



## WHEAT DATA INTEROPERABILITY GUIDELINES

**Authors:** Dzalé Yeumo Esther<sup>1</sup>, Fulss Richard<sup>2</sup>, Alaux Michael<sup>1</sup>, Aubin Sophie<sup>1</sup>, Arnaud Elizabeth<sup>3</sup>, Baumann Ute<sup>4</sup>, Cooper Laurel<sup>5</sup>, Davey Robert P.<sup>12</sup>, Laporte Marie-Angélique<sup>3</sup>, Larmande Pierre<sup>6</sup>, Letellier Thomas<sup>1</sup>, Pommier Cyril<sup>1</sup>, Protonotarios Vassilis<sup>7</sup>, Shrestha Rosemary<sup>2</sup>, Subirats Imma<sup>8</sup>, Venkatesan Aravind<sup>6</sup>, Whan Alex<sup>9</sup>, Wheat Data Interoperability working group<sup>11</sup>

**Contributors:** Lucas Hélène<sup>1</sup>, Quesneville Hadi<sup>1</sup>, Hologne Odile<sup>1</sup>, Keizer Johannes<sup>8</sup>, Mohellibi Nacer<sup>1</sup>, Buche Patrice<sup>1</sup>, Jonquet Clément<sup>10</sup>, Caracciolo Caterina<sup>8</sup>, Reverte Carme<sup>8</sup>

**Affiliations:** <sup>1</sup> L'Institut National de la Recherche Agronomique (INRA), <sup>2</sup> International Maize and Wheat Improvement Center (CIMMYT), <sup>3</sup> Bioversity International, <sup>4</sup> Australian Centre for Plant Functional Genomics (ACPGF), <sup>5</sup> Planteome, <sup>6</sup> L'Institut de Recherche pour le Développement, <sup>7</sup> Agroknow, <sup>8</sup> Food and Agriculture Organization (FAO), <sup>9</sup> Commonwealth Scientific and Industrial Research Organisation (CSIRO), <sup>10</sup> Laboratoire d'Informatique, de Robotique et de Microélectronique (LIRMM), <sup>11</sup> Research Data Alliance (RDA), <sup>12</sup> Earlham Institute (EI)

### ABSTRACT

[The International Wheat Initiative](#) has identified the easy access and interoperability of all wheat related data as a top priority for the wheat research community. An important goal is to make the best use of existing genetic, genomic, and phenotypic data in fundamental and applied wheat science. Hence, data interoperability has become a priority in this community, given the ever-growing data deluge coming from improvements in technologies and numeric methods for DNA (Deoxyribonucleic acid) and RNA (Ribonucleic acid) sequencing, high throughput genotyping and phenotyping, high throughput imaging and satellite monitoring.

The Wheat Data Interoperability working group (WDI-WG) was created in 2014 as a [Research Data Alliance \(RDA\)](#) working group under the umbrella of [Wheat Information System \(WheatIS\)](#) expert working group, which is endorsed by the International Wheat Initiative to build an international information system for genetic, genomic, and phenomic data. The primary mission of the WDI-WG is to provide a common framework for describing and representing data with respect to existing open standards. The WDI-WG is formed by data and information practitioners and scientists from different organizations and countries. The standpoint of the WDI-WG was to avoid the creation of new standards. Therefore, the WDI-WG started by surveying the practices of the wheat research community, then converged and agreed on specific recommendations. The proposed recommendations have been endorsed by the RDA and the WheatIS expert working group.

### SUMMARY OF THE GUIDELINES

→ Please visit the [website](#) for the full and updated version of the guidelines.

| Data types | Data exchange formats          | Metadata                       | Ontologies, vocabularies | Additional recommendations       |
|------------|--------------------------------|--------------------------------|--------------------------|----------------------------------|
| DNA        | The Variant Call Format (VCF). | Provide associated metadata to |                          | The currently most commonly used |

|                     |  |  |   |   |
|---------------------|--|--|---|---|
| sequence variations |  | contextualize the provenance of the SNPs and provide information about the SNP quality analysis.   |   | reference bread wheat sequence is the IWGSC (International Wheat Genome Sequencing Consortium) survey sequence (cv Chinese Spring), which is available at the <a href="#">IWGSC Sequence Repository</a> and <a href="#">The European Bioinformatics Institute (EBI)</a> .<br>When available, we encourage the use of the chromosome reference sequence. |
| Genome annotations  | The General Feature Format (GFF3).   |  | Use ontologies (Gene Ontology, Sequence Ontology) for functional annotation in column 9 in the GFF3 file.   | Provide comprehensive content description for column 9 in the GFF3 file.<br>Use external database cross references (Dbxref) in a consistent way.  |
| Gene sequences      | Use the existing format standards laid out by the repositories such as The <a href="#">National Center for Biotechnology Information (NCBI) (Gene Expression Omnibus, GEO)</a> , <a href="#">EBI Array Express</a> , and <a href="#">The European Nucleotide Archive (ENA)</a> . | Metadata is important for all gene expression studies – whether microarray or RNA Seq data. <ul style="list-style-type: none"> <li>• Microarray data: <a href="#">NCBI-GEO-MIAMI compliant</a></li> <li>• RNA-Seq data – follow the guidelines for Minimum Information About a Microarray Experiment (MIAME) as far as applicable (FastQ)</li> </ul> For <a href="#">BAM</a> files – additional info needed: <ul style="list-style-type: none"> <li>• mapping software</li> <li>• mismatch settings</li> <li>• reference sequences used such as IWGSC survey sequences or MIPS (The Munich Information Center for Protein Sequences) gene models, or transcriptome assembly</li> </ul> | Use ontologies and controlled vocabularies to annotate the required metadata <ul style="list-style-type: none"> <li>• <a href="#">Plant Ontology (PO)</a> terms to describe the plant tissues and developmental stage.</li> <li>• <a href="#">Plant Environment Ontology (EO)</a> terms to describe the experimental conditions.</li> <li>• <a href="#">Plant Stress Ontology (PSO)</a> terms to describe the treatments with pathogens, stress conditions (proposed).</li> <li>• <a href="#">Gene Ontology (GO)</a> is the standard for the functional analysis.</li> <li>• <a href="#">Microarray and Gene Expression Data Ontology (MGED) Ontology (MO)</a> terms mapped to The <a href="#">Ontology for Biomedical Investigations (OBI)</a> / <a href="#">The Open Biomedical Ontologies (OBO)</a> foundry ontology terms.</li> </ul> |   |
| Phenotype           | Use data format that follows   | Provide complete metadata for at least   | Use following ontologies and controlled   | We recommend sharing at least clean   |

|                |   |   |   |   |
|----------------|---|---|---|---|
| data           | minimum format principles with data matrices plus metadata for at least variables (trait along with method, units and scales or environmental ones) and germplasms. <a href="#">ISA-Tab</a> , which is compliant to the <a href="#">Minimum Information About a Plant Phenotyping Experiment (MIAPPE)</a> standard, is an implementation of this principle. | germplasm and observation variables.  | vocabularies to annotate nurseries and trials data: <ul style="list-style-type: none"> <li>• <a href="#">Crop Research Ontology (CO 715)</a> to describe the metadata related to nurseries and trials, field management, field environments, study design, etc. These metadata are actively adapted for a wider use.</li> <li>• <a href="#">Wheat Trait Ontology (CO 321)</a> describes the variables that are most commonly used in wheat phenotyping.</li> <li>• <a href="#">Biorefinery Ontology</a> describes the concepts and terms associated with biomass composition and characterization (crystallinity, surface area, particle size, porosity, etc.), physico-chemical pretreatments, enzymatic hydrolysis, and experimental processes descriptions.</li> </ul> | documented raw data, like plant height, leaf area, etc.   |
| Germplasm data | Use data matrices in csv (comma-separated values) or Excel. For examples: <a href="#">GnplS</a> , <a href="#">Breeding Management System (BMS)</a> .  | Use <a href="#">Multi-Crop Passport Descriptor (MCPD)</a> as metadata format for passport and other variables data required for characterization of germplasm | <a href="#">Wheat Trait Ontology (CO 321)</a> to describe germplasm characterization  | Identification of germplasm must be non-ambiguous at the scale of the world wide web. To achieve this, consider using URI and mainly follow the work of the <a href="#">DivSeek</a> project.  |
| Physical maps  | FPC format for physical map raw data<br>GFF3 format for data integration  |   |   | <b>Data submission:</b> INRA <a href="#">URGI</a> is the official repository for <a href="#">IWGSC</a> physical maps. We recommend using the data submission support at <a href="http://wheat-urgi.versailles.inra.fr/Seq-Repository/Support-to-assembly-and-data-submission">http://wheat-urgi.versailles.inra.fr/Seq-Repository/Support-to-assembly-and-data-submission</a> |

➔ In addition, the WDI-WG uses [AgroPortal](#) to maintain the list of the ontologies and vocabularies recommended by the working group. The WDI-WG encourages the developers of wheat related ontologies and vocabularies to share new resources through the same portal in order to facilitate their findability and accessibility by the community. AgroPortal allows users to search within ontologies, retrieve existing mappings between ontologies, visualize ontology content, comment about ontologies, concepts and mappings, annotate text data with ontology concepts, get a recommendation about which ontology to use for a text corpus or set of keywords, and store projects which use ontologies.

## MAINTENANCE OF THE GUIDELINES

The RDA WDI working group has now become a maintenance group since June 2017. The working group will maintain the guidelines in collaboration with the WheatIS expert working group and a steering group including representatives of the adopters. The guidelines will be revised twice (frequently) a year to consider the evolving landscape of data standards and practices.

## GETTING INVOLVED

There are many ways you can be part of this work:

- ➔ Give us your feedback: you can leave a comment on a specific page of the guidelines website in order to propose improvements or corrections and share best practices or useful tools.
- ➔ Join the maintenance group: the maintenance group is open and all contributions are welcome. You can freely join us at <https://www.rd-alliance.org/groups/wheat-data-interopability-wg.html>
- ➔ Become an adopter: Please let us know if you wish to adopt the guidelines by contacting us at [rda-wdinterop-wg@rda-groups.org](mailto:rda-wdinterop-wg@rda-groups.org)

### List of early adopters

| Organization   | Contact   |
|--|---|
| <a href="#">Australian Center for Plant Functional Genomics</a>  | Dr. Baumann, Ute, Bioinformatics Leader   |
| <a href="#">Earlham Institute,</a>   | Dr. Robert Davey, Data Infrastructure & Algorithms Group Leader   |
| <a href="#">EMBL European Bioinformatics Institute</a>   | Dr. Paul Kersey, Team Leader Non-vertebrate Genomics  |
| <a href="#">L'Institut National de la Recherche Agronomique (INRA), Unité de Recherche Génomique Info (URGI)</a> | Dr. Michael Alaux, Deputy leader of "Information System and data integration" team<br>Dr. Cyril Pommier, Deputy leader, Information System and Data integration team, Phenotype thematic leader |
| <a href="#">International Maize and Wheat Improvement Center (CIMMYT)</a>  | Richard Fulss, Head of Knowledge Management<br>Dr. Rosemary Shrestha, Data Coordinator  |
| <a href="#">James Hutton Institute</a>   | Dr. David Marshall, Information and Computational Sciences  |
| <a href="#">Munich Information Center for Protein Sequences (MIPS), Helmholtz Center Munich</a>                  | Dr. Klaus Mayer, Research Director MIPS   |
| <a href="#">Rothamsted Research</a>  | Dr. Christopher Rawlings, Head of Department Computational & Systems Biology Harpenden  |
| <a href="#">The National Institute of Agricultural Botany (NIAB)</a>   | Professor Dr. Mario Caccamo, Head of Crop Bioinformatics  |
| <a href="#">USDA ARS and Cold Spring Harbor Laboratory</a>   | Dr. Doreen Ware, Adjunct Associate Professor, Ohio State University   |