

# ILDG Metadata



## Working Groups

Hands-on Workshop

June 15, 2023

# Overview

1. Overview and Requirements
2. Digression: XML Technologies
3. Ensemble Metadata
4. Configuration Metadata
5. Concluding Remarks

# Content

1. Overview and Requirements
2. Digression: XML Technologies
3. Ensemble Metadata
4. Configuration Metadata
5. Concluding Remarks

# Metadata Content

- 1. Physics**

Description of action and simulation parameters

- 2. Algorithm**

Information about the used algorithm and algorithmic parameters

- 3. Source code**

Used code, compile time parameters, compilation software

- 4. Machine**

Machine used to produce the configuration

- 5. Data management**

Information related to data provenance and checksums

# Metadata Requirements

- **Standardised** description of the data (=**metadata**)
  - Solution strategy: **XML documents** which conform to an **XML schema**
- **Unique** description of the data
  - Avoid situation that action is described in two different ways,  
e.g. Iwasaki gauge action

$$S = \beta (c_0 \text{ plaquette} + c_1 \text{ rectangular})$$

or

$$S = \tilde{c}_0 \text{ plaquette} + \tilde{c}_1 \text{ rectangular}$$

- Forward compatible **extensibility** of the standard
  - A metadata document conforming to the current version of the standard should also conform to its future versions
- Broad **coverage** of data that can be described
  - In parts achieved by foreseeing flexibility

## Example (1/2)

- SU(3) Wilson plaquette action with gluon fields  $U$  and coupling  $\beta$ :

$$S_G = \beta \sum_{\text{plaquette}} \frac{1}{3} \text{Re} \text{Tr} (1 - U_{\text{plaquette}})$$

## Example (2/2)

- (Almost) ILDG-compliant mark-up:

```
1 <?xml version="1.0"?>
2 <markovChain>
3 ...
4   <plaquetteGluonAction>
5 ...
6     <gluonField>
7       <gaugeGroup>SU(3)</gaugeGroup>
8 ...
9   </gluonField>
10  <couplings>
11    <elem>
12      <beta>5.2</beta>
13    </elem>
14  </couplings>
15 </plaquetteGluonAction>
16 ...
17 </markovChain>
```

# Ensembles and Configurations

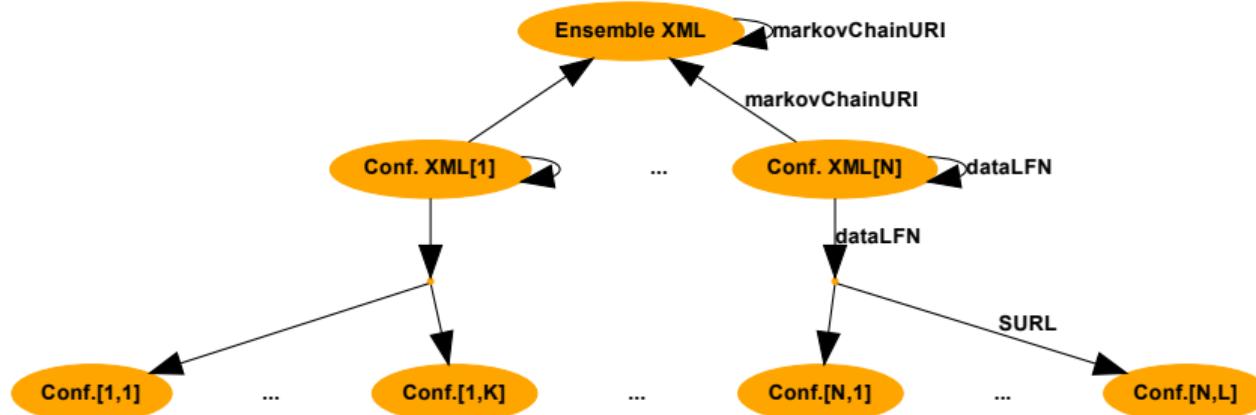
- **Aim:** Avoid replication of complicated data structures
- Markov chains are used to generate gauge field configurations

Markov chain     $\leftrightarrow$     Ensemble XML

Markov step     $\leftrightarrow$     Configuration XML

# Linking Metadata and Data

Objects	Links
Ensemble XML document	markovChainURI
Configuration XML document	dataLFN
Binary data file	



# Content

1. Overview and Requirements
2. Digression: XML Technologies
3. Ensemble Metadata
4. Configuration Metadata
5. Concluding Remarks

# Why XML?

- Both human-readable and machine-readable
- Standardized by the World Wide Web Consortium's (W3C) [XML 1.0 Specification](#)
- Availability of the XML schema (XSD) technology that allows to define the necessary metadata for interpreting and validating XML documents
- Availability of standardised technologies to query XML documents
  - W3C's XPath (current version: [XPath 3.1](#))
  - W3C's XQuery (current version: [XQuery 3.1](#))
- Many tools for processing XML documents
  - For parsing: libxml2 with interfaces for C, C++, Python, Perl
  - For formatting and validation: xmllint

# Background on XML

- Simple example:

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <dataEntry>
3   <key>7</key>
4   <value>Hello world!</value>
5 </dataEntry>
```

- Key elements:

- Tag: Markup construct that begins with < and ends with >
  - Start tag: <elem>
  - End tag: </elem>
  - Empty-element tag: <elem />
- Element: Logical document component starting/ending with a tag (or empty-element tag)
- Attribute

# XML Schema

- Technology to formally describe the elements in an XML document
- Simple example:

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema">
3   <xs:element name="dataEntry">
4     <xs:complexType>
5       <xs:sequence>
6         <xs:element name="key" type="xs:decimal"/>
7         <xs:element name="value" type="xs:string"/>
8       </xs:sequence>
9     </xs:complexType>
10   </xs:element>
11 </xs:schema>
```

# Validating XML Files (1/2)

- Document validation using the xmllint tool:

```
% xmllint --schema ./tst-schema.xml ./tst-doc-ok.xml
<?xml version="1.0" encoding="UTF-8"?>
<dataEntry>
  <key>7</key>
  <value>Hello world!</value>
</dataEntry>
./tst-doc-ok.xml validates
```

# Validating XML Files (2/2)

- Try document with incorrect tag:

```
% xmllint --schema ./tst-schema.xml ./tst-doc-bad-1.xml
<?xml version="1.0" encoding="UTF-8"?>
<dataEntry>
  <idx>7</idx>
  <value>Hello world!</value>
</dataEntry>
./tst-doc-bad-1.xml:3: element idx: Schemas validity error : Element 'idx': This element is not expected. Expected is ( key ).
./tst-doc-bad-1.xml fails to validate
```

- Try document with incorrect element content:

```
% xmllint --schema ./tst-schema.xml ./tst-doc-bad-2.xml
<?xml version="1.0" encoding="UTF-8"?>
<dataEntry>
  <key>seven</key>
  <value>Hello world!</value>
</dataEntry>
./tst-doc-bad-2.xml:3: element key: Schemas validity error : Element 'key': 'seven' is not a valid value of the atomic type 'xs:decimal'.
./tst-doc-bad-2.xml fails to validate
```

# Content

1. Overview and Requirements
2. Digression: XML Technologies
3. Ensemble Metadata
4. Configuration Metadata
5. Concluding Remarks

# Ensemble XML Overview

Top-level elements of the ensemble XML file:

- **markovChainURI**
  - A globally unique identifier that identifies
    - An Ensemble XML document
    - A set of configurations (metadata + data)
- **management**
  - Data relevant for tracing and recording the origins of the data (data provenance)
- **physics**
  - Area where providers must provide all information needed to define the physical parameters of the data set (ensemble)
- **algorithm**
  - Area where providers can provide all information needed to describe the used algorithms (except for configuration-specific information)

# Ensemble XML: markovChainURI

- Schema mandates an Uniform Resource Identifier (URI)
- ILDG convention not enforced by schema to ensure uniqueness:

mc://<reg>/<col>/<prj>/...

- Regional grid (reg)
- Collaboration (col)
- Project (prj)

- Warning: ILDG currently does not have mechanisms to check uniqueness of the markovChainURI

## Example:

```
1 <?xml version="1.0"?>
2 <markovChain>
3   <markovChainURI>
4     mc://ldg/qcdsf/clover_nf2p1/b7p20kp12450kp12450-16x32
5   </markovChainURI>
6   ...
7 </markovChain>
```

# Ensemble XML: management

- Collaboration
  - String unconstrained by the schema, but should be consistent and unique
- Project name
  - String unconstrained by the schema, but should be consistent and unique at collaboration level
- Archive history
  - List of entries for tracking data provenance
  - Each entry comprises
    - Action: add, replace, remove
    - Participant identified by name and institution
    - Date

## Example:

```
1 <?xml version="1.0"?>
2 <markovChain>
3 ...
4   <management>
5     <collaboration>qcdsf</collaboration>
6     < projectName>clover_nf2p1</projectName>
7     <archiveHistory>
8       <elem>
9         <revisionAction>add</revisionAction>
10        <participant>
11          <name>Yoshifumi Nakamura</name>
12          <institution>NIC/DESY</institution>
13        </participant>
14        <date>2007-03-30T14:34:29+02:00</date>
15      </elem>
16    </archiveHistory>
17    <replicate>deny</replicate>
18  </management>
19 ...
20 </markovChain>
```

# Ensemble XML: physics / Overview

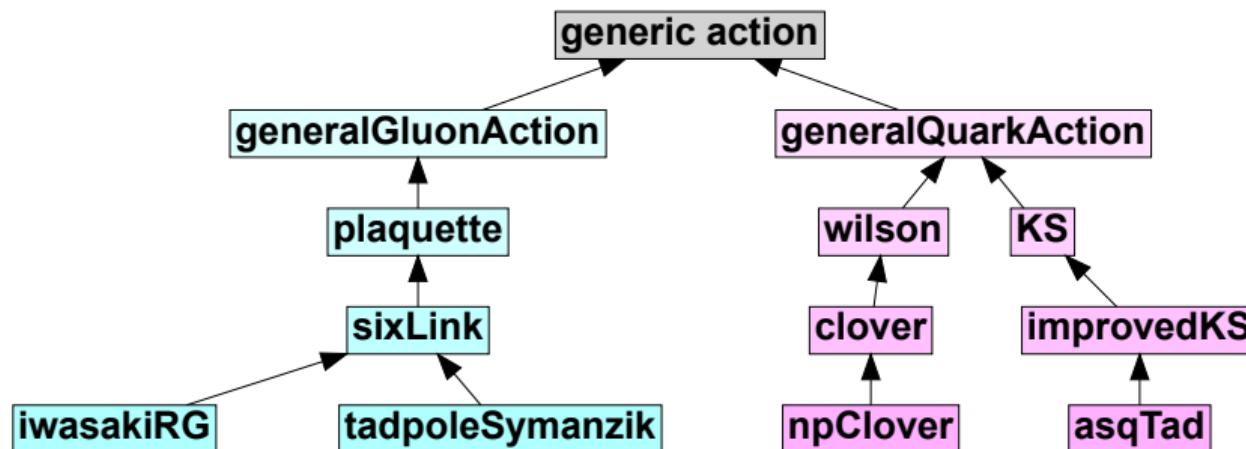
- size: Lattice size
  - Unbounded list of elements
  - Each element comprises of name and length
- action/gluon: Gluonic part of the action
- action/quark: Fermionic part of the action

## Example:

```
1 <?xml version="1.0"?>
2 <markovChain>
3 ...
4 <physics>
5   <size>
6     <elem>
7       <name>X</name>
8       <length>16</length>
9     </elem>
10    <elem>
11      <name>Y</name>
12      <length>16</length>
13    </elem>
14    <elem>
15      <name>Z</name>
16      <length>16</length>
17    </elem>
18    <elem>
19      <name>T</name>
20      <length>32</length>
21    </elem>
22  </size>
23 ...
24 </physics>
25 ...
26 </markovChain>
```

# Ensemble XML: physics / Action

- The description of action parameters is most critical to ensure metadata being unique and extensible
- Strategy: Hierarchy of actions



# Ensemble XML: physics / Gluonic Part (1/2)

- The element `gluon` contains an unbounded list of gluon actions
- Various gluon actions have been defined since the schema has been established
  - `plaquetteGluonAction`
  - `iwasakiRGGluonAction`
  - `tpLuescherWeiszGluonAction`
  - `anisotropicTpWilsonGluonActionType`
  - ...

## Example:

```
1 <?xml version="1.0"?>
2 <markovChain>
3 ...
4   <physics>
5   ...
6     <action>
7       <gluon>
8         <plaquetteGluonAction>
9           <glossary>http://...</glossary>
10          <gluonField>
11            <gaugeGroup>SU(3)</gaugeGroup>
12            <representation>fundamental</representation>
13            <boundaryCondition>
14              <elem>periodic</elem>
15              <elem>periodic</elem>
16              <elem>periodic</elem>
17              <elem>periodic</elem>
18            </boundaryCondition>
19          </gluonField>
20          <beta>7.2</beta>
21        </plaquetteGluonAction>
22      </gluon>
23 ...
24    </action>
25 ...
26  </physics>
27 ...
28 </markovChain>
```

# Ensemble XML: physics / Gluonic Part (2/2)

- All gluon actions must contain
  - glossary: URI to documentation
  - gluonField
  - topologyFixing (optional)
- The element gluonField must contain the following sub-elements
  - gaugeGroup
    - Allowed values: SU(3), SU(2), ...
  - representation:
    - Allowed values: fundamental, adjoint
  - boundaryCondition
    - List of elements with the following values: periodic, antiperiodic, dirichlet
- Other elements (e.g., beta) are action specific

## Example:

```
1 <?xml version="1.0"?>
2 <markovChain>
3 ...
4   <physics>
5   ...
6     <action>
7       <gluon>
8         <plaquetteGluonAction>
9           <glossary>http://...</glossary>
10          <gluonField>
11            <gaugeGroup>SU(3)</gaugeGroup>
12            <representation>fundamental</representation>
13            <boundaryCondition>
14              <elem>periodic</elem>
15              <elem>periodic</elem>
16              <elem>periodic</elem>
17              <elem>periodic</elem>
18            </boundaryCondition>
19          </gluonField>
20          <beta>7.2</beta>
21        </plaquetteGluonAction>
22      </gluon>
23      ...
24    </action>
25    ...
26  </physics>
27  ...
28 </markovChain>
```

# Ensemble XML: physics / Fermionic Part (1/3)

## Example:

- The element `quark` contains an unbounded list of quark actions
- Various quark actions have been defined since the schema has been established
  - `npCloverQuarkAction`
  - `overlapQuarkAction`
  - ...
- All quark actions must contain
  - `glossary`: URI to documentation
  - `quarkField`
  - `numberOfFlavours`
  - `linkSmearing` (optional)
- Other elements (e.g., `kappa`) are action specific

```
1  <?xml version="1.0"?>
2  <markovChain>
3  ...
4  <physics>
5  ...
6  <action>
7  ...
8  <quark>
9   <npCloverQuarkAction>
10    <glossary>...</glossary>
11    <quarkField>...</quarkField>
12    <numberOfFlavours>2</numberOfFlavours>
13    <linkSmearing>...</linkSmearing>
14    <kappa>0.12450</kappa>
15    <cSW>1.0</cSW>
16  </npCloverQuarkAction>
17  <npCloverQuarkAction>
18    <glossary>...</glossary>
19    <quarkField>...</quarkField>
20    <numberOfFlavours>1</numberOfFlavours>
21    <linkSmearing>...</linkSmearing>
22    <kappa>0.12450</kappa>
23    <cSW>1.0</cSW>
24  </npCloverQuarkAction>
25  ...
26  </action>
27  ...
28 </physics>
29 ...
```

# Ensemble XML: physics / Fermionic Part (2/3)

## Example for quarkField:

- The element `quarkField` must contain the following sub-elements
  - `normalisation`
    - Allowed values are not prescribed
  - `boundaryCondition`
    - List of elements with the following values: `periodic`, `antiperiodic`, `dirichlet`

```
1 <?xml version="1.0"?>
2 <markovChain>
3 ...
4   <physics>
5   ...
6     <action>
7     ...
8       <quark>
9         <npCloverQuarkAction>
10        ...
11        <quarkField>
12          <normalisation>sqrt2kappa</normalisation>
13          <boundaryCondition>
14            <elem>periodic</elem>
15            <elem>periodic</elem>
16            <elem>periodic</elem>
17            <elem>antiperiodic</elem>
18          </boundaryCondition>
19        </quarkField>
20        ...
21      </npCloverQuarkAction>
22      <npCloverQuarkAction>
23      ...
24      </npCloverQuarkAction>
25      ...
26    </action>
27    ...
28  </physics>
29  ...
```

# Ensemble XML: physics / Fermionic Part (3/3)

## Example for linkSmearing:

- The element `linkSmearing` must contain the following sub-elements
  - An element defining the link blocking, e.g.
    - `apeLinkBlocking`
    - `anisotropicApeLinkBlocking`
  - An element defining the link unitarization, e.g.
    - `stoutLinkUnitarization`
    - `invSqRootLinkUnitarization`
  - `numSmear`: integer value

```
1  <?xml version="1.0"?>
2  <markovChain>
3  ...
4  <physics>
5  ...
6  <action>
7  ...
8  <quark>
9   <npCloverQuarkAction>
10  ...
11  <linkSmearing>
12   <apeLinkBlocking>
13     <rho0>0.1</rho0>
14     <rho1>0.1</rho1>
15   </apeLinkBlocking>
16   <stoutLinkUnitarization/>
17   <numSmear>1</numSmear>
18 </linkSmearing>
19 ...
20 </npCloverQuarkAction>
21 <npCloverQuarkAction>
22 ...
23 </npCloverQuarkAction>
24 ...
25 </action>
26 ...
27 </physics>
28 ...
29 </markovChain>
```

# Ensemble XML: algorithm

- The element `algorithm` must contain the following sub-elements

- name: A string restricted to legal XML 1.0 names
- glossary: URI to documentation
- reference: String containing, e.g., a publication
- exact: Boolean value
- parameters (optional)
- Option to insert other XML elements using a different namespace

- parameters** is a list of pairs

- name: A string restricted to legal XML 1.0 names
- value: Any simple type

## Example:

```
1 <?xml version="1.0"?>
2 <markovChain>
3 ...
4   <algorithm>
5     <name>qcdsfAcceleratedHMC</name>
6     <glossary>http://www-zeuthen.desy.de/latfor/ldg/algorithmGloss...
7     <reference/>
8     <exact>true</exact>
9     <parameters>
10    <name>timeScaleRatio</name>
11    <value>3</value>
12  </parameters>
13 </algorithm>
14 </markovChain>
```

# Content

1. Overview and Requirements
2. Digression: XML Technologies
3. Ensemble Metadata
- 4. Configuration Metadata**
5. Concluding Remarks

# Configuration XML Overview

Top-level elements of the configuration XML file:

- management
- implementation
- algorithm
- precision
- markovStep

## Example:

```
1 <?xml version="1.0"?>
2 <gaugeConfiguration>
3   <management>
4     ...
5   </management>
6   <implementation>
7     ...
8   <algorithm>
9     ...
10  </algorithm>
11  <precision>...</precision>
12  <markovStep>
13  ...
14  </markovStep>
15 </gaugeConfiguration>
```

# Configuration XML: management

- CRC32 checksum of the ildg-binary-data part of an ILDG configuration based on the LIME format
- Archive history
  - List of entries for tracking data provenance
  - Each entry comprises
    - Action: generate, add, replace, remove
    - Participant identified by name and institution
    - Date

## Example:

```
1 <?xml version="1.0"?>
2 <gaugeConfiguration>
3   <management>
4     <crcCheckSum>2266949024</crcCheckSum>
5     <archiveHistory>
6       <elem>
7         <revisionAction>generate</revisionAction>
8         <participant>
9           <name>Hinnerk Stueben</name>
10          <institution>ZIB</institution>
11        </participant>
12        <date>2004-11-15T19:33:55+01:00</date>
13      </elem>
14    </archiveHistory>
15  </management>
16  <implementation>
17    ...
18  </implementation>
19  <algorithm>
20    ...
21  </algorithm>
22  <precision>...</precision>
23  <markovStep>
24    ...
25  </markovStep>
26 </gaugeConfiguration>
```

# Configuration XML: implementation

- machine
  - name: String referring to the name of the computer used to generate the configuration
  - institution: String with the name of the organisation hosting that computer
  - machineType: String showing the computer type
  - comment: String with optional additional information
- code
  - name: String with the name of the code
  - version: String with the code version
  - date: Date of the code release
  - comment: String with optional additional information

## Example:

```
1 <?xml version="1.0"?>
2 <gaugeConfiguration>
3   <management>
4     ...
5   </management>
6   <implementation>
7     <machine>
8       <name>sr8000</name>
9       <institution>LRZ Munich</institution>
10      <machineType>Hitachi SR8000-F1</machineType>
11    </machine>
12    <code>
13      <name>BQCD</name>
14      <version>3.0.1</version>
15      <date>2003-11-28T13:00:00+01:00</date>
16    </code>
17  </implementation>
18  <algorithm>
19    ...
20  </algorithm>
21  <precision>...</precision>
22  <markovStep>
23    ...
24  </markovStep>
25 </gaugeConfiguration>
```

# Configuration XML: algorithm

- The optional element `parameters` comprises a list of pairs
  - name: A string restricted to legal XML 1.0 names
  - value: Any simple type
- It is allowed to insert other XML elements using a different namespace

## Example:

```
1 <?xml version="1.0"?>
2 <gaugeConfiguration>
3   <management>
4   ...
5   </management>
6   <implementation>
7   ...
8   </implementation>
9   <algorithm>
10  <parameters>
11    <name>targetResidue</name>
12    <value>1e-07</value>
13  </parameters>
14  </algorithm>
15  <precision>...</precision>
16  <markovStep>
17  ...
18  </markovStep>
19 </gaugeConfiguration>
```

# Configuration XML: precision

- Allowed values:
  - single: Single precision calculations
  - double: Double precision calculations
  - mixed: Mix of single and double precision calculations
- Note that configurations may be stored in different precision

## Example:

```
1 <?xml version="1.0"?>
2 <gaugeConfiguration>
3   <management>
4     ...
5   </management>
6   <implementation>
7     ...
8   </implementation>
9   <algorithm>
10    ...
11  </algorithm>
12  <precision>double</precision>
13  <markovStep>
14    ...
15  </markovStep>
16 </gaugeConfiguration>
```

# Configuration XML: markovStep

- **markovChainURI**: Reference to the ensemble
- **series**: Whitespace-replaced and collapsed string identifying a run
- **update**: Any simple type that references an update within a run
- **avePlaquette**: Average plaquette value (double)
- **dataLFN**: URI that is a globally unique reference to a configuration

## Example:

```
1 <?xml version="1.0"?>
2 <gaugeConfiguration>
3   <management>
4     ...
5   </management>
6   <implementation>
7     ...
8   </implementation>
9   <algorithm>
10    ...
11 </algorithm>
12 <precision>...</precision>
13 <markovStep>
14   <markovChainURI>mc://...</markovChainURI>
15   <series>563</series>
16   <update>1310</update>
17   <avePlaquette>0.5610635491</avePlaquette>
18   <dataLFN>lfn://reg/col/prj/...</dataLFN>
19 </markovStep>
20 </gaugeConfiguration>
```

# Content

1. Overview and Requirements
2. Digression: XML Technologies
3. Ensemble Metadata
4. Configuration Metadata
5. Concluding Remarks

# Concluding Remarks

- The ILDG metadata standards have been a robust framework for describing LQCD configurations data
  - Forward compatibility has been important
  - But: Evolution of the standard stagnated during the last 10 years
- The complexity of the underlying metadata schema is in practice not exposed to end-users
  - Most metadata documents are derived from existing documents with small modifications
- The restart with ILDG 2.0 is an opportunity to identify gaps in terms of
  - Training and tools for generating metadata
  - Coverage of data that can be described by ILDG-compliant metadata