# The International Lattice Data Grid

**ILDG**

Hubert Simma (DESY) and Carsten Urbach (Uni Bonn)

PUNCH4NFDI workshop on "open data"
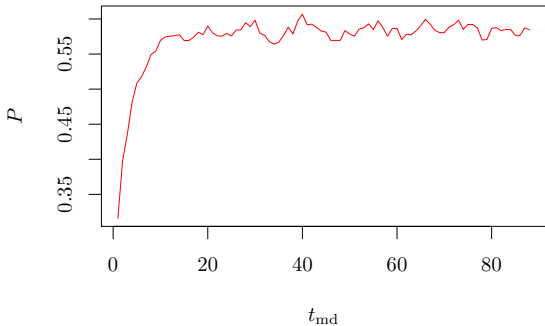
Feb 11, 2021

# LQCD: Studying QCD on a Lattice

### 1. Simulation Code

```
gaugefield g(L, T);
hamiltonianfield h(L, T);
integrate_md(h, g, latticeparams);
...
```
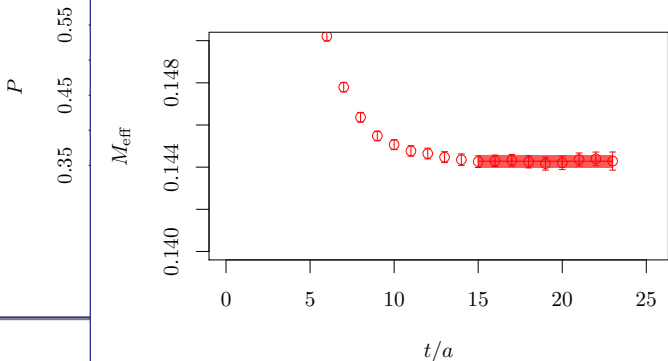
# LQCD: Studying QCD on a Lattice

## 1. Simulation Code

```
gaugefi
hamilto
integra
...
```

## 2. Run Simulation

# LQCD: Studying QCD on a Lattice



**1. Simulation Code**

```
gaugefi
hamilto
integra
...
```

**2. Run Simulation**

**3. Measure**

# LQCD: Studying QCD on a Lattice

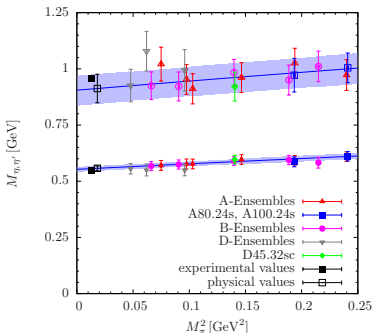**1. Simulation Code**

```
gaugefi
hamilto
integra
...
```

**2. Run Simulation**

**3. Measure**
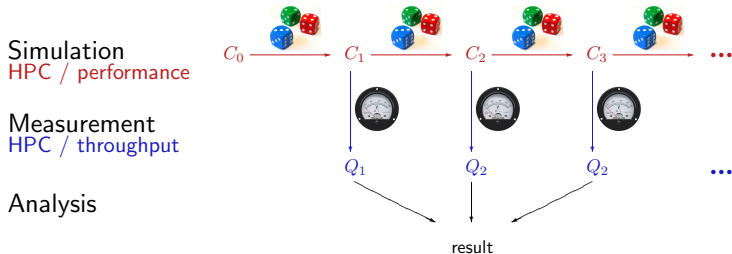


**4. Combine and Extrapolate**

# Lattice QCD, from the data perspective

Given few physical and simulation parameters (coupling, quark masses, ...)

- a Markov Chain ensemble $\{C_i\}$ is generated using MCMC, $i = 1, ..., N_{\text{cfg}}$

- each element $C_i$ (configuration) up to $O(200)$ Gigabytes large
  no compression possible

- ensemble ($N_{\text{cfg}} \sim 10^4$) **takes years** to generate on largest HPC systems
  $O(100000)$ core hours per $C_i$

- many ensembles $O(20)$ are needed to obtain physical results

$\Rightarrow$ ensembles (our *raw data*) **highly valuable** (and costly)!

$\Rightarrow$ ensembles the basis for a **plethora of physical observables** $Q$

**Community concluded: ensembles need to be shared and preserved!**

# Lattice QCD Workflow
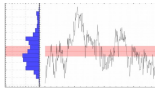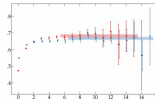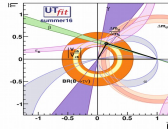


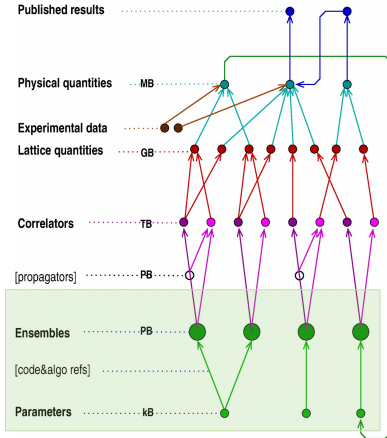(1) **"Simulation"** or "Gauge field generation" by MCMC

(2) **"Measurement"** and average of (primary) observables $Q$

(3) **"Statistical Analysis"**

- extract physical quantities e.g. $\langle Q \rangle \sim e^{-tm}$
- combine different ensembles, fit, extrapolate, ...  [1]

# LQCD Workflow towards PUNCH

In practice, workflow more complex

$\mapsto$ tune, calibrate, cross-check, optimise, ...

# International Lattice Data Grid (ILDG)

**Motivation**

**Highly valuable configurations $C_i$ should be shared internationally!**

**Time Line**

| | | |
|---|---|---|
| $\approx 2001$ | Plans to organize and set up infrastructures for community-wide sharing of configs | |
| $\approx 2003$ | ILDG architecture as a grid of (regional) grids | `hep-lat/0309029` |
| . | Development of common metadata schemata | `hep-lat/0409055` |
| . | " " standard data format for configs | |
| . | Setup of regional grids and Virtual Organization | |
| . | Specification of interoperable web services | |
| | | `hep-lat/0609012` |
| 2008 | Fully operational infrastructure and services | `arXiv/0910.1692` |
| 2013 | Last (minimal) revision of QCDmlEnsemble schema | |

# ILDG vs. Regional Grids (RG)

**ILDG**

- defines a **grid of inter-operable (data-) grids**

  (initially 5 regional grids: Japan, EU, UK, US, Australia)

- formally consists of
  - Virtual Organization (VO)
  - Specification of **Services** (File Catalog, Metadata Catalog)
  - Specification of **Data Formats** (Lime container)
  - Specification of **Metadata Standards** (Ensembles, Configs)
  - URLs of Services of each Regional Grid
  - Metadata and Middleware working groups, governance

- is implemented by
  - VOMS (VO)
  - Website (specifications and docu)
  - Board and working groups

## Regional Grids

e.g. the **Latfor Data Grid (LDG)** in continental Europe

**Implement and operate** the following services

- Webpage = RG-specific info
- Metadata Catalog (MDC) = LDG-specific implementation
- File Catalog (FC) = WLCG middleware
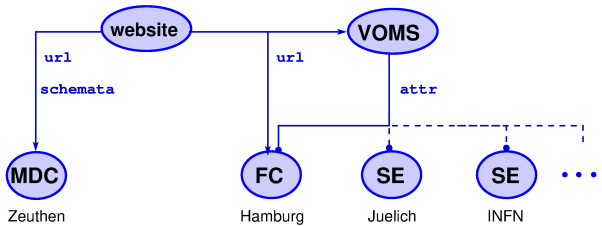- Data Storage (SE) = shared with WLCG

Backward compatibility is a **must** for ILDG!

# Distributed Architecture of ILDG

# Distributed Architecture of ILDG

# ILDG Metadata (MD)

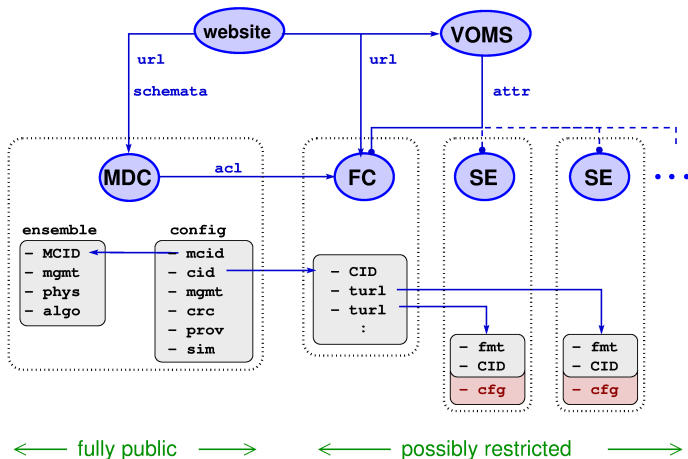## MD Schemata

Rich and extensible ($20\ldots\infty$ elements)

- Unique ID's                                MCID, CID
- Revisions (who, when, why)                 mgmt
- Provenance (where, when, who)              prov
- Integrity [and access]                     crc, acc
- LQCD-specific info                         phys, algo, sim

## MD Catalog (LDG Implementation)

- searchable (Xpath)              $\rightarrow$ **F** indable
- publicly accessible (https)     $\rightarrow$ **A** ccessible
- standard SOAP interface (WSDL/axis)  $\rightarrow$ **I** nteroperable
- free MD schemata (XSD)          $\rightarrow$ **R** eusable
- scalable data base (eXist)
  $O(10^3)$ ensembles, $O(10^6)$ configs

---

# Summary

- completed a long and difficult community-wide process to **converge** to rich and flexible MD schemata (still adequate and FAIR compliant)

- the developed schemata and standards can serve as a **blueprint / seed** for the PUNCH community

- can provide a **working prototye** for building blocks and architecture of research data infrastructure

- needs a thorough **redesign** based on modern data-lake concepts and up-to-date (cloud and web service) technologies

# Future (I)LDG directions within PUNCH

- exchange and share experiences and knowledge for open data
  $\rightarrow$ integration into PUNCH community

- incorporate derived LQCD data types
  $\rightarrow$ towards reproducible analysis chains

- design and converge to PUNCH-wide authentication services

- Joint development of modernised (RESTfull) Web Services