 50k | 26 |
| Robert | RBC-UKQCD | 1 | 0 | 41 | ? | ? |
| Christian | HotQCD | 1 | 2 | 58 | 15M | 2200 |
| Wolfgang | CLS | 1 | 2 | 55 | 125k | 1000 |
| Bartosz | ETMC | 1 | 2/3 | 21 | 100k | 1500 |
| Takumi | HAL | 1 | 2 | 1 | 1.4k | 70 |
| James | QCDSF-UKQCD-CSSM | 1 | 2/3 | 60 | 90k | 300 |

Public availability: $0 =$ no, $1 =$ yes, but after some embargo time, $2 =$ yes, already now
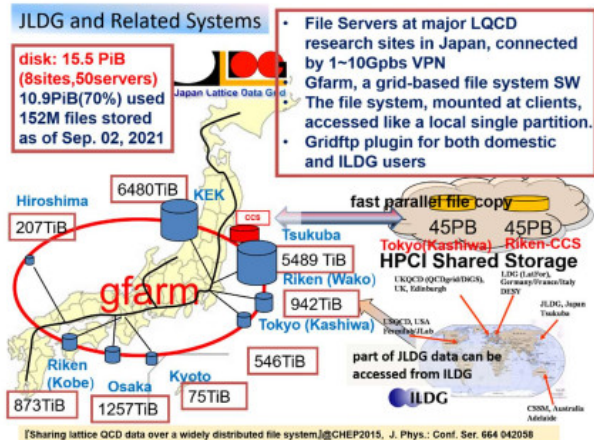ILDG: $0 =$ no interest, $1 =$ interest, $2 =$ planned, $3 =$ already using

## Command-line tools

```
lls [-g <grid>]              list all ensembles (of specified RG)
lls <uri>                    list configs of ensemble <uri>

lfind [-g <grid>] -e <xpath> Xpath search in ensemble MD (of specified RG)
lfind [-g <grid>] -c <xpath> Xpath search in config MD (of specified RG)

lget <uri>                   download MD of ensemble <uri>
lget <lfn>                   download MD of config <lfn>
lget -d <lfn>                download data of config <lfn>
```
```
linit ...                    register new ensemble
lput ...                     upload ...
ladm ...                     manage access control
```

# JLDG architecture

- Single federated storage system (GFARM)
- JLDG internal write access
- Fast read access to ILDG data available for VO members
- Transition to token-based authentication



T. Yoshie

# USQCD ideas



K. Chard et al. 2017