

# **Sample and Data Relationship Format for Proteomics (SDRF- Proteomics)**

Version 1.1.0-dev, 2026-01-18

# Table of Contents

<b>1. Status of this document</b>	1
<b>2. Abstract</b>	2
<b>3. Motivation</b>	3
<b>4. Quick Start</b>	5
4.1. Minimal Example .....	5
4.2. Key Concepts .....	5
4.3. Format Requirements .....	5
4.4. Scope .....	6
4.5. Getting Started Steps .....	6
<b>5. Validating SDRF Files</b>	7
<b>6. Specification structure</b>	8
6.1. Versioning .....	8
6.2. Notational Conventions .....	9
6.3. Relationship to other specifications .....	9
<b>7. SDRF-Proteomics specification</b>	10
7.1. Format rules .....	10
7.2. SDRF file-level metadata .....	11
7.3. Table Column headers .....	12
7.4. Table Cell values .....	13
<b>8. SDRF-Proteomics: Samples metadata</b>	14
8.1. BioSamples database integration .....	15
8.2. Encoding sample technical and biological replicates .....	15
8.3. Pooled samples .....	17
8.3.1. Allowed values for characteristics[pooled sample] .....	17
8.3.2. Example with simple pooled annotation .....	17
8.3.3. Example with detailed pooled reference .....	18
8.4. Spiked-in samples .....	18
8.5. Sample Metadata Guidelines .....	19
<b>9. SDRF-Proteomics: data files metadata</b>	21
9.1. CV Term Format for Data File Metadata .....	21
9.2. Sample Preparation and Fragmentation .....	22
9.3. Data acquisition .....	22
9.4. Data File Metadata Guidelines .....	23
<b>10. Additional SDRF Rules</b>	24
10.1. Row Uniqueness Requirements .....	24
<b>11. Ontologies and Controlled Vocabularies</b>	25
<b>12. Core Templates</b>	27
12.1. What is a Template? .....	27
12.2. Template Inheritance and Composition .....	27
12.3. YAML Template Definitions .....	27
12.4. Choosing a Template .....	28
12.5. Default Template .....	29

12.6. Human Template .....	30
12.7. Vertebrates Template .....	31
12.8. Invertebrates Template.....	32
12.9. Plants Template .....	33
12.10. Experiment-Type Templates.....	34
12.11. Column Cardinality.....	34
12.12. Extending SDRF with Custom Columns .....	35
12.12.1. When to Add Custom Columns .....	35
12.12.2. Rules for Custom Columns .....	35
12.12.3. Common Additional Columns .....	36
12.12.4. Example: Adding Study-Specific Columns.....	36
<b>13. Factor Values (Study Variables)</b>	<b>38</b>
13.1. Column Format .....	38
13.2. When to Use Factor Values .....	38
13.3. Rules .....	38
13.4. Example.....	38
<b>14. Examples of Annotated Datasets</b>	<b>39</b>
<b>15. Ongoing template discussions</b>	<b>40</b>
<b>16. Intellectual Property Statement</b>	<b>41</b>
<b>17. Copyright Notice</b>	<b>42</b>
<b>18. How to cite</b>	<b>43</b>
<b>References</b>	<b>44</b>

# Chapter 1. Status of this document

This document provides information to the proteomics community about a proposed standard for sample metadata annotations in public repositories called Sample and Data Relationship Format (SDRF)-Proteomics. Distribution is unlimited.

**Version v1.1.0 - 2025-01**

## Chapter 2. Abstract

The Human Proteome Organisation (HUPO) Proteomics Standards Initiative (PSI) defines community standards for data representation in proteomics to facilitate data comparison, exchange, and verification. This document presents a specification for the Sample and Data Relationship Format (SDRF-Proteomics).

Further detailed information, including any updates to this document, implementations, and examples is available at [SDRF GitHub Repository](#). The official PSI web page for the document is: [HUPO-PSI SDRF](#).

## Chapter 3. Motivation

Many resources have emerged that provide raw or integrated proteomics data in the public domain. If these are valuable individually, their integration through re-analysis represents a huge asset for the community [1].

Unfortunately, proteomics experimental design and sample related information are often missing in public repositories or stored in very diverse ways and formats. For example:

- The [CPTAC Consortium](#) provides for every dataset a set of Excel files with the information on [each sample](#) including tumor size, origin, but also how every sample is related to a specific raw file (e.g. instrument configuration parameters).
- As a resource routinely re-analysing public datasets, [ProteomicsDB](#) captures for each sample in the database a minimum number of properties to describe the sample and the related experimental protocol such as tissue, digestion method and instrument.

Such heterogeneity often prevents data interpretation, reproducibility, and integration of data from different resources. For every proteomics dataset we propose to capture at least three levels of metadata:

- (i) dataset description
- (ii) the sample metadata and data files acquisition metadata.
- (iii) The relation between the sample and the data files. The experimental design.

The general description includes minimum information to describe the study overall: [title](#), [description](#), [date of publication](#), [type of experiment](#). In ProteomeXchange partners this metadata is captured at the dataset level, in other omics resources this is captured as IDF file format (e.g. MAGE-TAB). Currently, all ProteomeXchange partners mandate this information for each dataset. However, the information regarding the sample and its relation to the data files (**Figure 1**) is mostly missing [1].



**Figure 1:** SDRF-Proteomics file format stores the information of the sample and its relation to the data files in the dataset. The file format includes not only information about the sample but also about how the data was acquired and processed.

Here, we introduced the Sample and Data Relationship Format (SDRF-Proteomics) to capture the sample metadata and its relation to the data files for proteomics experiments. The SDRF-Proteomics format is a tab-delimited file format that describes the sample characteristics and the relationships between samples and data files included in a dataset.

This specification, which is a community effort, aims to provide a standard for the proteomics community to annotate the sample metadata and its relation to the data files.

# Chapter 4. Quick Start

If you're new to SDRF-Proteomics, here's a minimal example to get you started. An SDRF file is a tab-separated file where each row represents a sample-to-data-file relationship.

## 4.1. Minimal Example

source name	characteristics[organism]	characteristics[disease]	...	assay name	comment[label]	comment[instrument]	comment[data file]	...	factor value[disease]
sample_1	homo sapiens	normal	...	run_1	label free sample	Q Exactive HF	sample_1.raw	...	normal
sample_2	homo sapiens	hepatocellular carcinoma	...	run_2	label free sample	Q Exactive HF	sample_2.raw	...	hepatocellular carcinoma

*Blue: Sample metadata | Green: Data file metadata | Orange: Factor values*

## 4.2. Key Concepts

1. **Sample metadata** uses `characteristics[...]` columns (e.g., organism, disease)
2. **Data file metadata** uses `comment[...]` columns (e.g., instrument, label)
3. **Factor values** use `factor value[...]` columns to indicate variables under study
4. Each row links one sample to one data file

## 4.3. Format Requirements

The SDRF-Proteomics format has the following core requirements:

- The SDRF file is a **tab-delimited format** where each row corresponds to a relationship between a Sample and a Data file.
- Each column MUST correspond to an attribute/property of the Sample or the Data file.
- Each cell value MUST be the property value for the corresponding Sample or Data file.
- The file MUST start with columns describing sample properties (e.g., organism, disease), followed by data file properties (e.g., label, fraction identifier, data file).
- Unknown values MUST be handled using `not available` (value is unknown), `not applicable` (property doesn't apply), or `pooled` (value is a mixture from multiple samples).

## 4.4. Scope

The SDRF-Proteomics format aims to capture the **sample metadata** and its **relationship with data files** (e.g., raw files from mass spectrometers).

**IMPORTANT**

SDRF-Proteomics does **not** aim to capture downstream analysis details, including: which samples were compared to which other samples, how samples are combined into study variables, or analysis parameters such as FDR thresholds or p-value cutoffs.

## 4.5. Getting Started Steps

1. Choose a [core template](#) (Human, Vertebrates, Plants, etc.)
2. Fill in sample metadata (characteristics columns)
3. Fill in data file metadata (comment columns)
4. Add factor values for your experimental variables
5. Validate your file using [sdrf-pipelines](#)

For detailed guidance, continue reading the full specification below.

## Chapter 5. Validating SDRF Files

The official validator for SDRF-Proteomics files is **sdrf-pipelines**, a Python tool that checks your SDRF file for errors and compliance with the specification.

Installation:

```
pip install sdrf-pipelines
```

Basic Validation:

```
# Validate an SDRF file
parse_sdrf validate-sdrf --sdrf_file your_file.sdrf.tsv

# Validate with a specific template
parse_sdrf validate-sdrf --sdrf_file your_file.sdrf.tsv --template human
```

For more information, visit: [sdrf-pipelines on GitHub](#)

# Chapter 6. Specification structure

This document describes the main specification of SDRF-Proteomics, the structure of the specification (**Figure 2**), how to contribute, and extend the specification. SDRF-Proteomics uses a three-tier system for organizing metadata requirements:

- **The SDRF-Proteomics core specification:** This document contains the main specification, requirements and rules for the SDRF-Proteomics format. It also includes the notational conventions and the relationship to other specifications.
- **Core templates:** Organism-based templates (human, vertebrates, plants, etc.) that define base schemas for common proteomics experiments. See the [Core Templates](#) section.
- **Specialized templates:** Complete schemas for specific experiment types and acquisition methods (DDA acquisition, DIA acquisition, cell-lines, single-cell, affinity-proteomics, crosslinking, immunopeptidomics, metaproteomics). Each template has its own directory containing:
  - A detailed README.adoc with checklists, and examples.
  - A template file ({name}-template.sdrf.tsv) with column headers.
- **Annotation guidelines:** Detailed documentation for specific metadata annotations (e.g., patient pre-existing condition, sample metadata, data file metadata).

[Logo] | [images/sdrf-guidelines-structure.png](#)

**Figure 2:** SDRF-Proteomics specification structure. The main specification defines the core rules and is extended by specific experiment templates and annotation guidelines.

**NOTE**

The main specification is in the `sdrf-proteomics` directory. Core templates (organism-based) are in `sdrf-proteomics/core-templates/` and specialized templates (experiment-type-specific) are in `sdrf-proteomics/templates/`. Templates are extensions of the core specification, and should follow all the rules and requirements in the main specification. If a template rule is in conflict with the specification, a note should be done in the main specification to reflect the extension or conflict.

The official website for SDRF-Proteomics project is <https://github.com/bigbio/proteomics-metadata-standard>. New use cases, changes to the specification and examples can be added by using Pull requests or issues in GitHub to reach the bigbio team.

A set of examples and annotated projects from ProteomeXchange can be found here: <https://github.com/bigbio/proteomics-metadata-standard/tree/master/annotated-projects>

Multiple tools have been implemented to validate, annotate and convert SDRF-Proteomics files. The official validator of SDRF-Proteomics is sdrf-pipelines (Python - <https://github.com/bigbio/sdrf-pipelines>). This tool allows to validate an SDRF-Proteomics file. In addition, it allows converting SDRF to other popular pipelines and software configure files such as MaxQuant or OpenMS.

## 6.1. Versioning

The SDRF-Proteomics specification is versioned using the Semantic Versioning 2.0.0

(<https://semver.org/>) scheme. The version number is in the format MAJOR.MINOR.PATCH, where:

- MAJOR version is incremented for incompatible changes to the specification, when major changes are done to the specification.
- MINOR version is incremented for new features that are backward compatible with the previous version. Guidelines and templates are added or modified.
- PATCH version is incremented for bug fixes and minor changes that do not affect the specification or the templates. This includes typos, formatting changes, and other minor updates.

Every change in the specification should be done in GitHub using pull requests into the dev branch. The pull request should include a description of the changes and the reason for the changes. The pull request will be reviewed by the community and merged into the main branch when approved. After the merge, the version number will be updated according to the changes made, the release will be performed, and the Zenodo record will be updated.

**NOTE**

We added the prefix v to the version number to indicate that it is the version of the specification that was used to create the file. Examples: v1.1.0, v2.0.0, v3.0.0.

## 6.2. Notational Conventions

The key words “MUST”, “MUST NOT”, “REQUIRED”, “SHALL”, “SHALL NOT”, “SHOULD”, “SHOULD NOT”, “RECOMMEND/RECOMMENDED”, “MAY”, “COULD BE”, and “OPTIONAL” are to be interpreted as described in RFC 2119 (<https://www.rfc-editor.org/rfc/rfc2119>).

## 6.3. Relationship to other specifications

SDRF-Proteomics is fully compatible with the SDRF file format part of [MAGE-TAB](#). MAGE-TAB is the file format used to store metadata and sample information for transcriptomics experiments. When the ProteomeXchange project file is converted to idf file (project description in MAGE-TAB) and is combined with the SDRF-Proteomics a valid MAGE-TAB is obtained.

SDRF-Proteomics sample information can be embedded into mzTab metadata files. The sample metadata in mzTab contains properties as the columns in the SDRF-Proteomics and values as Sample cell values.

The SDRF-Proteomics aims to capture the sample metadata and its relationship with the data files (e.g. raw files from mass spectrometers). The SDRF-Proteomics do not aim to capture the downstream analysis part of the experimental design such as what samples should be compared, how they can be combined or parameters for the downstream analysis (FDR or p-values thresholds). The HUPO-PSI community will work in the future to include this information in other file formats such as mzTab or a new type of file format.

## Chapter 7. SDRF-Proteomics specification

The SDRF-Proteomics file format describes the sample characteristics and the relationships between samples and data files. The file format is a tab-delimited one where each **ROW** corresponds to a relationship between a Sample and a Data file (in an ms proteomics experiment the data file containing the mass spectra), each **COLUMN** corresponds to an attribute/property of the Sample, the Data file, or the Factor values; and the value in each **CELL** is the specific value of the property for a given Sample/Data file/Factor value (**Figure 3**).

source name	characteristics[organism]	characteristics[disease]	assay name	comment[instrument]	comment[data file]	factor value[disease]
sample_1	homo sapiens	normal	run_1	Q Exactive HF	sample_1.raw	normal
sample_2	homo sapiens	liver cancer	run_2	Q Exactive HF	sample_2.raw	liver cancer

*Blue: Sample metadata | Green: Data file metadata | Orange: Factor values*

**Figure 3:** SDRF-Proteomics in a nutshell. Each **row** links a sample to a data file. **Columns** represent sample properties (characteristics), data file properties (comments), or experimental variables (factor values).

The SDRF-Proteomics format contains three main sections:

- The first section contains the [sample metadata](#).
- The second section contains the [data file metadata](#).
- The third section contains the [factor values](#) properties.

### 7.1. Format rules

There are general scenarios/use cases that are addressed by the following rules:

- **Unknown values:** In some cases, the column is mandatory in the format, but for some samples the corresponding value is unknown. In those cases, users **SHOULD** use **not available**.
- **Not Applicable values:** In some cases, the column is mandatory, but for some samples the corresponding value is not applicable. In those cases, users **SHOULD** use **not applicable**.
- **Pooled values:** In some cases, the sample is a pool of multiple samples (e.g., TMT reference channels), and the value cannot be represented as a single value. In those cases, users **SHOULD** use **pooled**.

*Table 1. Special values for SDRF cells*

Term	Meaning	Example	Use Case
not available	Value exists but is unknown or could not be determined	characteristics[age] = not available	Patient age was not recorded in the study
not applicable	Value or concept does not apply to this sample	characteristics[age] = not applicable	Synthetic peptide library has no age
pooled	Value represents a mixture of multiple samples	characteristics[biological replicate] = pooled	TMT reference channel pooled from multiple replicates

- **Case sensitivity:** Text values are case-insensitive, but **column names are case-sensitive**. Use lowercase for all column names (e.g., `source name`, `characteristics[organism]`, `comment[label]`). Incorrect casing like `Source Name` or `Characteristics[organism]` will cause validation failures.
- **Space sensitivity:** The SDRF is sensitive to spaces in column names (`sourcename ↪ source name`). Column names must include appropriate spaces (e.g., `source name`, not `sourcename`) but must NOT have a space before the bracket (e.g., `characteristics[organism]`, not `characteristics [organism]`).
- **Column order:** The SDRF columns follows some structure; first the sample metadata columns in [Chapter 8](#); then the data file metadata columns in [Chapter 9](#); followed by the factor values columns in [Chapter 13](#).
- **Extension:** The extension of the SDRF file SHOULD be `sdrf.tsv` (preferred) or `.txt`.

## 7.2. SDRF file-level metadata

Since version 1.1.0, SDRF-Proteomics supports optional file-level metadata using header comments at the beginning of the file. These header comments provide information about the SDRF file itself, such as the format version, template used, and validation status. This approach is inspired by other omics formats such as VCF (Variant Call Format) file headers and is fully compatible with pandas and other tabular data processing tools.

Header comments MUST:

- Start with `#` (single hash) followed by a key-value pair
- Appear at the very beginning of the file, before the column header row
- Use the format `#key=value`

The following header fields are supported:

Key	Description	Example	Requirement	Ontology Term
<code>file_format</code>	Identifier for the file format	SDRF	RECOMMENDED	PRIDE:0000831

Key	Description	Example	Requirement	Ontology Term
version	SDRF-Proteomics specification version used	v1.1.0	RECOMMENDED	NCIT:C25714
template	Name of the template used	human, cell_lines, default	OPTIONAL	PRIDE:0000832
template_version	Version of the template	v1.0.0	OPTIONAL	PRIDE:0000833
source	Origin or creator of the file	PRIDE, user-generated	OPTIONAL	NCIT:C25683
validation_hash	Hash from validator certification	sha256:abc123...	OPTIONAL	PRIDE:0000834

Example of an SDRF file with header comments (simplified example showing only select columns; see [Chapter 12](#) for complete required columns):

```
#file_format=SDRF
#version=v1.1.0
#template=human
#template_version=v1.0.0
#source=PRIDE
source name characteristics[organism]    characteristics[organism part]
characteristics[disease]    assay name comment[data file]
sample_1      homo sapiens      liver    normal   run_1    sample_1.raw
```

**NOTE**

Header comments are OPTIONAL. SDRF files without header comments are still valid. When present, header comments provide valuable provenance information and enable tools to handle version-specific features appropriately. Header property names use underscores (e.g., `file_format`, `template_version`) rather than spaces to maintain consistency with the tab-delimited nature of SDRF files and avoid ambiguity when parsing.

### 7.3. Table Column headers

Depending on each section the column headers (property names) will be prefixed with the following prefixes:

- `characteristics`: Sample metadata (e.g. `characteristics[organism]`)
- `comment`: Data file metadata (e.g. `comment[data file]`)
- `factor value`: Factor values properties (e.g. `factor value[disease]`)

Each property name MUST be a valid ontology term or a valid controlled vocabulary term. Each section will have some specific order for column headers.

**NOTE**

A list of all controlled vocabularies and ontologies supported are in the [Chapter 11](#)

section. On each section we also provide a list of properties that are supported.

## 7.4. Table Cell values

The value for each property, (e.g. characteristics, comment, factor value) corresponding to each sample or data file can be represented in multiple ways.

- Free Text (Human readable): In the free text representation, the value is provided as text without Ontology support (e.g. colon or providing accession numbers). This is only RECOMMENDED when the text inserted in the table is the exact name of an ontology/CV term in EFO. If the term is not in EFO, other ontologies can be used.

source name	characteristics[organism]
sample 1	homo sapiens
sample 2	homo sapiens

- Ontology url (Computer readable): Users can provide the corresponding URI (Uniform Resource Identifier) of the ontology/CV term as a value. This is recommended for enriched files where the user does not want to use intermediate tools to map from free text to ontology/CV terms.

source name	characteristics[organism]
Sample 1	<a href="http://purl.obolibrary.org/obo/NCBITaxon_9606">http://purl.obolibrary.org/obo/NCBITaxon_9606</a>
Sample 2	<a href="http://purl.obolibrary.org/obo/NCBITaxon_9606">http://purl.obolibrary.org/obo/NCBITaxon_9606</a>

- Key=value representation (Human and Computer readable): The current representation aims to provide a mechanism to represent the complete information of the ontology/CV term including Accession, Name and other additional properties. In the key=value pair representation, the Value of the property is represented as an Object with multiple properties, where the key is one of the properties of the object and the value is the corresponding value for the particular key. An example of key value pairs is post-translational modification (see [Protein Modifications](#)):

```
NT=Glu->pyro-Glu;MT=fixed;PP=Anywhere;AC=Unimod:27;TA=E
```

## Chapter 8. SDRF-Proteomics: Samples metadata

The Sample metadata section provides information about the samples of origin and their characteristics. Each sample contains a *source name* (unique identifier) and a set of *characteristics* columns. The first column of the file should be the *source name* and the following columns should be the characteristics of the sample. For example, for any proteomics experiment (human, vertebrate, cell line), the following characteristics should be provided:

- **source name:** Unique sample name (it can be present multiple times if the same sample is used several times in the same dataset)
- **characteristics[organism]:** The organism of the Sample of origin. Values MUST come from [NCBI Taxonomy](#).
- **characteristics[organism part]:** The part of organism's anatomy or substance arising from an organism from which the biomaterial was derived (e.g., liver). Values SHOULD come from [UBERON](#) or [BTO](#).
- **characteristics[disease]:** The disease under study in the Sample. Values SHOULD come from [MONDO](#), [EFO](#), or [DOID](#). For healthy/control samples, use [normal](#) ([PATO:0000461](#)) - see [Disease Annotation Guidelines](#).
- **characteristics[cell type]:** A cell type is a distinct morphological or functional form of cell (e.g., epithelial, glial). Values SHOULD come from [Cell Ontology \(CL\)](#) or [BTO](#).

Example:

source name	characteristics[organism]	characteristics[organism part]	characteristics[disease]	characteristics[cell type]
sample_treat	homo sapiens	liver	liver cancer	not available
sample_control	homo sapiens	liver	liver cancer	not available

**NOTE**

Additional characteristics can be added depending on the type of the experiment and sample. The [SDRF-Proteomics templates](#) defines a set of templates and checklists of properties that should be provided depending on the proteomics experiment. In the core guidelines and templates, main document of SDRF-Proteomics, we explain the major sample properties for different experiments. However, SDRF-Proteomics can be extended using guidelines for specific experiments.

Some important notes:

- Each characteristic name in the column header SHOULD be a CV term from the EFO ontology. For example, the header *characteristics[organism]* corresponds to the ontology term Organism. However the values could be from EFO or other ontologies. For example, we RECOMMEND to use MONDO for diseases because it has better coverage than EFO. For healthy samples, use [normal](#) ([PATO:0000461](#)) - see [Disease Annotation Guidelines](#).
- Multiple values (columns) for the same characteristics term are allowed in SDRF-Proteomics. However, it is RECOMMENDED not to use the same column in the same file. If you have multiple phenotypes, you can specify what it refers to or use another more specific term, e.g.,

"immunophenotype".

## 8.1. BioSamples database integration

BioSamples provides persistent identifiers for biological samples that enable cross-database linking [5]. Use the optional *characteristics[biosample accession number]* column to link samples to BioSamples entries (EBI format: SAMEA\*, NCBI format: SAMN\*).

**Column:** `characteristics[biosample accession number]`

**Examples:** `SAMN12345678` (NCBI), `SAMEA12345678` (EBI)

Example usage:

source name	characteristics[bio sample accession number]	characteristics[organism]	characteristics[organism part]	characteristics[disease]
sample_001	SAMN12345678	homo sapiens	liver	liver cancer
sample_002	SAMN12345679	homo sapiens	liver	normal
sample_003	SAMEA12345680	mus musculus	brain	normal

The BioSamples accession number enables connection to sample resources across databases, linking proteomics data with genomics, transcriptomics, and other omics datasets.

**NOTE**

BioSample accession numbers from NCBI follow the format `SAMNxxxxxxxxxx` and from EBI follow the format `SAMEAxxxxxxxxxx`, where `x` represents digits. Either NCBI or EBI BioSample accession numbers can be used depending on where the sample is registered. The *characteristics[biosample accession number]* column is optional, but when available, providing BioSample accession numbers is RECOMMENDED to enhance data integration and reusability. Users must first request BioSample accession numbers from the appropriate service (NCBI or EBI) before including them in their SDRF files.

## 8.2. Encoding sample technical and biological replicates

SDRF-Proteomics uses two columns to track replicates [4]:

- **characteristics[biological replicate]:** Identifies independent biological samples within each experimental condition (factor value group). Replicate numbers are assigned per condition, so if you have 2 cancer samples and 2 healthy samples, both groups would have biological replicates numbered 1 and 2.
- **comment[technical replicate]:** Identifies repeated measurements of the same sample (e.g., multiple LC-MS/MS injections)

**IMPORTANT**

Biological replicate numbering restarts for each experimental condition (factor

value). For example, in a disease study with factor value[disease], the cancer samples would be numbered 1, 2, 3... and the normal samples would independently be numbered 1, 2, 3... This establishes the relationship between each sample and its experimental condition.

The following example shows 2 biological replicates, each with 2 fractions and 2 technical replicates:

source name	characteristics[biological replicate]	assay name	comment[label]	comment[fraction identifier]	comment[technical replicate]	comment[data file]
patient_001_sample	1	run_01	label free sample	1	1	P001_F1_TR 1.raw
patient_001_sample	1	run_02	label free sample	2	1	P001_F2_TR 1.raw
patient_001_sample	1	run_03	label free sample	1	2	P001_F1_TR 2.raw
patient_001_sample	1	run_04	label free sample	2	2	P001_F2_TR 2.raw
patient_002_sample	2	run_05	label free sample	1	1	P002_F1_TR 1.raw
patient_002_sample	2	run_06	label free sample	2	1	P002_F2_TR 1.raw
patient_002_sample	2	run_07	label free sample	1	2	P002_F1_TR 2.raw
patient_002_sample	2	run_08	label free sample	2	2	P002_F2_TR 2.raw

In this example:

- Biological replicates:** `patient_001_sample` and `patient_002_sample` are different biological samples (different source names), annotated with `characteristics[biological replicate]` values 1 and 2
- Technical replicates:** Each biological sample is measured twice (`comment[technical replicate] = 1` and `2`)
- Fractions:** Each technical replicate has 2 fractions (`comment[fraction identifier] = 1` and `2`)

### IMPORTANT

Both `characteristics[biological replicate]` and `comment[technical replicate]` columns are REQUIRED. When no replicates are performed in a study, set both columns to 1 (i.e., each sample is biological replicate 1 and technical replicate 1). For **pooled reference samples** (e.g., TMT reference channels), use `pooled` for biological replicate since these samples are mixtures of multiple replicates and assigning a specific replicate number would be misleading.

Some examples with explicit annotation of the biological replicates can be found here:

- <https://github.com/bigbio/proteomics-metadata-standard/blob/c3a56b076ef381280dfcb0140d2520126ace53ff/annotated-projects/PXD006401/sdrf.tsv>

## 8.3. Pooled samples

When multiple samples are pooled into one, the general approach is to annotate them separately, abiding by the general rule: one row stands for one sample-to-file relationship. In this case, multiple rows are created for the corresponding data file, much like in multiplexed labeling experiments (see [Label Annotations](#)).

One possible exception is made for the case when one channel (e.g., in a TMT/iTRAQ multiplexed experiment) is used for a sample pooled from all other channels, typically for normalization purposes. In this case, it is not necessary to repeat all sample annotations. Instead, the *characteristics[pooled sample]* column SHOULD be used.

### 8.3.1. Allowed values for characteristics[pooled sample]

The *characteristics[pooled sample]* column accepts the following values:

Value	Description	When to Use
not pooled	Sample is not pooled, represents a single biological sample	Regular individual samples
pooled	Sample is pooled but individual source samples cannot be annotated	When pooling details are unknown or samples are from external sources
SN=sample1;SN=sample2;...	Structured format listing source names of pooled samples	When individual samples are known and annotated in the same SDRF file

**NOTE**

The `SN` key stands for "source name" and lists the `source_name` values of samples that are annotated in the same file and used in the same experiment and same MS run. Use semicolons to separate multiple entries.

### 8.3.2. Example with simple pooled annotation

When pooling details are unknown or samples come from external sources, use the simple `pooled` value:

source name	characteristics [pooled sample]	characteristics [organism]	assay name	comment[label]	comment[data file]
sample_1	not pooled	homo sapiens	run_1	TMT126	file01.raw

source name	characteristics [pooled sample]	characteristics [organism]	assay name	comment[label]	comment[data file]
sample_2	not pooled	homo sapiens	run_1	TMT127N	file01.raw
pooled_ref	pooled	homo sapiens	run_1	TMT131C	file01.raw

### 8.3.3. Example with detailed pooled reference

When pooled samples are known and annotated in the same SDRF file, use the `SN=` format:

source name	characteristics[pooled sample]	characteristics[organism]	characteristics[age]	characteristics[sex]	assay name	comment [label]	comment [data file]
sample_1	not pooled	homo sapiens	45Y	male	run_1	TMT126	file01.raw
sample_2	not pooled	homo sapiens	52Y	female	run_1	TMT127N	file01.raw
sample_3	not pooled	homo sapiens	38Y	male	run_1	TMT127C	file01.raw
pooled_ref	SN=sample_1;SN=sample_2;SN=sample_3	homo sapiens	pooled	pooled	run_1	TMT131C	file01.raw

**TIP**

For pooled reference samples (e.g., TMT reference channels), use `pooled` for individual-specific fields including **biological replicate**, age, sex, and individual. This clearly indicates that the value represents a mixture rather than a single sample. If all pooled samples share a value (e.g., all females, or age range 40Y-50Y), that shared value MAY be used instead.

## 8.4. Spiked-in samples

There are multiple scenarios when a sample is spiked with additional analytes. Peptides, proteins, or mixtures can be added to the sample as controlled amounts to provide a standard or ground truth for quantification, or for retention time alignment, etc.

To include information about the spiked compounds, use `characteristics[spiked compound]`. The information is provided in key-value pairs. Here are the keys and values that SHOULD be provided:

Key	Meaning	Examples	Peptide	Protein	Mixture	Other
CT	Compound type	protein, peptide, mixture, other	Required	Required	Required	Required

Key	Meaning	Examples	Peptide	Protein	Mixture	Other
QY	Quantity (molar or mass)	10 mg, 20 nmol	Required	Required	Required	Required
PS	Peptide sequence	PEPTIDESEQ	Required	-	-	-
AC	Uniprot Accession	A9WZ33	-	Required	-	-
CN	Compound name	iRT mixture, substance name	Optional	Optional	Optional	Optional
SP	Species	Escherichia coli K-12	Optional	Optional	Optional	Optional
CV	Compound vendor	in-house or vendor name	Optional	Optional	Required	Optional
CS	Compound specification URI	<a href="http://vendor.web.site/specs/kit.xlsx">http://vendor.web.site/specs/kit.xlsx</a>	Optional	Optional	Optional	Optional
CF	Compound formula	C2H2O	-	-	-	Optional

In addition to specifying the component and its quantity, the injected mass of the main sample SHOULD be specified as *characteristics[mass]*.

An example of SDRF-Proteomics for a sample spiked with a peptide would be:

characteristics[mass]	characteristics[spiked compound]
1 ug	CT=peptide;PS=PEPTIDESEQ;QY=10 fmol

For multiple spiked components, the column *characteristics[spiked compound]* may be repeated.

If the spiked component is another biological sample (e.g. *E. coli* lysate spiked into human sample), then the spiked component MUST be annotated in its own row. Both components of the sample SHOULD have *characteristics[mass]* specified. Inclusion of *characteristics[spiked compound]* is optional in this case; if provided, it SHOULD be the string *spiked* for the spiked sample.

## 8.5. Sample Metadata Guidelines

For detailed guidance on annotating sample metadata, refer to the following conventions documents:

- [Sample Metadata Guidelines](#) - Detailed guidelines for age, sex, disease, organism part, cell type, developmental stage, and other sample characteristics

- [Human Sample Metadata Guidelines](#) - Human-specific metadata including disease staging, treatment history, demographics, and lifestyle factors

# Chapter 9. SDRF-Proteomics: data files metadata

The connection between samples and data files is done using properties annotated with the **comment** prefix. All properties referring to a data file (e.g., MS run file) are annotated with the category [comment](#). This differentiates data file properties from sample properties (characteristics).

## 9.1. CV Term Format for Data File Metadata

For data file metadata (comment columns) that reference ontology terms, use the structured format:  
**NT={term name};AC={accession}**

Examples: [NT=HCD;AC=PRIDE:0000590](#), [NT=Orbitrap;AC=MS:1000484](#)

This format enables automated validation and software extraction from raw files. Sample metadata (characteristics) can use simple term names since they are typically human-annotated.

The following properties MUST be provided for each data file:

Column	Requirement	Description	Ontology
assay name	REQUIRED	Unique identifier for an MS run/data file	Free text
technology type	REQUIRED	Technology used to capture the data	Fixed values
comment[proteomics data acquisition method]	REQUIRED	DDA, DIA, PRM, SRM	<a href="#">PRIDE:0000659</a>
comment[label]	REQUIRED	Label applied to sample (or "label free sample")	<a href="#">PRIDE - Labels</a>
comment[instrument]	REQUIRED	Mass spectrometer model	<a href="#">PSI-MS - Instruments</a>
comment[cleavage agent details]	REQUIRED	Enzyme information (use "not applicable" for top-down/undigested samples)	<a href="#">PSI-MS - Cleavage agents</a>
comment[fraction identifier]	REQUIRED	Fraction number (1 if not fractionated)	Integer
comment[technical replicate]	REQUIRED	Technical replicate number (1 if none)	Integer
comment[data file]	REQUIRED	Name of the raw file	Free text

Example:

source name	assay name	technology type	comment[proteomics data acquisition method]	comment[label]	comment[instrument]	comment[data file]
sample_1	sample1_ru_n1	proteomic profiling by mass spectrometry	data-dependent acquisition	label free sample	Q Exactive HF	sample1.raw

## 9.2. Sample Preparation and Fragmentation

For detailed documentation of sample preparation and MS/MS fragmentation properties, see the [Data File Metadata Guidelines](#):

- **Sample preparation:** depletion, reduction reagent, alkylation reagent
- **Fractionation:** fractionation method (used with fraction identifier)
- **Fragmentation:** collision energy, dissociation method

## 9.3. Data acquisition

Proteomics data acquisition method can happen in multiple ways: Data Dependent Acquisition (DDA), Data Independent Acquisition (DIA), and targeted approaches. The SDRF-Proteomics file format REQUIRES capturing the method used for the data acquisition in the *comment[proteomics data acquisition method]* column. The values MUST be children of the PRIDE ontology term [proteomics data acquisition method \(PRIDE:0000659\)](#). The following values are commonly used:

- [data-dependent acquisition](#)
- [data-independent acquisition](#)
  - [diaPASEF](#)
  - [SWATH MS](#)
- [parallel reaction monitoring](#)
- [selected reaction monitoring](#)

**IMPORTANT**

The *comment[proteomics data acquisition method]* column is REQUIRED for all mass spectrometry-based SDRF files. This field must be explicitly specified and cannot be omitted or assumed.

You can find an example of a DIA experiment in the following link: [DIA example](#)

**TIP**

For DIA experiments, additional properties like MS1 scan range can be captured. See [DIA Scan Window Limits](#) in the Data File Metadata Guidelines.

## 9.4. Data File Metadata Guidelines

For detailed guidance on data file metadata, refer to the conventions document:

- [Data File Metadata Guidelines](#) - Detailed guidelines for labels, instruments, modifications, cleavage agents, mass tolerances, RAW file URLs, and other data file properties

# Chapter 10. Additional SDRF Rules

## 10.1. Row Uniqueness Requirements

SDRF files must satisfy specific uniqueness constraints to ensure data integrity and enable proper indexing by analysis tools.

**Error-level constraint (validation fails):** The combination of `source name + assay name + comment[label]` MUST be unique across all rows in the SDRF file. If two rows have identical values for all three columns, validation will fail with an error. This constraint ensures that each sample-run-label combination can be uniquely identified.

**Warning-level constraint (validation warns):** The combination of `source name + assay name` SHOULD be unique across all rows. Non-unique combinations will generate a warning during validation. This constraint helps identify potential issues where the same sample appears to have multiple entries for the same MS run without distinguishing labels.

**Assay name uniqueness:** Each distinct MS run/data file MUST have exactly one globally unique `assay name`, and no two different data files may share an assay name. To ensure uniqueness, it is RECOMMENDED to incorporate sample-specific information in assay names, such as sample IDs or replicate numbers (e.g., "sample1\_run1", "sample1\_run2", "patient001\_fraction01").

**NOTE**

For multiplexed experiments (e.g., TMT, iTRAQ), multiple SDRF rows will share the same `assay name` because multiple samples are analyzed in a single MS run. In these cases, the `comment[label]` column distinguishes between different samples within the same run, and the combination of `source name + assay name + comment[label]` remains unique.

Example of valid multiplexed experiment:

source name	...	assay name	comment[label]	...	comment[data file]
sample_A	...	TMT_batch1_ru n1	TMT126	...	batch1_run1.ra w
sample_B	...	TMT_batch1_ru n1	TMT127N	...	batch1_run1.ra w
sample_C	...	TMT_batch1_ru n1	TMT127C	...	batch1_run1.ra w
sample_D	...	TMT_batch1_ru n1	TMT128N	...	batch1_run1.ra w

In this example, all four rows share the same `assay name` and `comment[data file]` because they represent different samples multiplexed in a single MS run. The combination of `source name + assay name + comment[label]` is unique for each row.

# Chapter 11. Ontologies and Controlled Vocabularies

SDRF-Proteomics uses ontologies and controlled vocabularies (CVs) to standardize metadata values. The following ontologies are supported:

Category	Ontology/CV	Description	Notes
<b>General Purpose</b>			
General	<a href="#">Experimental Factor Ontology (EFO)</a>	General experimental metadata	
General	<a href="#">PATO</a>	Phenotype and Trait Ontology	
General	<a href="#">NCI Thesaurus (NCIT)</a>	Biomedical terminology	
General	<a href="#">PRIDE Controlled Vocabulary</a>	Proteomics-specific terms	
<b>Organism and Taxonomy</b>			
Taxonomy	<a href="#">NCBI Taxonomy (NCBITaxon)</a>	Organism classification	
<b>Anatomy and Cell Types</b>			
Anatomy	<a href="#">UBERON</a>	Cross-species anatomy ontology	
Cell Type	<a href="#">Cell Ontology (CL)</a>	Cell type classification	
Anatomy	<a href="#">BRENDA Tissue Ontology (BTO)</a>	Tissues and cell lines	
Anatomy	<a href="#">Plant Ontology (PO)</a>	Plant anatomy and development	For plant samples
Anatomy	<a href="#">FlyBase Anatomy (FBbt)</a>	Drosophila anatomy	For Drosophila samples
Anatomy	<a href="#">WormBase Anatomy (WBbt)</a>	C. elegans anatomy	For C. elegans samples
Anatomy	<a href="#">Zebrafish Anatomy (ZFA)</a>	Zebrafish anatomy and development	For zebrafish samples
<b>Disease (see <a href="#">Disease Annotation Guidelines</a>)</b>			
Disease	<a href="#">Mondo Disease Ontology (MONDO)</a>	Unified disease ontology	RECOMMENDED
Disease	<a href="#">Experimental Factor Ontology (EFO)</a>	Disease terms from EFO	
Healthy samples	<a href="#">Phenotype And Trait Ontology (PATO)</a>	Use <a href="#">normal</a> (PATO:0000461) for healthy samples	
<b>Cell Lines</b>			
Cell Lines	<a href="#">Cellosaurus</a>	Cell line knowledge resource	RECOMMENDED

Category	Ontology/CV	Description	Notes
Cell Lines	<a href="#">Cell Line Ontology (CLO)</a>	Cell line ontology	Legacy support only
<b>Mass Spectrometry and Proteomics</b>			
MS/Proteomics	<a href="#">PSI Mass Spectrometry CV (PSI-MS)</a>	Instruments, methods, parameters	
Modifications	<a href="#">Unimod</a>	Protein modifications database	
Modifications	<a href="#">PSI-MOD CV</a>	Protein modifications ontology	
<b>Other</b>			
Chemistry	<a href="#">ChEBI</a>	Chemical Entities of Biological Interest	
Environment	<a href="#">Environment Ontology (ENVO)</a>	Environmental sample classification	For metaproteomics
Ancestry	<a href="#">Human Ancestry Ontology (HANCESTRO)</a>	Human ancestry categories	For human samples

# Chapter 12. Core Templates

## 12.1. What is a Template?

A **template** in SDRF-Proteomics is a predefined set of metadata columns (both required and recommended) that ensures consistent and complete annotation for a specific type of experiment or sample. Templates serve the same purpose as **metadata checklists**, **minimum information standards** (like MIAPE), or **validation schemas** in other data standards—they define what information must be captured to make a dataset FAIR (Findable, Accessible, Interoperable, Reusable).

Each template includes both **sample metadata** (characteristics columns describing the biological sample) and **data file metadata** (comment columns describing the MS data files)—everything needed to create a complete, validated SDRF file.

## 12.2. Template Inheritance and Composition

Templates in SDRF-Proteomics follow a **hierarchical inheritance model**:

- **Core templates** (default, human, vertebrates, invertebrates, plants) define organism-specific requirements
- **Experiment-type templates** (DDA, DIA, single-cell, immunopeptidomics, crosslinking, metaproteomics, affinity-proteomics) define acquisition or methodology-specific requirements
- Templates can be **combined**: a human single-cell proteomics experiment would use both the human template AND the single-cell template
- Templates can be **extended**: add custom columns beyond the template requirements for study-specific metadata

When combining templates, include all required columns from each applicable template. The validator will check compliance with all specified templates.

**TIP**

**For developers and maintainers:** When creating new templates, follow the inheritance principle—child templates must not weaken parent requirements. For example, if the default template has `organism` as REQUIRED, the human template (which inherits from it) must also keep `organism` as REQUIRED. Child templates may add new required columns or promote recommended columns to required, but never downgrade required columns to optional.

## 12.3. YAML Template Definitions

Templates are implemented as YAML files that define validation rules, column requirements, and inheritance relationships. These YAML definitions are used by the [sdrf-pipelines](#) validator to check SDRF files for compliance.

**Basic YAML template structure:**

```
name: default          # Template identifier
```

```

description: Default SDRF template for general proteomics experiments
version: 1.1.0
extends: minimum # Parent template (inheritance)

validators: # File-level validators
  - validator_name: min_columns
    params:
      min_columns: 12

columns: # Column definitions
  - name: characteristics[disease]
    description: Disease state of the sample
    requirement: required # required | recommended | optional
    allow_not_applicable: true
    allow_not_available: true
    validators: # Column-level validators
      - validator_name: ontology
        params:
          ontologies:
            - mondo
            - efo
        error_level: warning

```

### Key properties:

- [name](#): Template identifier used in validation commands
- [extends](#): Parent template to inherit from (e.g., `minimum`, `default`)
- [requirement](#): Column requirement level (`required`, `recommended`, `optional`)
- [allow\\_not\\_applicable](#) / [allow\\_not\\_available](#): Whether special values are permitted
- [validators](#): Validation rules (ontology checks, patterns, value constraints)

Template files are located in the [core-templates](#) and [templates](#) directories.

## 12.4. Choosing a Template

Choose the appropriate core template based on your sample organism:

- [Default template](#): Basic template for any proteomics experiment
- [Human template](#): For human samples with additional clinical metadata (age, sex, ancestry)
- [Vertebrates template](#): For non-human vertebrate species (mouse, rat, zebrafish)
- [Invertebrates template](#): For insects (Drosophila), nematodes (C. elegans), and other invertebrates
- [Plants template](#): For plant species (Arabidopsis, crops)

Then, if applicable, also apply an [experiment-type template](#) for specialized methodologies.

For detailed explanations of each column, see [sample metadata](#) for sample properties and [data file metadata](#) for data file properties.

## 12.5. Default Template

The default template is the most basic template that can be used for any proteomics experiment. Use this when no organism-specific template is applicable.

### Checklist:

Column	Category	Requirement	Example
<b>Sample Metadata</b>			
source name	Sample	Required	sample_1
characteristics[organism]	Sample	Required	homo sapiens
characteristics[organism part]	Sample	Required	liver
characteristics[disease]	Sample	Required	normal
characteristics[biological replicate]	Sample	Required	1
<b>Data File Metadata</b>			
assay name	Data	Required	run_1
technology type	Data	Required	proteomic profiling by mass spectrometry
comment[proteomics data acquisition method]	Data	Required	data-dependent acquisition
comment[label]	Data	Required	label free sample
comment[instrument]	Data	Required	Q Exactive HF
comment[cleavage agent details]	Data	Required	NT=Trypsin;AC=MS:1001 251
comment[fraction identifier]	Data	Required	1
comment[technical replicate]	Data	Required	1
comment[data file]	Data	Required	sample_1.raw

Template file: [sdrf-default.sdrf.tsv](#)

**NOTE**

The `characteristics[cell type]` column is RECOMMENDED across all templates when the cell type is known or can be determined. Use `not available` if the cell type cannot be determined (e.g., whole tissue samples, mixed cell populations). For cell line experiments, use the cell-lines template which provides more specific guidance.

## 12.6. Human Template

The human template extends the default template with clinical and demographic metadata required for human samples.

### Checklist:

Column	Category	Requirement	Example
<b>Sample Metadata</b>			
source name	Sample	Required	patient_001
characteristics[organism]	Sample	Required	homo sapiens
characteristics[organism part]	Sample	Required	liver
characteristics[disease]	Sample	Required	hepatocellular carcinoma
characteristics[cell type]	Sample	Recommended	hepatocyte
characteristics[biological replicate]	Sample	Required	1
characteristics[age]	Sample	Required	45Y
characteristics[sex]	Sample	Required	male
characteristics[ancestry category]	Sample	Recommended	European
characteristics[individual]	Sample	Recommended	P001
<b>Data File Metadata</b>			
assay name	Data	Required	patient_001_run1
technology type	Data	Required	proteomic profiling by mass spectrometry
comment[proteomics data acquisition method]	Data	Required	data-dependent acquisition
comment[label]	Data	Required	label free sample
comment[instrument]	Data	Required	Orbitrap Exploris 480
comment[cleavage agent details]	Data	Required	NT=Trypsin;AC=MS:1001 251
comment[fraction identifier]	Data	Required	1
comment[technical replicate]	Data	Required	1

Column	Category	Requirement	Example
comment[data file]	Data	Required	patient_001.raw

Template file: [sdrf-human.sdrf.tsv](#)

**NOTE** The characteristics[individual] column is optional and is only used when a code for the individual is available; if not, use **not available**. For age encoding format, see [Sample Metadata Guidelines](#).

## 12.7. Vertebrates Template

The vertebrates template is used for non-human vertebrate species such as mouse, rat, zebrafish, and other model organisms.

### Checklist:

Column	Category	Requirement	Example
<b>Sample Metadata</b>			
source name	Sample	Required	mouse_001
characteristics[organism]	Sample	Required	mus musculus
characteristics[organism part]	Sample	Required	brain
characteristics[disease]	Sample	Required	normal
characteristics[cell type]	Sample	Recommended	neuron
characteristics[biological replicate]	Sample	Required	1
characteristics[developmental stage]	Sample	Recommended	adult
<b>Data File Metadata</b>			
assay name	Data	Required	mouse_001_run1
technology type	Data	Required	proteomic profiling by mass spectrometry
comment[proteomics data acquisition method]	Data	Required	data-dependent acquisition
comment[label]	Data	Required	label free sample
comment[instrument]	Data	Required	timsTOF Pro
comment[cleavage agent details]	Data	Required	NT=Trypsin;AC=MS:1001 251

Column	Category	Requirement	Example
comment[fraction identifier]	Data	Required	1
comment[technical replicate]	Data	Required	1
comment[data file]	Data	Required	mouse_001.raw

Template file: [sdrf-vertebrates.sdrf.tsv](#)

## 12.8. Invertebrates Template

The invertebrates template is used for non-vertebrate animal species such as insects (*Drosophila*), nematodes (*C. elegans*), and other invertebrate model organisms.

### Checklist:

Column	Category	Requirement	Example
<b>Sample Metadata</b>			
source name	Sample	Required	fly_sample_1
characteristics[organism]	Sample	Required	drosophila melanogaster
characteristics[organism part]	Sample	Required	head
characteristics[disease]	Sample	Required	normal
characteristics[cell type]	Sample	Recommended	neuron
characteristics[biological replicate]	Sample	Required	1
<b>Data File Metadata</b>			
assay name	Data	Required	fly_sample_1_run1
technology type	Data	Required	proteomic profiling by mass spectrometry
comment[proteomics data acquisition method]	Data	Required	data-dependent acquisition
comment[label]	Data	Required	label free sample
comment[instrument]	Data	Required	Q Exactive Plus
comment[cleavage agent details]	Data	Required	NT=Trypsin;AC=MS:1001 251

Column	Category	Requirement	Example
comment[fraction identifier]	Data	Required	1
comment[technical replicate]	Data	Required	1
comment[data file]	Data	Required	fly_sample_1.raw

Template file: [sdrf-invertebrates.sdrf.tsv](#)

**NOTE**

For Drosophila samples, use [FBbt](#) (FlyBase anatomy ontology) for organism part. For *C. elegans*, use [WBbt](#) (WormBase anatomy ontology).

## 12.9. Plants Template

The plants template is used for plant species including model organisms like *Arabidopsis thaliana* and crop species.

### Checklist:

Column	Category	Requirement	Example
<b>Sample Metadata</b>			
source name	Sample	Required	arabidopsis_col0_1
characteristics[organism]	Sample	Required	arabidopsis thaliana
characteristics[organism part]	Sample	Required	leaf
characteristics[disease]	Sample	Required	normal
characteristics[cell type]	Sample	Recommended	guard cell
characteristics[biological replicate]	Sample	Required	1
<b>Data File Metadata</b>			
assay name	Data	Required	arabidopsis_col0_run1
technology type	Data	Required	proteomic profiling by mass spectrometry
comment[proteomics data acquisition method]	Data	Required	data-dependent acquisition
comment[label]	Data	Required	label free sample
comment[instrument]	Data	Required	Orbitrap Fusion Lumos

Column	Category	Requirement	Example
comment[cleavage agent details]	Data	Required	NT=Trypsin;AC=MS:1001 251
comment[fraction identifier]	Data	Required	1
comment[technical replicate]	Data	Required	1
comment[data file]	Data	Required	arabidopsis_col0.raw

Template file: [sdrf-plants.sdrf.tsv](#)

**NOTE** For plant samples, use [PO](#) (Plant Ontology) for organism part and cell type annotations.

## 12.10. Experiment-Type Templates

In addition to core templates, SDRF-Proteomics provides specialized templates. These templates extend the core templates with methodology-specific columns.

- **DDA Acquisition:** Data-dependent acquisition experiments. Includes columns for dissociation method, collision energy, fractionation method, modification parameters, and mass tolerances.
- **DIA Acquisition:** Data-independent acquisition experiments. Includes columns for scan window limits, isolation window width, DIA method, and spectral library information.
- **Cell Lines:** Experiments using cell line samples. Includes Cellosaurus integration for cell line identification.
- **Single-Cell Proteomics:** Single-cell proteomics experiments. Includes columns for cell isolation method, carrier proteome, and single-cell identifiers.
- **Immunopeptidomics:** MHC peptide immunopeptidomics. Includes columns for MHC class, HLA typing, and enrichment methods.
- **Crosslinking MS:** Cross-linking mass spectrometry experiments. Includes columns for crosslinking reagents and enrichment methods.
- **Metaproteomics:** Environmental and microbiome proteomics. Includes columns for environmental sample type and geographic location.
- **Affinity Proteomics:** Affinity-based proteomics (Olink, SomaScan). Includes columns specific to these platforms.

## 12.11. Column Cardinality

Some columns can appear multiple times for the same sample. The cardinality rules are:

- **Single (1):** Column appears exactly once per sample (e.g., biological replicate)
- **Multiple (\*):** Column can appear multiple times (e.g., organism part can specify both "heart" and

"heart left ventricle")

Example of multiple organism part columns:

source name	...	characteristics[organism part]	characteristics[organism part]	...
sample-1	...	heart	heart left ventricle	...

The template files can be downloaded from the [core-templates](#) folder.

## 12.12. Extending SDRF with Custom Columns

Templates define the minimum required and recommended columns for a given experiment type. However, SDRF-Proteomics is designed to be **extensible** - you can add any additional columns to capture study-specific metadata beyond what templates define. This section provides guidance on adding custom columns.

### 12.12.1. When to Add Custom Columns

Add custom columns when your experiment requires metadata that is:

- Not covered by existing templates but important for data interpretation
- Specific to your experimental design or technology
- Required for cross-study integration or data reuse
- Needed to comply with domain-specific standards (e.g., clinical trials, environmental studies)

### 12.12.2. Rules for Custom Columns

#### 1. Use appropriate column prefixes:

- `characteristics[...]` - For sample-related metadata (properties of the biological material)
- `comment[...]` - For technical or protocol-related metadata (data acquisition, processing)
- `factor value[...]` - For experimental variables (see [Chapter 13](#))

#### 2. Follow naming conventions:

- Use lowercase for column names inside brackets
- Use descriptive, specific names (e.g., `characteristics[tumor grade]` not `characteristics[grade]`)
- Use ontology terms when they exist (e.g., `characteristics[patient bmi]` maps to EFO term)

#### 3. Use controlled vocabularies when available:

- Reference existing ontologies (EFO, MONDO, UBERON, etc.) for values
- Prefer standardized terms over free text when possible
- For new terms, follow the key=value format if ontology reference is needed

### 12.12.3. Common Additional Columns

The following columns are not required by templates but are commonly used and recommended for specific study types:

Column	Category	Description	Example Values	When to Use
characteristics[material type]	Sample	Type of biological material (derived from <a href="#">MAGE-TAB</a> )	tissue, cell, cell line, organism part, whole organism, synthetic	Cross-omics integration, clarifying sample origin
characteristics[treatment]	Sample	Treatment applied to the sample	dexamethasone, vehicle control, untreated	Drug treatment studies
characteristics[time point]	Sample	Time of sample collection	0h, 24h, 7d, baseline	Time-course experiments
characteristics[dose ]	Sample	Dose of treatment if applicable	10 mg/kg, 100 nM, high dose	Dose-response studies
characteristics[patient bmi]	Sample	Body mass index (for human studies)	25.3 kg/m <sup>2</sup> , 30.1 kg/m <sup>2</sup>	Metabolic or obesity-related studies
characteristics[smoking status]	Sample	Smoking history of the patient	never smoked, current smoker, former smoker	Lung or cardiovascular studies
comment[sample preparation]	Data	Sample preparation method details	in-solution digestion, FASP, SP3	Detailed protocol documentation
comment[enrichment method]	Data	Enrichment or depletion strategy	phosphopeptide enrichment, glycopeptide enrichment	PTM-focused studies

See [Material Type Guidelines](#) for detailed guidance on the `characteristics[material type]` column.

### 12.12.4. Example: Adding Study-Specific Columns

For a drug treatment time-course study, you might add columns beyond the template requirements:

source name	characteristics[organism]	characteristics[disease]	characteristics[treatment]	characteristics[time point]	characteristics[dose]	...	factor value[treatment]	factor value[time point]
sample_ctrl_0h	homo sapiens	normal	vehicle control	0h	not applicable	...	vehicle control	0h

source name	characteristics[organism]	characteristics[disease]	characteristics[treatment]	characteristics[time point]	characteristics[dosage]	...	factor value[treatment]	factor value[time point]
sample_d rug_24h	homo sapiens	normal	dexamethasone	24h	100 nM	...	dexamethasone	24h

**TIP**

When adding custom columns, check the [SDRF Terms Reference](#) and [EBI Ontology Lookup Service \(OLS\)](#) to find existing terms that match your metadata needs. Using standardized terms improves data interoperability.

# Chapter 13. Factor Values (Study Variables)

Factor values identify the experimental variables being studied - the conditions you want to compare in your analysis. They highlight which sample characteristics are the focus of your experiment.

## 13.1. Column Format

```
factor value[{variable name}]
```

## 13.2. When to Use Factor Values

Use factor values to indicate:

- The primary variable(s) under investigation
- Conditions being compared (e.g., disease vs. normal, treated vs. untreated)
- Variables that define experimental groups

**NOTE**

Use `normal` (not "control") in the disease field for healthy samples. "Control" is an experimental design concept, not a disease state. See [Disease Annotation Guidelines](#) for details.

## 13.3. Rules

- Factor value columns SHOULD appear after all characteristics and comment columns
- Multiple factor values can be used when studying multiple variables
- The value in a factor value column typically mirrors a characteristics column value

## 13.4. Example

In an experiment comparing tumor vs. normal tissue across different cancer stages:

source name	...	characteristi cs[disease]	characteristi cs[disease stage]	...	factor value[diseas e]	factor value[diseas e stage]
tumor_sample_1	...	breast carcinoma	stage II	...	breast carcinoma	stage II
normal_sample_1	...	normal	not applicable	...	normal	not applicable
tumor_sample_2	...	breast carcinoma	stage III	...	breast carcinoma	stage III

In this example, both `disease` and `disease stage` are factor values because the experiment aims to compare expression differences between disease states and across cancer stages.

## Chapter 14. Examples of Annotated Datasets

The following table provides links to example SDRF files for different experiment types. These can serve as references when creating your own SDRF files.

Experiment Type	Dataset	Description	SDRF URL
Label-free	PXD008934	Human proteome label-free quantification	<a href="#">View SDRF</a>
TMT	PXD017710	TMT-labeled quantitative proteomics	<a href="#">View SDRF</a>
SILAC	PXD000612	SILAC-based quantification	<a href="#">View SDRF</a>
DIA	PXD018830	Data-independent acquisition	<a href="#">View SDRF</a>
Phosphoproteomics	PXD000759	PTM enrichment study	<a href="#">View SDRF</a>
Cell lines	PXD001819	Cell line proteomics	<a href="#">View SDRF</a>

A comprehensive collection of annotated projects is available at: [Annotated Projects Repository](#)

## Chapter 15. Ongoing template discussions

We have created a file in GitHub [Ongoing template discussions](#) where we aggregate all the ongoing discussions about the format and new templates.

## Chapter 16. Intellectual Property Statement

The PSI takes no position regarding the validity or scope of any intellectual property or other rights that might be claimed to pertain to the implementation or use of the technology described in this document or the extent to which any license under such rights might or might not be available; neither does it represent that it has made any effort to identify any such rights. Copies of claims of rights made available for publication and any assurances of licenses to be made available or the result of an attempt made to obtain a general license or permission for the use of such proprietary rights by implementers or users of this specification can be obtained from the PSI Chair.

The PSI invites any interested party to bring to its attention any copyrights, patents or patent applications, or other proprietary rights which may cover technology that may be required to practice this recommendation. Please address the information to the PSI Chair (see contacts information at PSI website).

## Chapter 17. Copyright Notice

Copyright © Proteomics Standards Initiative (2020). All Rights Reserved.

This document and translations of it may be copied and furnished to others, and derivative works that comment on or otherwise explain it or assist in its implementation may be prepared, copied, published, and distributed, in whole or in part, without the restriction of any kind, provided that the above copyright notice and this paragraph are included on all such copies and derivative works. However, this document itself may not be modified in any way, such as by removing the copyright notice or references to the PSI or other organizations, except as needed for the purpose of developing Proteomics Recommendations in which case the procedures for copyrights defined in the PSI Document process must be followed, or as required to translate it into languages other than English.

The limited permissions granted above are perpetual and will not be revoked by the PSI or its successors or assigns.

This document and the information contained herein is provided on an "AS IS" basis and THE PROTEOMICS STANDARDS INITIATIVE DISCLAIMS ALL WARRANTIES, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO ANY WARRANTY THAT THE USE OF THE INFORMATION HEREIN WILL NOT INFRINGE ANY RIGHTS OR ANY IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE."

## Chapter 18. How to cite

Please cite this document as:

Dai C, Füllgrabe A, Pfeuffer J, Solovyeva EM, Deng J, Moreno P, Kamatchinathan S, Kundu DJ, George N, Fexova S, Grüning B, Föll MC, Griss J, Vaudel M, Audain E, Locard-Paulet M, Turewicz M, Eisenacher M, Uszkoreit J, Van Den Bossche T, Schwämmle V, Webel H, Schulze S, Bouyssie D, Jayaram S, Duggineni VK, Samaras P, Wilhelm M, Choi M, Wang M, Kohlbacher O, Brazma A, Papatheodorou I, Bandeira N, Deutsch EW, Vizcaíno JA, Bai M, Sachsenberg T, Levitsky LI, Perez-Riverol Y. A proteomics sample metadata representation for multiomics integration and big data analysis. Nat Commun. 2021 Oct 6;12(1):5854. doi: 10.1038/s41467-021-26111-3. PMID: 34615866; PMCID: PMC8494749. [Manuscript - <https://www.nature.com/articles/s41467-021-26111-3>]

## References

- [1] Y. Perez-Riverol, S. European Bioinformatics Community for Mass, Toward a Sample Metadata Standard in Public Proteomics Repositories, *J Proteome Res* 19(10) (2020) 3906-3909.  
[doi:10.1021/acs.jproteome.0c00376](https://doi.org/10.1021/acs.jproteome.0c00376)
- [2] A. Gonzalez-Beltran, E. Maguire, S.A. Sansone, P. Rocca-Serra, linkedISA: semantic representation of ISA-Tab experimental metadata, *BMC Bioinformatics* 15 Suppl 14 (2014) S4.  
[doi:10.1186/1471-2105-15-S14-S4](https://doi.org/10.1186/1471-2105-15-S14-S4)
- [3] T.F. Rayner, P. Rocca-Serra, P.T. Spellman, et al., A simple spreadsheet-based, MIAME-supportive format for microarray data: MAGE-TAB, *BMC Bioinformatics* 7 (2006) 489.  
[doi:10.1186/1471-2105-7-489](https://doi.org/10.1186/1471-2105-7-489)
- [4] P. Blainey, M. Krzywinski, N. Altman, Points of significance: replication, *Nat Methods* 11(9) (2014) 879-80. [doi:10.1038/nmeth.3091](https://doi.org/10.1038/nmeth.3091)
- [5] D. Gupta, I. Liyanage, Y. Perez-Riverol, et al., BioSamples database: the global hub for sample metadata and multi-omics integration, *Nucleic Acids Res* (2025). [doi:10.1093/nar/gkaf1133](https://doi.org/10.1093/nar/gkaf1133)