

# Package: phiperio (via r-universe)

June 4, 2026

**Title** PhIP-Seq Data Import and Validation Tools

**Version** 0.5.0

**Description** Provides utilities to import, validate, and manage PhIP-Seq datasets, including standardized conversion pipelines, data checks, and access to cached peptide metadata.

**License** GPL (>= 3)

**Encoding** UTF-8

**Depends** R (>= 4.1.0)

**Imports** chk, cli, DBI, dbplyr, dplyr, duckdb, rlang, tibble, tools, withr, yaml

**Suggests** testthat, mockery, knitr, rmarkdown

**Roxygen** list(markdown = TRUE)

**RoxygenNote** 7.3.3

**URL** <https://github.com/Polymerase3/phiperio>,  
<https://polymerase3.github.io/phiperio/>

**BugReports** <https://github.com/Polymerase3/phiperio/issues>

**VignetteBuilder** knitr

**Config/pak/sysreqs** libicu-dev xz-utils

**Repository** <https://polymerase3.r-universe.dev>

**Date/Publication** 2026-02-04 13:42:22 UTC

**RemoteUrl** <https://github.com/Polymerase3/phiperio>

**RemoteRef** main

**RemoteSha** cd3f886b5dfed8be36e32ed628cb48c40a7993d0

## Contents

add_exist	2
close.phip_data	3
convert_legacy	3

convert_standard . . . . .	5
create_data . . . . .	7
expand_data . . . . .	9
export_parquet . . . . .	10
get_counts . . . . .	11
get_example_path . . . . .	11
get_meta . . . . .	12
get_peptide_library . . . . .	12
load_example_data . . . . .	13
merge.phip_data . . . . .	14
phip_data_join . . . . .	15

<b>Index</b>	<b>16</b>
--------------	-----------

---

add_exist	<i>Ensure an existence flag (all ones) on data_long</i>
-----------	---

---

## Description

Appends/overwrites a column (default: "exist") filled with 1L on the lazy data\_long table. Preserves laziness; no collection is forced.

## Usage

```
add_exist(phip_data, exist_col = "exist", overwrite = FALSE)
```

## Arguments

phip_data	A <phip_data> object.
exist_col	Name of the existence column to append/overwrite.
overwrite	If FALSE and the column exists, abort with a phiperio-style error.

## Value

Modified <phip\_data> with updated data\_long.

## Examples

```
pd <- load_example_data()
pd <- add_exist(pd, overwrite = TRUE) # overwrites if present
```

---

close.phip_data	<i>Close phip_data connections</i>
-----------------	------------------------------------

---

### Description

Closes any open database connections held by a `phip_data` object. This includes the main `data_long` backend connection and any peptide-library connection stored in attributes or metadata. The method is idempotent and safe to call multiple times.

### Usage

```
## S3 method for class 'phip_data'  
close(con, ...)
```

### Arguments

<code>con</code>	A valid <code>phip_data</code> object.
<code>...</code>	Unused (for S3 generic compatibility).

### Value

The input `phip_data` object, invisibly.

### Examples

```
pd <- load_example_data()  
close(pd)
```

---

convert_legacy	<i>Convert legacy Carlos-style input to a modern <b>phip_data</b> object</i>
----------------	--

---

### Description

`convert_legacy()` ingests the original three-file PhIP-Seq input (binary *exist* matrix, *samples* metadata, optional *timepoints* map). Paths can be supplied directly or via a single YAML config; explicit arguments always override the YAML. The function normalises the chosen DuckDB storage, validates every file, and returns a ready-to-use `phip_data` object.

**Usage**

```

convert_legacy(
  exist_file = NULL,
  fold_change_file = NULL,
  samples_file = NULL,
  input_file = NULL,
  hit_file = NULL,
  timepoints_file = NULL,
  extra_cols = NULL,
  output_dir = NULL,
  peptide_library = TRUE,
  n_cores = 8,
  materialise_table = TRUE,
  config_yaml = NULL
)

```

**Arguments**

exist_file	Path to the <b>exist</b> CSV (peptide x sample binary matrix). <i>Required unless given in config_yaml.</i>
fold_change_file	Path to the <b>fold_change</b> CSV (peptide x sample numeric matrix). <i>Required unless given in config_yaml.</i>
samples_file	Path to the <b>samples</b> CSV (sample metadata). <i>Required unless given in config_yaml.</i>
input_file, hit_file	Paths to the <b>raw_counts</b> CSV (peptide x sample integer matrix). <i>Required unless given in config_yaml.</i>
timepoints_file	Path to the <b>timepoints</b> CSV (subject <-> sample mapping). Optional for cross-sectional data.
extra_cols	Character vector of extra metadata columns to retain.
output_dir	<i>Deprecated.</i> Ignored with a warning.
peptide_library	logical, defining if the peptide_library is to be downloaded from the official phiperio GitHub
n_cores	Integer >= 1. Number of CPU threads DuckDB may use while reading and writing files.
materialise_table	Logical. If FALSE the result is registered as a <b>view</b> ; if TRUE the table is fully <b>materialised</b> and stored on disk, trading higher load time and storage for faster repeated queries.
config_yaml	Optional YAML file containing any of the above parameters (see example).

**Details**

Input files are validated in two stages:

- **Fast-fail** checks (paths, extensions, and required arguments) run during path resolution.
- **Data validation** (required columns, uniqueness, value ranges, etc.) is centralized in `validate_phip_data()`.

## Value

A validated `hip_data` object whose `data_long` slot is backed by a DuckDB connection.

## Examples

```
## 1. Direct-path usage (package example files)
ext <- system.file("extdata", package = "phiperio")
pd <- convert_legacy(
  exist_file = file.path(ext, "exist.csv"),
  samples_file = file.path(ext, "samples_meta.csv"),
  timepoints_file = file.path(ext, "samples2ind_timepoints.csv"),
  peptide_library = FALSE
)

## 2. YAML-driven usage (explicit args override YAML)
pd <- convert_legacy(
  config_yaml = file.path(ext, "config.yaml"),
  peptide_library = FALSE
)
```

---

<code>convert_standard</code>	<i>Convert raw PhIP-Seq output into a <code>hip_data</code> object</i>
-------------------------------	--

---

## Description

`convert_standard()` ingests a "long" table of PhIPsSeq read counts / enrichment statistics, optionally expands it to the full `sample_id` x `peptide_id` grid, and registers the result in DuckDB. The function returns a fully initialised `hip_data` object that can be queried with the tidy API used throughout the package.

## Usage

```
convert_standard(
  data_long_path,
  sample_id = NULL,
  peptide_id = NULL,
  subject_id = NULL,
  timepoint = NULL,
  exist = NULL,
  fold_change = NULL,
  counts_input = NULL,
  counts_hit = NULL,
```

```

sample_id_from_filenames = FALSE,
n_cores = 8,
materialise_table = TRUE,
auto_expand = FALSE,
peptide_library = TRUE
)

```

## Arguments

- data\_long\_path** Character scalar. File or directory containing the *long-format* PhIP-Seq data. Allowed extensions are .csv and .parquet. Directories are treated as partitions of a parquet set.
- sample\_id, peptide\_id, subject\_id, timepoint, exist, fold\_change, counts\_input, counts\_hit** Optional character strings. Supply these only if your column names differ from the defaults ("sample\_id", "peptide\_id", "subject\_id", "timepoint", "exist", "fold\_change", "counts\_input", "counts\_hit"). Each argument should contain the *name* of the column in the incoming data; NULL lets the default stand.
- sample\_id\_from\_filenames** Logical. If TRUE and data\_long\_path is a **directory of files** (CSV or Parquet), automatically derive sample\_id from each filename (basename without extension). Requires that no sample\_id mapping is provided and that the input files do not already contain a sample\_id column. Default: FALSE.
- n\_cores** Integer >= 1. Number of CPU threads DuckDB may use while reading and writing files.
- materialise\_table** Logical. If FALSE the result is registered as a **view**; if TRUE the table is fully **materialised** and stored on disk, trading higher load time and storage for faster repeated queries.
- auto\_expand** Logical. If TRUE and the incoming data are **not** a complete Cartesian product of sample\_id x peptide\_id, missing combinations are generated:
- Columns that are constant within each sample\_id (metadata) are copied to the new rows.
  - Non-recyclable measurement columns (fold\_change, exist, counts\_input, counts\_hit, etc.) are initialised to 0. The expanded table replaces the original *in place*.
- peptide\_library** Logical. If TRUE (default) convert\_standard() will attempt to locate and attach the matching peptide-library metadata for downstream annotation. Set to FALSE to skip this step.

## Details

*Paths are resolved to absolute form* before any work begins, and explicit checks confirm existence as well as extension validity.

**Value**

An S3 object of class `phip_data` containing:

`data_long` The (possibly expanded) long-format table.

`peptide_library` Loaded peptide-library metadata (if `peptide_library = TRUE`).

`meta` List with DuckDB connection handles.

**See Also**

- `create_data()` for the object constructor.
- `dplyr::tbl()` to query DuckDB tables lazily.

**Examples**

```
# Basic import, auto-detecting default column names
phip_obj <- convert_standard(
  data_long_path = get_example_path("phip_mixture"),
  n_cores = 4,
  materialise_table = TRUE
)

# Import a CSV and rename columns
tmp_csv <- tempfile(fileext = ".csv")
utils::write.csv(
  data.frame(
    sample = c("s1", "s1"),
    pep = c("p1", "p2"),
    exist = c(1, 0),
    stringsAsFactors = FALSE
  ),
  tmp_csv,
  row.names = FALSE
)
phip_mem <- convert_standard(
  data_long_path = tmp_csv,
  sample_id      = "sample",
  peptide_id     = "pep",
  peptide_library = FALSE,
  materialise_table = FALSE
)
```

---

`create_data`*Construct a **phip\_data** object*

---

**Description**

Creates a fully-validated S3 object that bundles the tidy PhIP-Seq counts (`data_long`), a peptide-library annotation table, and other metadata. The data itself is validated via `validate_phip_data()`.

## Usage

```
create_data(  
  data_long,  
  peptide_library = TRUE,  
  auto_expand = TRUE,  
  materialise_table = TRUE,  
  meta = list()  
)
```

## Arguments

data_long	A tidy data frame (or <code>tbl_lazy</code> ) with one row per <code>peptide_id</code> x <code>sample_id</code> combination. <b>Required.</b>
peptide_library	A data frame with one row per <code>peptide_id</code> and its annotations. If <code>NULL</code> , the package's current default library is used.
auto_expand	Logical. If <code>TRUE</code> and the input is <b>not</b> already the full Cartesian product of <code>sample_id</code> x <code>peptide_id</code> , the function fills in the missing combinations. <ul style="list-style-type: none"><li>• Columns that are constant within a <code>sample_id</code> (metadata) are duplicated to the newly created rows.</li><li>• Measurement columns such as <code>fold_change</code>, <code>exist</code>, raw counts, or any other non-recyclable fields are initialised to 0. The expanded table replaces <code>data_long</code> in place.</li></ul>
materialise_table	Logical. If <code>FALSE</code> (default) the result is registered as a <b>view</b> . If <code>TRUE</code> the result is fully <b>materialised</b> and stored as a physical table, which speeds up repeated queries at the cost of extra memory/disk.
meta	Optional named list of metadata flags to pre-populate the <code>meta</code> slot (rarely needed by users).

## Value

An object of class "pchip\_data".

## Examples

```
## minimal constructor call  
tidy_counts <- data.frame(  
  sample_id = c("s1", "s1"),  
  peptide_id = c("p1", "p2"),  
  exist = c(1, 0),  
  stringsAsFactors = FALSE  
)  
pd <- create_data(  
  data_long = tidy_counts,  
  peptide_library = FALSE,  
  materialise_table = FALSE  
)
```

---

expand_data	<i>Expand to a full sample_id * peptide_id grid</i>
-------------	---

---

### Description

Create the full Cartesian product of samples and peptides and join back per-sample metadata. For rows introduced by the expansion, numeric/integer columns are filled with 0 and logical columns with FALSE, unless overridden via `fill_override`.

### Usage

```
expand_data(
  x,
  key_col = "sample_id",
  id_col = "peptide_id",
  fill_override = NULL,
  add_exist = FALSE,
  exist_col = "exist",
  validate = TRUE,
  ...
)
```

### Arguments

<code>x</code>	A <phip_data> object.
<code>key_col</code>	Name(s) of the sample identifier column(s). Character scalar or vector, e.g. "sample_id" or <code>c("subject_id", "timepoint_factor")</code> .
<code>id_col</code>	Name of the peptide identifier column. Default "peptide_id".
<code>fill_override</code>	Optional named list of fill values for <b>introduced</b> rows, e.g. <code>list(present = 0L, fold_change = NA_real_)</code> . User-provided entries take precedence over the defaults.
<code>add_exist</code>	If TRUE, add an integer existence flag (0/1) marking whether a row was present before the expansion.
<code>exist_col</code>	Name for the existence flag. If this column already exists, it will be <b>overwritten</b> .
<code>validate</code>	Logical; if TRUE, perform input checks for required columns and uniqueness. Set to FALSE when these checks were already performed upstream (e.g., inside <code>validate_phip_data()</code> ).
<code>...</code>	Reserved for future extensions; currently unused.

### Details

Updates `x$data_long` in place (preserving laziness unless you later `compute()` / `collect()`).

### Value

The updated <phip\_data> object.

### Examples

```
pd <- load_example_data()
pd <- expand_data(pd, fill_override = list(fold_change = NA_real_))
```

---

export_parquet	<i>Export a phip_data Table to Parquet</i>
----------------	--

---

### Description

Exports the data\_long table from a **phip\_data** object to disk in Apache Parquet format.

### Usage

```
export_parquet(x, path)
```

### Arguments

x	A <phip_data> object or a data frame.
path	File path (character) to save the output .parquet file.

### Value

NULL (invisibly).

### Note

The export is performed directly and efficiently from the database/lazy table without reading all data into memory.

### Examples

```
pd <- load_example_data()
out_path <- tempfile(fileext = ".parquet")
export_parquet(pd, out_path)
unlink(out_path)
```

---

get_counts	<i>Retrieve the main PhIP-Seq counts table</i>
------------	--

---

**Description**

Quick accessor for the data\_long slot of a **phip\_data** object.

**Usage**

```
get_counts(x)
```

**Arguments**

x                    A valid phip\_data object.

**Value**

A tibble or lazy table with one row per peptide \* sample pair.

**Examples**

```
pd <- load_example_data()
tbl <- get_counts(pd)
```

---

get_example_path	<i>Path to example PhIP-Seq datasets shipped with phiperio</i>
------------------	--

---

**Description**

Path to example PhIP-Seq datasets shipped with phiperio

**Usage**

```
get_example_path(name = c("phip_mixture"))
```

**Arguments**

name                    Character scalar. Name of the example dataset. Currently supported: "phip\_mixture".

**Value**

A character scalar with an absolute path to the file.

**Examples**

```
sim_path <- get_example_path("phip_mixture")
# phip_obj <- convert_standard(sim_path)
```

---

get_meta	<i>Retrieve the metadata list</i>
----------	-----------------------------------

---

**Description**

Accesses the meta slot, which holds flags such as whether the table is a full peptide \* sample grid, the available outcome columns, etc.

**Usage**

```
get_meta(x)
```

**Arguments**

x                    A valid phip\_data object.

**Value**

A named list.

**Examples**

```
pd <- load_example_data()
meta <- get_meta(pd)
```

---

get_peptide_library	<i>Retrieve the peptide metadata table into DuckDB, forcing atomic types</i>
---------------------	--

---

**Description**

This function uses the phiperio logging utilities for consistent, ASCII-only progress messages and timing. Long-running steps are bracketed with `.ph_with_timing()`, and informational/warning/error messages are emitted via `.ph_log_info()`, `.ph_log_ok()`, `.ph_warn()`, and `.ph_abort()`.

- Downloads the RDS once, sanitizes types (logical, character, numeric), and writes into a DuckDB cache on disk.
- Subsequent calls return a lazy `tbl_dbi` without loading into R memory.

**Usage**

```
get_peptide_library(force_refresh = FALSE)
```

**Arguments**

force\_refresh    Logical. If TRUE, re-downloads and rebuilds the cache.

## Details

**Caching:** A persistent DuckDB database is created under the user cache directory (via `tools::R_user_dir("phiperio", "cache")`). You can override this location with `options(phiperio.cache_dir = "...")`. The `force_refresh` argument bypasses the fast path and rebuilds the cache.

**Sanitization:** Columns are stripped of attributes, list-columns are flattened, textual "NaN" and numeric NaN are coerced to NA. Binary 0/1 fields are converted to `logical`, "TRUE"/"FALSE" (case-insensitive) are converted to `logical`, and numeric-looking character columns (beyond trivial 0/1) are converted to `numeric`. All other atomic types are preserved.

**Integrity check:** If a SHA-256 checksum is provided, a warning is logged when the downloaded file's checksum does not match the expected value.

## Value

A `dplyr::tbl_dbi` pointing to the `peptide_meta` table. The returned object carries an attribute `"duckdb_con"` with the open DBI connection.

## See Also

`dplyr::tbl()`, `DBI::dbConnect()`, `duckdb::duckdb()`

## Examples

```
lib <- get_peptide_library()
```

---

load_example_data	<i>Load Example PhIP-Seq Dataset as &lt;phip_data&gt;</i>
-------------------	---

---

## Description

Convenience helper to quickly load a shipped example dataset ("phip\_mixture") into a `<phip_data>` object, suitable for downstream analysis and visualization. This function wraps `convert_standard`, automatically supplying the correct parameters for the included example data.

## Usage

```
load_example_data(name = c("phip_mixture", "small_mixture"))
```

## Arguments

`name` Character scalar. Name of the shipped example dataset. Currently supported: "phip\_mixture", "small\_mixture".

## Value

A `<phip_data>` object created from the chosen example dataset.

## Examples

```
# Load the example data shipped with the package:
ex <- load_example_data()
# ex is now a <phip_data> object ready for analysis

# Specify the dataset name explicitly
ex2 <- load_example_data("small_mixture")

# Use with downstream analysis/plotting functions as needed
```

---

merge.phip_data	<i>Merge or join a phip_data object</i>
-----------------	---

---

## Description

Merge or join a phip\_data object

## Usage

```
## S3 method for class 'phip_data'
merge(x, y, ...)
```

## Arguments

x	A phip_data object.
y	A data-frame-like object <i>or</i> another phip_data.
...	Arguments forwarded to either <code>base::merge()</code> or the chosen <b>dplyr</b> join (e.g. <code>by =</code> , <code>suffix =</code> , etc.).

## Value

A new phip\_data whose data\_long contains the merged / joined tibble.

## Examples

```
pd <- load_example_data()
merged <- merge(pd, pd, by = c("sample_id", "peptide_id"))
```

---

phip_data_join	<i>dplyr joins for phip_data</i>
----------------	----------------------------------

---

### Description

dplyr joins for phip\_data

### Usage

```
## S3 method for class 'phip_data'
left_join(x, y, ...)

## S3 method for class 'phip_data'
right_join(x, y, ...)

## S3 method for class 'phip_data'
inner_join(x, y, ...)

## S3 method for class 'phip_data'
full_join(x, y, ...)

## S3 method for class 'phip_data'
semi_join(x, y, ...)

## S3 method for class 'phip_data'
anti_join(x, y, ...)
```

### Arguments

x	A phip_data object.
y	A phip_data or a data frame / tbl.
...	Passed to the corresponding dplyr::<join> function.

### Value

A phip\_data object with updated data\_long.

### Examples

```
pd <- load_example_data()
joined <- dplyr::left_join(pd, pd, by = c("sample_id", "peptide_id"))
```

# Index

`add_exist`, 2  
`anti_join.phip_data (phip_data_join)`, 15  
  
`base::merge()`, 14  
  
`close.phip_data`, 3  
`convert_legacy`, 3  
`convert_standard`, 5, 13  
`create_data`, 7  
  
`DBI::dbConnect()`, 13  
`dplyr::tbl()`, 13  
`duckdb::duckdb()`, 13  
  
`expand_data`, 9  
`export_parquet`, 10  
  
`full_join.phip_data (phip_data_join)`, 15  
  
`get_counts`, 11  
`get_example_path`, 11  
`get_meta`, 12  
`get_peptide_library`, 12  
  
`inner_join.phip_data (phip_data_join)`,  
15  
  
`left_join.phip_data (phip_data_join)`, 15  
`load_example_data`, 13  
  
`merge.phip_data`, 14  
  
`phip_data_join`, 15  
  
`right_join.phip_data (phip_data_join)`,  
15  
  
`semi_join.phip_data (phip_data_join)`, 15