

Package ‘artoo’

June 24, 2026

Title Lossless CDISC-Native Input and Output for Clinical Datasets

Version 0.1.1

Description Reads and writes clinical-trial datasets losslessly across 'SAS' XPORT (XPT), Clinical Data Interchange Standards Consortium (CDISC) Dataset-JSON, and 'Apache Parquet', applying a specification to produce submission-ready Study Data Tabulation Model (SDTM) and Analysis Data Model (ADaM) datasets. A single canonical metadata model carries labels, CDISC data types, lengths, 'SAS' display formats, controlled-terminology references, and sort keys identically across every format, so conversion between any two formats is lossless by construction. Pure 'R' and lightweight, with no external 'SAS' or 'Java' runtime. Implements the published format specifications for CDISC Dataset-JSON

(<<https://cdisc-org.github.io/DataExchange-DatasetJson/doc/dataset-json1-1.html>>) and 'SAS' XPORT (<<https://www.loc.gov/preservation/digital/formats/fdd/fdd000466.shtml>>).

License MIT + file LICENSE

URL <https://vthanik.github.io/artoo/>, <https://github.com/vthanik/artoo>

BugReports <https://github.com/vthanik/artoo/issues>

Encoding UTF-8

Language en-US

Depends R (>= 4.2)

Imports cli, hms, jsonlite, nanoparquet (>= 0.3.0), rlang (>= 1.1.0), S7, utf8

Suggests quarto, readxl, testthat (>= 3.0.0), tibble, withr, writexl, xml2

VignetteBuilder quarto

Config/testthat/edition 3

LazyData true

Config/roxygen2/version 8.0.0

NeedsCompilation no

Author Vignesh Thanikachalam [aut, cre, cph]

Maintainer Vignesh Thanikachalam <about.vignesh@gmail.com>

Repository CRAN

Date/Publication 2026-06-24 09:50:18 UTC

Contents

apply_spec	3
artoo_checks	6
artoo_encodings	8
artoo_formats	10
artoo_spec	10
cdisc_adae	13
cdisc_adsl	13
cdisc_dm	14
cdisc_spec	14
cdisc_specs	15
cdisc_suppdm	16
cdisc_ts	17
cdisc_vs	17
check_spec	18
check_study	19
columns	21
conformance	22
decode_column	24
get_meta	26
is_artoo_checks	27
is_artoo_meta	28
is_artoo_spec	29
members	30
read_dataset	31
read_json	32
read_ndjson	34
read_parquet	35
read_rds	36
read_spec	38
read_xpt	39
repair_spec	41
set_meta	42
set_type	44
spec_codelists	45
spec_comments	46
spec_datasets	47
spec_documents	48
spec_keys	49

spec_methods	50
spec_standard	51
spec_study	52
spec_variables	53
sync_meta	54
validate_spec	55
write_dataset	57
write_json	58
write_ndjson	60
write_parquet	61
write_rds	63
write_spec	64
write_xpt	66
xpt_members	68

Index	70
--------------	-----------

apply_spec	<i>Conform a data frame to its spec</i>
------------	---

Description

Run the ordered, transactional artoo pipeline that turns a raw analysis data frame into one conformed to its specification and carrying artoo_meta. This is the middle of the workflow (spec -> apply_spec -> read_/write_): the conformed frame is ready for any write_*() codec, and the metadata it now carries makes that write lossless. The input is never mutated; if any step aborts, the call leaves your data untouched.

Usage

```
apply_spec(
  x,
  spec,
  dataset,
  conformance = c("warn", "abort", "off"),
  na_position = c("first", "last"),
  extra = c("keep", "drop"),
  on_coercion_loss = c("error", "keep")
)
```

Arguments

x	<i>The raw data frame to conform.</i> <data.frame>: required.
spec	<i>The specification to conform to.</i> <artoo_spec>: required.
dataset	<i>The dataset whose rules apply.</i> <character(1)>: required. Must name a dataset in spec.
conformance	<i>What to do with conformance findings.</i> <character(1)>. One of:

- "warn" (default) run `check_spec()`, attach the findings (read them with `conformance()`), warn on any error-severity finding.
- "abort" abort with `ar too_error_conformance` on any error-severity finding.
- "off" skip the check entirely.

Note: this governs only the *findings* disposition — what is *reported*. Pipeline errors are a different category and abort under every setting, including "off": an unknown dataset, and lossy coercion (`ar too_error_type`). Lossy coercion has its own governed gate, `on_coercion_loss`, not this argument; when the spec is the problem, retype it with `set_type()`.

`na_position` *Where missing key values sort.* <character(1)>. One of "first" (default) or "last". "first" matches SAS PROC SORT (and the FDA submission convention) by ordering missings before present values; "last" matches R's `order()` and the pandas/Polars default. Both are lossless; pick the one your comparison target uses.

`extra` *What happens to undeclared columns.* <character(1)>. An "extra" is a column of `x` the spec does not declare — typically a derivation temporary. One of:

- "keep" (default) extras ride along after the declared columns, reported by the `extra_variable` finding.
- "drop" the returned frame carries exactly the spec's columns; the drop is announced (`ar too_message_apply`), and that message is the audit trail of what was removed.

Interaction: the drop runs *before* the check, so `conformance()` reports only the columns the returned frame keeps. Under `conformance = "abort"` an error-severity finding still aborts (those findings arise only on spec-declared columns, which the drop never touches) and the input is never mutated, so the trim cannot mask a failure.

Note: "keep" is the default deliberately. `ar too` is a lossless carrier, so the metadata step never silently discards a column; extras are surfaced every run (the `extra_variable` finding, and a warning under `conformance = "warn"`), making "drop" a conscious opt-in.

`on_coercion_loss`

What to do when coercion would lose data. <character(1)>. The governed gate for an integer data type whose data truncates (fractions) or overflows (R's 32-bit range). One of:

- "error" (default) abort with `ar too_error_type` before any value is touched, refusing to damage the data.
- "keep" skip coercion for the offending column, leaving it at its wider source type. The values are preserved and the mismatch is reported as an `integer_fraction / integer_overflow` finding (read with `conformance()`), never silently truncated.

Interaction: independent of `conformance`. "error" aborts even under `conformance = "off"`; under "keep" the finding it leaves is surfaced by `conformance = "warn"` (the default) and suppressed by "off".

Tip: "keep" is the iterate stance (preserve the data, flag the spec); "error" is the submission stance. To fix the spec itself, see `set_type()` and `repair_spec()`.

Details

Ordered pipeline. Four fixed steps run in order: coerce each column to its CDISC dataType, reorder columns to the spec, sort rows by the dataset keys, then stamp the metadata. A spec variable the data lacks is never fabricated as an empty column: artoo is a lossless carrier, not a deriver. It is reported instead, an informational heads-up at apply time plus a `missing_variable` finding (when mandatory) or `missing_permmissible` (when not), and left absent, so the conformed frame carries only the columns the data actually had.

Extras are kept by default. A column the spec does not declare survives the pipeline (ordered after the declared ones), is *reported* by the `extra_variable` conformance finding, and round-trips through every `write_*`() codec with metadata inferred from its R class — membership reported, never enforced by silent destruction. Keeping is the default because artoo is lossless by construction: a metadata-application step that silently discarded columns would break that contract, so trimming data is always an explicit, announced choice rather than a default side effect. `extra = "drop"` opts in to trim-to-spec (the returned frame carries exactly the spec's columns): the undeclared columns are removed *before* the check, so the findings describe exactly the returned frame (a dropped column is never reported as `extra_variable`), and the drop itself is always announced (`artoo_message_apply`) as the audit trail of what was removed — even under `conformance = "off"`.

Lossless or abort, your call. A coercion that would damage values — an integer dataType truncating fractions or overflowing R's 32-bit range — aborts with `artoo_error_type` before any value is touched, under the default `on_coercion_loss = "error"`. This gate is independent of conformance: `conformance = "off"` does not bypass it. When the data (not the spec) is right, set `on_coercion_loss = "keep"`: the column keeps its wider source type and the divergence is reported as an `integer_fraction / integer_overflow` finding, never silently truncated. When the spec is wrong, retype it with `set_type()` (or `repair_spec()` from the findings). The error abort carries the offending rows as data: `cnid$variables` is a data frame with columns `variable`, `data_type`, `n`, and `reason` ("truncated" / "overflowed"), so a pipeline can collect every mismatch in one `tryCatch(..., artoo_error_type = function(cnid) cnid$variables)` pass. The NA-introduction warning (`artoo_warning_coercion`) carries the same frame with `reason = "na_introduced"`, and a `conformance = "abort"` failure carries the complete findings frame as `cnid$findings`.

Values are never translated. Coded variables keep their submission values (SEX stays "M"); codelist translation is its own verb, `decode_column()`.

Value

A *conformed* <data.frame> carrying `artoo_meta` (read it with `get_meta()`) and, unless `conformance = "off"`, the findings frame `conformance()` reads back. Hand it to any `write_*`() codec.

See Also

Check: `check_spec()` for the findings; `conformance()` to read them back.

Fix the spec: `set_type()` to retype a variable the data disagrees with, `repair_spec()` to apply every integer fix from a findings frame.

Translate: `decode_column()` for codelist value mapping.

Metadata: `get_meta()` / `set_meta()` for what the stamp attaches.

Examples

```
# ---- Example 1: conform ADSL, then read its metadata ----
#
# The bundled adam_spec describes ADSL; the raw frame is coerced,
# ordered, sorted, and stamped with the CDISC metadata get_meta() reads
# back. Variables the spec declares but this extract never derived are
# reported (not added), readable via conformance().
adsl <- apply_spec(cdisc_adsl, adam_spec, "ADSL")
get_meta(adsl)@dataset$records

# ---- Example 2: extras are kept and reported, or dropped on request ----
#
# By default a column outside the spec rides along (reported by the
# extra_variable finding) and still writes losslessly; extra = "drop"
# trims to the spec, announced and still reported. DM is SDTM, so it
# conforms against the bundled sdtm_spec.
raw <- cdisc_dm
raw$DERIVED <- seq_len(nrow(raw))
dm <- apply_spec(raw, sdtm_spec, "DM")
findings <- conformance(dm)
findings[findings$check == "extra_variable", c("variable", "message")]
trimmed <- apply_spec(raw, sdtm_spec, "DM", extra = "drop")
"DERIVED" %in% names(trimmed)
```

artoo_checks

Control which conformance checks run

Description

Build a reusable control that selects which dimensions `check_spec()` evaluates. Construct one per study and thread it through every `check_spec()` call so the conformance surface is consistent. Each toggle is validated at construction, so a mistyped name or value aborts early rather than being silently ignored.

Usage

```
artoo_checks(
  missing_variable = TRUE,
  missing_permissible = TRUE,
  extra_variable = TRUE,
  type_mismatch = TRUE,
  length_overflow = TRUE,
  char_length_limit = TRUE,
  codelist_membership = TRUE,
  codelist_membership_extensible = TRUE,
  label_match = TRUE,
  key_uniqueness = TRUE,
```

```

display_format = TRUE,
variable_name = TRUE,
dataset_name = TRUE,
label_length = TRUE,
integer_overflow = TRUE,
integer_fraction = TRUE,
iso8601_format = TRUE
)

```

Arguments

missing_variable

Flag mandatory spec variables absent from the data. <logical(1)>: default TRUE.

missing_permissible

Flag permissible (non-mandatory) spec variables absent from the data. <logical(1)>: default TRUE.

extra_variable *Flag data columns the spec does not declare.* <logical(1)>: default TRUE.

type_mismatch *Flag columns whose storage differs from the spec dataType.* <logical(1)>: default TRUE.

length_overflow

Flag character values longer than the spec length. <logical(1)>: default TRUE.

char_length_limit

Flag character values longer than the SAS XPORT v5 / FDA 200-byte limit.
<logical(1)>: default TRUE.

codelist_membership

Flag values outside their closed codelist. <logical(1)>: default TRUE.

codelist_membership_extensible

Flag values outside an extensible codelist's enumerated terms. <logical(1)>: default TRUE.
A codelist whose extended flag is TRUE allows sponsor terms, so a non-member is a note, never an error; this toggle silences those notes independently of codelist_membership.

label_match *Flag a column whose label attribute differs from the spec label.* <logical(1)>: default TRUE.

key_uniqueness *Flag a dataset whose spec key variables do not uniquely identify its rows.* <logical(1)>: default TRUE.

display_format *Flag a date/datetime/time variable whose displayFormat is not a recognized SAS format of that family.* <logical(1)>: default TRUE.

variable_name *Flag a data column name that violates the XPORT naming rules.* <logical(1)>: default TRUE.
Over 8 characters (the v5 limit), over 32 (the v8 limit), or containing anything but ASCII letters, digits, and underscore.

dataset_name *Flag a dataset name that violates the XPORT naming rules.* <logical(1)>: default TRUE.
Same limits as variable_name.

label_length *Flag a column label attribute over the 40-byte XPORT v5 / FDA limit.* <logical(1)>: default TRUE.

integer_overflow

Flag an integer-typed variable holding values beyond R's 32-bit integer range.
<logical(1)>: default TRUE. Such values become NA under coercion, so this is an error, not a warning.

integer_fraction

Flag an integer-typed variable holding fractional values. <logical(1)>: default TRUE.
Coercion would truncate them (162.6 becomes 162) — a data-integrity event; fix the spec dataType (float / decimal) or the data before conforming.

`iso8601_format` *Flag a character date/datetime/time variable whose values are not valid ISO 8601 text.* `<logical(1)>`: default TRUE. A character column under a temporal data type is the CDISC --DTC form; complete values, right-truncated partials ("1951", "1951-12"), and SDTMIG hyphen placeholders ("2003---15") all pass, while "12NOV2019" or an impossible calendar date is flagged.

Details

Selection, not severity. This control decides which findings are *produced*; `apply_spec()`'s conformance argument (warn, abort, off) decides what to *do* with the findings its full-default check raises. A disabled dimension is skipped entirely, so the findings frame stays clean.

Value

A `<artoo_checks>` control object. Pass it as the checks argument to `check_spec()`.

See Also

`check_spec()`, which consumes it; `apply_spec()` for the findings disposition.

Examples

```
# ---- Example 1: the default runs every conformance dimension ----
#
# With no arguments, every conformance dimension is enabled.
artoo_checks()

# ---- Example 2: silence one dimension for a whole study ----
#
# Turn off the length check (e.g. while a spec's lengths are provisional)
# and reuse the control across every dataset.
spec <- artoo_spec(cdisc_sdtm_datasets, cdisc_sdtm_variables, codelists = cdisc_codelists)
ck <- artoo_checks(length_overflow = FALSE)
nrow(check_spec(cdisc_dm, spec, "DM", checks = ck))
```

Description

List the character encodings clinical data actually travels in, with the name each ecosystem uses for the same thing: the R name (the standard IANA name, which `iconv()` and the wider R ecosystem use), the SAS session-encoding name, and the Python codec. Any spelling from the `r` or `sas` column works as the encoding argument of every artoo reader and writer.

Usage

```
artoo_encodings()
```

Details

What an encoding is. Text is stored as bytes; an encoding is the rule that maps those bytes to characters. Plain A-Z digits and punctuation are the same bytes in every encoding listed here — the differences only show in accented letters (a-umlaut, e-acute), special symbols (micro, degree), and non-Latin scripts. Reading bytes with the wrong rule is what turns a degree sign into garbage.

Which one do I have? In SAS, run `PROC OPTIONS OPTION=ENCODING; RUN;` and look up the reported name in the `sas` column. Most US/EU Windows SAS installs report `WLATIN1` — that is `windows-1252` here.

Which one should I write? Usually none: `write_*(encoding = NULL)` (the default) inherits the encoding recorded when the data was read, so a round-trip is byte-faithful. The regulatory defaults artoo applies when nothing is recorded: SAS XPORT writes US-ASCII (the FDA Study Data Technical Conformance Guide expectation) and Dataset-JSON / NDJSON write UTF-8 (required by CDISC and RFC 8259). A value that cannot be represented in the target encoding aborts loudly — see `on_invalid` on the writers.

Note: in memory, artoo text is always UTF-8 (NFC-normalised) — encodings only matter at the file boundary, exactly as in Python 3.

Value

A `<data.frame>` with one row per encoding and columns `r` (the R name — the standard IANA name `iconv()` uses, and what artoo records in the metadata), `sas` (the SAS session-encoding name), `python` (the Python codec name), and `description`.

See Also

Use it: the encoding argument of `read_xpt()`, `write_xpt()`, `read_json()`, and the other readers/writers.

Formats: `artoo_formats()` for the codec registry.

Examples

```
# ---- Example 1: the full cross-ecosystem table ----
#
# One row per encoding; the same byte rule under each ecosystem's name.
artoo_encodings()

# ---- Example 2: look up a SAS session encoding ----
#
# PROC OPTIONS reported WLATIN1: find the R and Python names for the
# same bytes (the sas and r spellings both work as encoding=).
enc <- artoo_encodings()
enc[enc$sas == "WLATIN1", ]
```

`artoo_formats`*Report which formats are available*

Description

List every registered codec and whether it can read and write in this session. The pure-R formats (xpt, json, rds) are always available; optional-engine formats (parquet) report FALSE until their package is installed. Purely informational, modelled on the diagnostic helpers in the wider ecosystem; it never aborts.

Usage

```
artoo_formats()
```

Value

A `<data.frame>` with one row per format and columns `format`, `read`, `write` (logical), and `extensions`.

See Also

[read_dataset\(\)](#) and [write_dataset\(\)](#) which use the registry.

Examples

```
# ---- Example 1: see what this session can read and write ----  
#  
# rds is always available; the table shows the extensions each codec claims.  
artoo_formats()
```

`artoo_spec`*Construct a CDISC specification*

Description

Build and validate a `artoo_spec` from `dataset`, `variable`, and `codelist` tables. Each table is coerced to a plain data frame, missing optional columns are filled with typed NAs, every variable type is canonicalised to the CDISC `dataType` vocabulary, and cross-slot integrity (dataset and codelist references) is checked before the object is returned. The spec is the lingua franca the rest of artoo reads, applies, and serialises.

Usage

```
artoo_spec(
  datasets = NULL,
  variables = NULL,
  codelists = NULL,
  study = NULL,
  values = NULL,
  methods = NULL,
  comments = NULL,
  documents = NULL,
  standard = NULL
)
```

Arguments

datasets	<i>Dataset-level metadata table.</i> <data.frame>: required. One row per dataset; must carry a dataset column. Optional columns label, class, structure, keys are filled with NA when absent.
variables	<i>Variable-level metadata table.</i> <data.frame>: required. One row per variable; must carry dataset, variable, and data_type. The data_type column is canonicalised to a CDISC dataType (e.g. "text" becomes "string"). Requirement: every dataset value must appear in datasets.
codelists	<i>Controlled-terminology terms.</i> <data.frame> NULL. Must carry codelist_id and term when supplied. Interaction: every codelist_id referenced by variables must resolve here.
study	<i>Study-level metadata.</i> <data.frame> NULL. A single row of named study fields. Well-known fields are canonicalised to study_name, study_description, and protocol_name (aliases such as StudyName or studyid resolve automatically); other fields pass through verbatim. A standard field, when present, is consumed into @standard.
values	<i>Value-level (VLM) metadata.</i> <data.frame> NULL.
methods	<i>Derivation methods.</i> <data.frame> NULL. The Define-XML method definitions variables reference by method_id; must carry method_id when supplied. Completeness (e.g. a referenced method has a description) is checked by validate_spec() , not here.
comments	<i>Comment definitions.</i> <data.frame> NULL. Referenced by comment_id; must carry comment_id when supplied.
documents	<i>Document references.</i> <data.frame> NULL. Referenced by document_id; must carry document_id when supplied.
standard	<i>The CDISC standard the spec implements.</i> <character(1)> NULL. E.g. "AdMIG 1.1" or "SDTMIG 3.2". When NULL (default) it is resolved from datasets\$standard or study\$standard; absent everywhere, @standard is NA. Restriction: all sources must agree on one value; conflicting standards abort with artoo_error_spec.

Details

Coerce, then validate. Each table is first coerced to a plain data frame (a tibble is accepted and demoted); known columns are cast to their storage mode and absent optional columns are added as typed NA, so every downstream reader can trust the schema. Validation runs only after coercion, on the completed slots.

Type canonicalisation. `variables$data_type` is mapped through the closed CDISC `dataType` vocabulary (`string`, `integer`, `decimal`, `float`, `double`, `boolean`, `date`, `datetime`, `time`, `URI`). Common SAS / P21 spellings resolve automatically ("`text`", "`Char`", "`integer (8)`", ...); an unrecognised token aborts with `artoo_error_type`.

Cross-slot integrity. Construction fails (`artoo_error_spec`) if a variable names a dataset absent from datasets, or references a `codelist_id` absent from codelists.

One spec, one standard. A `artoo_spec` carries exactly one CDISC standard, stored as the scalar `@standard` property. The constructor resolves it from the `standard` argument, a `standard` column in datasets (the P21 workbook shape), and a `standard` field in `study` (the Define-XML shape) — those columns are consumed, so `@standard` is the single home. More than one distinct value aborts with `artoo_error_spec`; scope the source to one standard (e.g. `read_spec(path, datasets = ...)`) instead of mixing.

One study vocabulary. Well-known study fields are canonicalised to the CDISC ODM `GlobalVariables` names, `snake_cased`: `study_name`, `study_description`, `protocol_name`. Source spellings resolve automatically (`StudyName`, `studyid`, ...); fields the vocabulary does not know pass through verbatim. Aliases that disagree on a value abort with `artoo_error_spec`.

Value

A *validated* `artoo_spec` object. Inspect it with `spec_datasets()` / `spec_variables()`, or check it with `validate_spec()`.

See Also

Inspect: `spec_datasets()`, `spec_variables()`, `spec_codelists()`, `spec_keys()`, `spec_study()`.

Check: `validate_spec()`. **Predicate:** `is_artoo_spec()`.

Examples

```
# ---- Example 1: build a spec from the bundled CDISC-pilot tables ----
#
# `cdisc_sdtm_datasets` and `cdisc_sdtm_variables` hold the CDISC pilot SDTM
# metadata in the shape artoo_spec() expects; the constructor
# canonicalises every type and checks cross-slot integrity.
spec <- artoo_spec(cdisc_sdtm_datasets, cdisc_sdtm_variables, codelists = cdisc_codelists)
spec_datasets(spec)

# ---- Example 2: a focused spec for a single dataset ----
#
# Slice the bundled tables to one dataset (DM) to build a smaller spec.
dm_ds <- cdisc_sdtm_datasets[cdisc_sdtm_datasets$dataset == "DM", ]
dm_var <- cdisc_sdtm_variables[cdisc_sdtm_variables$dataset == "DM", ]
dm_spec <- artoo_spec(dm_ds, dm_var, codelists = cdisc_codelists)
```

```
head(spec_variables(dm_spec, "DM")[, c("variable", "label", "data_type")])
```

`cdisc_adae`*Demo adverse events analysis dataset (ADaM ADAE)*

Description

A 60-row sample of the CDISC pilot ADaM adverse events analysis dataset (ADAE): one row per reported event, with treatment-emergent flags, severity, and coding variables (labels preserved as attributes).

Usage

```
cdisc_adae
```

Format

A data frame with 60 rows (STUDYID, USUBJID, AETERM, AESEV, TRTEMFL, ASTDT, ...).

Source

First 60 rows of the CDISC pilot adam/cdisc/adae.xpt from the PHUSE Test Data Factory (phuse-org/phuse-scripts).

`cdisc_adsl`*Demo subject-level analysis dataset (ADaM ADSL)*

Description

A 60-subject sample of the CDISC pilot ADaM subject-level analysis dataset (ADSL): one row per subject, with treatment, demographic, baseline, and disposition variables (labels preserved as column attributes).

Usage

```
cdisc_adsl
```

Format

A data frame with 60 rows and 48 variables (STUDYID, USUBJID, TRT01P, AGE, SEX, RACE, SAFFL, TRTSDT, ...).

Source

First 60 subjects of the CDISC pilot adam/cdisc/adsl.xpt from the PHUSE Test Data Factory (phuse-org/phuse-scripts).

cdisc_dm	<i>Demo demographics dataset (SDTM DM)</i>
----------	--

Description

A 60-subject sample of the CDISC pilot SDTM demographics domain (DM): one row per subject, with the standard DM variables (labels preserved as attributes).

Usage

```
cdisc_dm
```

Format

A data frame with 60 rows and 25 variables (STUDYID, DOMAIN, USUBJID, AGE, SEX, RACE, ARM, COUNTRY, ...).

Source

First 60 subjects of the CDISC pilot sdtm/TDF_SDTM_v1.0/dm.xpt from the PHUSE Test Data Factory (phuse-org/phuse-scripts).

cdisc_spec	<i>CDISC demo specification tables (one standard per pair)</i>
------------	--

Description

The constructor-shaped metadata tables for the bundled demo data, split by CDISC standard because a `artoo_spec` carries exactly one: `cdisc_adam_datasets` + `cdisc_adam_variables` describe ADaM (ADaMIG 1.1), and `cdisc_sdtm_datasets` + `cdisc_sdtm_variables` describe DM (SDTMIG 3.1.2). Each variables table is *derived from the data* (names, labels, inferred CDISC types, byte lengths) by `data-raw/`. Pass one standard's pair to `artoo_spec()`; passing both pairs together aborts with `artoo_error_spec` — mixing standards in one spec is the mistake the split exists to prevent.

Usage

```
cdisc_adam_datasets
cdisc_adam_variables
cdisc_sdtm_datasets
cdisc_sdtm_variables
cdisc_codelists
```

Format

Each *_datasets table is a data frame with one row per dataset:

dataset Dataset name ("ADSL" or "DM").

label Dataset label.

standard The CDISC standard, consumed into the spec's `spec_standard()`.

Each *_variables table is a data frame with one row per variable:

dataset Owning dataset name.

variable Variable name.

label Variable label (from the data's label attribute).

data_type CDISC dataType inferred from the column's class.

length Storage length (max byte width for character, 8 for numeric).

order Variable order within the dataset.

codelist_id NCI codelist reference ("C66731" on SEX).

cdisc_codelists is a data frame of controlled-terminology terms (the real NCI codelist C66731 for SEX):

codelist_id Codelist identifier ("C66731").

term Submission value ("M", "F", ...).

decode Decoded value ("Male", "Female", ...).

order Term order.

Source

Derived from the CDISC pilot .xpt files in the public PHUSE Test Data Factory (phuse-org/phuse-scripts) by data-raw/bundle-demo.R.

cdisc_specs

Bundled CDISC specifications (ADaM and SDTM)

Description

Ready-made artoo_spec objects built from the official CDISC Define-XML 2.1 release examples: adam_spec (ADaMIG 1.1; datasets ADSL, ADAE) and sdtm_spec (SDTMIG 3.1.2; datasets TS, DM, VS, SUPPDM). Every bundled demo dataset conforms to its spec under `apply_spec(conformance = "abort")` — the pairing is gated at build time. The same specs ship as P21 workbooks in `system.file("extdata", "adam-spec.xlsx", package = "artoo")` and `"sdtm-spec.xlsx"`, written by `write_spec()`.

Usage

adam_spec

sdtm_spec

Format

A validated `artoo_spec()` object; inspect it with `spec_datasets()`, `spec_variables()`, and `spec_standard()`.

Details

Demo adaptations (each an ADR in `data-raw/bundle-spec.R`): the sponsor-defined codelists (`CL.ARM`, `CL.ARMCD`, `CL.BMICAT`, and the extensible NCI VS codelists) are marked extended; `VISITNUM` is typed `float` (the pilot data has fractional visit numbers); `VS` declares the `SDTMIG` timepoint variables `VSTPT/VSTPTNUM`, with `VSTPTNUM` in the `VS` key.

Source

The CDISC Define-XML 2.1 release example defines (ADaM + SDTM), pinned by `sha256` in `data-raw/bundle-spec.R`; data from the PHUSE Test Data Factory.

`cdisc_suppdm`*Demo supplemental qualifiers dataset (SDTM SUPPDM)*

Description

A 60-row sample of the CDISC pilot SDTM supplemental qualifiers for DM (SUPPDM): the non-standard qualifier values that ride alongside the DM domain.

Usage`cdisc_suppdm`**Format**

A data frame with 60 rows (`STUDYID`, `RDOMAIN`, `USUBJID`, `QNAM`, `QVAL`, ...).

Source

First 60 rows of the CDISC pilot `sdm/cdiscpilot01/suppdm.xpt` from the PHUSE Test Data Factory (`phuse-org/phuse-scripts`).

cdisc_ts	<i>Demo trial summary dataset (SDTM TS)</i>
----------	---

Description

The CDISC pilot SDTM trial summary domain (TS): one row per trial characteristic (33 rows in the pilot), the study-design parameters a submission carries.

Usage

cdisc_ts

Format

A data frame with 33 rows (STUDYID, TSPARMCD, TSPARM, TSVAL, ...).

Source

The CDISC pilot sdtm/cdiscpilot01/ts.xpt from the PHUSE Test Data Factory (phuse-org/phuse-scripts).

cdisc_vs	<i>Demo vital signs dataset (SDTM VS)</i>
----------	---

Description

A 60-row sample of the CDISC pilot SDTM vital signs domain (VS): repeated measurements per subject across visits, positions, and planned timepoints.

Usage

cdisc_vs

Format

A data frame with 60 rows (STUDYID, USUBJID, VSTESTCD, VSORRES, VISITNUM, VSPOS, VSTPTNUM, ...).

Source

First 60 rows of the CDISC pilot sdtm/cdiscpilot01/vs.xpt from the PHUSE Test Data Factory (phuse-org/phuse-scripts).

check_spec	<i>Check a dataset against its spec</i>
------------	---

Description

Compare a data frame to one dataset's specification and report where they diverge. This is the data-conformance check at the end of the artoo workflow (spec -> apply_spec -> check_spec): it reuses the metadata the spec already carries (variables, types, lengths, codelists, keys). It is distinct from [validate_spec\(\)](#), which checks the spec's own internal integrity rather than the data. Both report findings keyed to the same open rule catalog.

Usage

```
check_spec(
  x,
  spec,
  dataset,
  decode = c("none", "to_decode", "to_code"),
  checks = NULL
)
```

Arguments

x	<i>The data frame to check.</i> <data.frame>: required. Typically the output of apply_spec() , but any frame works.
spec	<i>The specification to check against.</i> <artoo_spec>: required.
dataset	<i>The dataset whose rules apply.</i> <character(1)>: required. Restriction: must name a dataset in spec (see spec_datasets()).
decode	<i>Which codelist column membership is checked against.</i> <character(1)>. One of "none" (default), "to_decode", or "to_code".
checks	<i>Which conformance dimensions to evaluate.</i> <artoo_checks> NULL. When NULL (default) every dimension runs; build a control with artoo_checks() to disable some.

Details

Findings, not enforcement. `check_spec()` never modifies data; it returns every divergence it finds. [apply_spec\(\)](#) runs it and decides what to do via its conformance argument (warn, abort, off). The dimensions checked are: missing variables (split into mandatory, an error, and permissible, a warning), extra variables (data column the spec does not declare), type mismatch, ISO 8601 validity of character date/datetime/time values (CDISC partials pass; "12NOV2019" does not), fractional values and 32-bit overflow under an integer dataType (both would corrupt data at coercion), character length overflow, the hard 200-byte XPORT v5 / FDA character limit, codelist membership, label drift against the spec, key uniqueness, and displayFormat validity.

Decode-aware membership. `decode` selects which codelist column the data is checked against, matching [apply_spec\(\)](#)'s `decode` step: "none"/"to_code" check against the codelist terms, "to_decode"

against the decodes. `apply_spec()` threads its own decode through, so a decoded column is not wrongly flagged.

Fatal vs informational coercion checks. Only `integer_fraction` and `integer_overflow` carry error severity: they mark data an integer `data` type cannot hold without loss, which `apply_spec()` refuses to coerce (its `on_coercion_loss` governs that gate). `type_mismatch` is a note: a column stored more widely than the spec declares (an integer-valued double, for instance) coerces cleanly, so it is informational, not a blocker.

Value

A *findings data frame* with columns `check`, `dimension`, `severity` ("error", "warning", or "note"), `dataset`, `variable`, and `message`, one row per divergence. Zero rows means the data conforms.

See Also

`apply_spec()` which runs this; `check_study()` for the same check across a whole study; `artoo_checks()` to select dimensions; `validate_spec()` for spec integrity.

Examples

```
# ---- Example 1: a conformed frame surfaces only the genuine gaps ----
#
# apply_spec() coerces and orders to spec but never fabricates a variable
# the data lacks; checking the result reports the permissible variables
# this extract never derived (here, six) instead of hiding them as empty
# columns.
adsl <- apply_spec(cdisc_adsl, adam_spec, "ADSL", conformance = "off")
nrow(check_spec(adsl, adam_spec, "ADSL"))

# ---- Example 2: raw data surfaces divergences ----
#
# Checking a raw frame with an undeclared column flags the extras.
raw <- cdisc_adsl
raw$NOTASPEC <- 1
head(check_spec(raw, adam_spec, "ADSL")[, c("check", "variable", "severity")])
```

check_study

Check a whole study against its spec

Description

Run `check_spec()` over every dataset in a study and return one stacked findings frame. Where `check_spec()` answers "does this dataset conform?", `check_study()` answers "is my whole study submittable?" in a single pass, surfacing every dataset's divergences at once instead of one abort at a time. The result is an ordinary findings frame underneath, so filter it by severity or hand it straight to `repair_spec()`.

Usage

```
check_study(
  spec,
  data,
  decode = c("none", "to_decode", "to_code"),
  checks = NULL
)
```

Arguments

spec	<i>The specification to check against.</i> <artoo_spec>: required.
data	<i>The study's datasets.</i> <named list of data.frame>: required. One entry per dataset, named by the dataset (e.g. list(ADSL = adsl, ADAE = adae)). Every name must be a dataset in spec.
decode	<i>Which codelist column to check against.</i> <character(1)>. Passed to check_spec() ; one of "none" (default), "to_decode", "to_code".
checks	<i>Which conformance dimensions to run.</i> <artoo_checks> NULL. Passed to check_spec() ; NULL (default) runs every dimension. Build a subset with artoo_checks() .

Details

One row per divergence, every dataset stacked. Each dataset's findings carry its name in the dataset column, so the frame is the union of the per-dataset [check_spec\(\)](#) results. Printing renders the dataset-by-check count matrix (the study-level summary); the underlying frame is unchanged.

Data-requiring, like [check_spec\(\)](#). [check_study\(\)](#) checks data against the spec, so it needs the data frames. For the spec's own structural integrity (no data), use [validate_spec\(\)](#).

Value

A <artoo_study_findings> *data frame* with the same columns as [check_spec\(\)](#) (check, dimension, severity, dataset, variable, message), one row per divergence across all datasets. Zero rows means the whole study conforms. Print it for the count matrix; treat it as an ordinary data frame otherwise.

See Also

One dataset: [check_spec\(\)](#). **Spec structure only:** [validate_spec\(\)](#).

Repair: [repair_spec\(\)](#) to apply the integer fixes the matrix surfaces.

Examples

```
# ---- Example 1: scan a whole study in one pass ----
#
# Loop the conformance check over every dataset's data. A fractional AGE
# (the spec types it integer) surfaces as an integer_fraction finding; the
# print is a dataset-by-check count matrix.
```

```

adsl <- cdisc_adsl
adsl$AGE <- adsl$AGE + 0.5
check_study(adam_spec, list(ADSL = adsl, ADAE = cdisc_adae))

# ---- Example 2: feed the findings straight into repair_spec() ----
#
# The result is an ordinary findings frame, so repair_spec() consumes it to
# flip every integer_fraction / integer_overflow variable across the study.
findings <- check_study(adam_spec, list(ADSL = adsl))
fixed <- repair_spec(adam_spec, findings)
spec_variables(fixed, "ADSL")$data_type[
  spec_variables(fixed, "ADSL")$variable == "AGE"
]

```

columns

View a dataset's variable attributes, SAS-style

Description

Return a one-row-per-variable attribute table — the pane a SAS programmer reads in PROC CONTENTS or the Universal Viewer: position, name, Char/Num type, length, format, informat, label, and the CDISC key sequence. This is the quick look after `apply_spec()` stamps a frame, or on any dataset file artoo can read.

Usage

```
columns(x, member = NULL)
```

Arguments

x	<i>What to describe.</i> <data.frame> <character(1)>: required. A stamped frame (carries artoo_meta), any plain data frame, or a path to a dataset file (.xpt, .json, .ndjson, .parquet, .rds).
member	<i>XPORT member to describe.</i> <character(1)> NULL. Only meaningful when x is a path to a multi-member .xpt file.

Details

Every real column shows. The table covers the *frame's* columns: a column the spec never declared (which `apply_spec()` keeps, never drops) still appears, its attributes inferred from the R class. A plain, never-stamped data frame works the same way — every attribute is inferred.

Len is physical storage. The pane mirrors what PROC CONTENTS shows and what the writers store, not a spec digit-width. A Char column always carries a byte Len (the declared length, else inferred from the data); a numeric Len is blank, because a numeric stores as an 8-byte IEEE double with no character width (a Define-XML numeric Length is a digit-width, kept in the metadata for the Define / P21 surface). Format and informat names render uppercase; the metadata keeps the source spelling.

A path reads through the codec. A file path is dispatched by extension through the same registry as `read_dataset()`, so the attributes come from the one lossless reader (an unknown extension aborts with the registry's known-extensions message).

Tip: a multi-member XPORT file needs `member =`; without one the xpt reader aborts and points at `xpt_members()` for the listing.

Note: an `.xpt` path shows a blank Key: the XPORT byte layout stores only name, label, length, and formats, so `keySequence` (like `codelist` and `origin`) cannot ride in the file. The metadata-carrying formats (`.json`, `.ndjson`, `.parquet`, `.rds`) and the in-session conformed frame show it; re-apply the spec after an xpt read to restore it.

Value

A `<artoo_columns> data frame` with columns `#`, `Variable`, `Type`, `Len`, `Format`, `Label`, `Key`, printed left-aligned. The `Informat` column appears only when at least one variable carries an informat (most clinical panes have none). It is an ordinary data frame underneath — filter or inspect it like one.

See Also

Members: `xpt_members()` lists a multi-member XPORT file.

Metadata: `get_meta()` for the full `artoo_meta`; `apply_spec()` which stamps it.

Examples

```
spec <- artoo_spec(cdisc_adam_datasets, cdisc_adam_variables, codelists = cdisc_codelists)

# ---- Example 1: the column pane of a conformed frame ----
#
# apply_spec() stamps ADSL with its metadata; columns() reads it back as
# the SAS-style attribute table.
adsl <- apply_spec(cdisc_adsl, spec, "ADSL", conformance = "off")
columns(adsl)

# ---- Example 2: straight off a file ----
#
# Write the conformed frame to any format and point columns() at the
# path; the codec reads it back and the attributes are identical.
p <- tempfile(fileext = ".json")
write_json(adsl, p)
columns(p)
```

Description

Pull the conformance findings `apply_spec()` attached to a conformed data frame — the readable answer to "what did the check find?". The result is the same findings frame `check_spec()` returns (one row per divergence), with a print method that renders a sectioned report, so `conformance(ads1)` at the console is the inspection step the `artoo_warning_conformance` warning points you at.

Usage

```
conformance(x)
```

Arguments

`x` *A data frame produced by `apply_spec()`.* `<data.frame>`: required.
Requirement: the conformance check must have run: a frame from `apply_spec(..., conformance = "off")` (or one rebuilt by a transform that dropped attributes) carries no findings and aborts with `artoo_error_input`.

Value

A `<artoo_findings>` *data frame* with columns `check`, `dimension`, `severity` ("error", "warning", or "note"), `dataset`, `variable`, and `message`. Zero rows means the data conformed. Print it for the sectioned report; treat it as an ordinary data frame for programmatic use.

See Also

`apply_spec()` which attaches the findings; `check_spec()` for the same check on demand; `artoo_checks()` to select dimensions.

Examples

```
spec <- artoo_spec(
  cdisc_sdtm_datasets, cdisc_sdtm_variables,
  codelists = cdisc_codelists
)

# ---- Example 1: inspect what the conform step found ----
#
# Conforming raw DM records the findings on the result; conformance()
# renders them as a report instead of a raw attribute.
dm <- suppressWarnings(apply_spec(cdisc_dm, spec, "DM"))
conformance(dm)

# ---- Example 2: gate a pipeline on error-severity findings ----
#
# The findings frame is an ordinary data frame: filter by severity to
# drive your own logic.
f <- conformance(dm)
nrow(f[f$severity == "error", ])
```

 decode_column

Derive or translate a variable through its codelist

Description

Map one column's values through a spec codelist — code to decode or decode to code — writing the result to a new variable or in place. This is the everyday companion to `apply_spec()`'s whole-dataset decode step: deriving RACEN from RACE, recovering submission codes from decoded values, or decoding a single variable for display, without re-running the pipeline. When the target variable is declared in the spec, the result is also coerced to its `dataType` and labelled, so the new column lands conformed.

Usage

```
decode_column(
  x,
  spec,
  dataset,
  from,
  to = from,
  direction = c("to_decode", "to_code"),
  no_match = c("error", "keep", "na"),
  trim = TRUE,
  ignore_case = FALSE
)
```

Arguments

<code>x</code>	<i>The data frame to extend.</i> <code><data.frame></code> : required.
<code>spec</code>	<i>The specification carrying the codelists.</i> <code><artoo_spec></code> : required.
<code>dataset</code>	<i>The dataset whose variables apply.</i> <code><character(1)></code> : required. Must name a dataset in spec.
<code>from</code>	<i>The source column.</i> <code><character(1)></code> : required. Must be a column of <code>x</code> .
<code>to</code>	<i>The destination variable.</i> <code><character(1)></code> : default <code>from</code> . Defaults to translating in place. A <code>to</code> declared in the spec gets its <code>dataType</code> coercion and label; an undeclared <code>to</code> is plain character.
<code>direction</code>	<i>Which way to map.</i> <code><character(1)></code> . One of: <ul style="list-style-type: none"> <code>"to_decode"</code> (default) map codes to their decoded values ("M" becomes "Male"). <code>"to_code"</code> map decoded values to their submission codes — the RACEN-from-RACE derivation.
<code>no_match</code>	<i>Policy for values absent from the codelist.</i> <code><character(1)></code> . One of <code>"error"</code> (default; abort with <code>artoo_error_codelist</code> , naming the unmatched values and the codelist), <code>"keep"</code> (carry the source value through), or <code>"na"</code> .

`trim` *Match after trimming whitespace.* <logical(1)>: default TRUE.

`ignore_case` *Match case-insensitively.* <logical(1)>: default FALSE. Case differences are usually genuine CT violations, so this is opt-in.

Details

Which codelist applies. The codelist attached to `to` in the spec wins (the natural direction for RACEN-style derivations, where the numeric variable owns the code/decode pairs); when `to` declares none, `from`'s codelist is used. If neither variable references a codelist the call aborts — there is nothing to map through.

Mismatched surfaces chain. A single call maps through ONE codelist, so the winning codelist's terms (or decodes) must line up with the `from` values — the CDISC *N convention guarantees this for RACEN-style pairs, whose decodes are the character variable's submission values. When the two codelists share no value surface (say `SEXN`'s decodes are "Female"/"Male" but `SEX` holds "F"/"M"), the unmatched values hit the `no_match` policy; translate in two hops instead — decode through `from`'s codelist first, then `to_code` through the destination's:

```
dm |>
  decode_column(spec, "DM", from = "SEX", to = "SEXDECD") |>
  decode_column(spec, "DM", from = "SEXDECD", to = "SEXN",
                direction = "to_code")
```

Soft matches are reported, never silent. Values that match only after trimming whitespace (or case-folding, when `ignore_case = TRUE`) still map, with a `artoo_warning_codelist` naming the variants — `check_spec()` always compares exactly, so clean the source for submission.

Value

The data frame `x` with the `to` column added (at the end) or replaced (in place), ready for the next pipeline step.

See Also

Whole-dataset decode: `apply_spec()` with `decode =`.

Inspect the terms: `spec_codelists()`. **Check membership:** `check_spec()`.

Examples

```
spec <- artoo_spec(cdisc_sdtm_datasets, cdisc_sdtm_variables, codelists = cdisc_codelists)

# ---- Example 1: decode a coded variable into a display column ----
#
# SEX is coded against C66731; map the codes to their decodes in a new
# column, leaving the submission values untouched.
dm <- decode_column(cdisc_dm, spec, "DM", from = "SEX", to = "SEXDECD")
table(dm$SEX, dm$SEXDECD)

# ---- Example 2: the RACEN pattern, a coded numeric from its decode ----
#
```

```

# Declare SEXN as an integer variable owning a numeric codelist, then
# derive it from SEX's decoded values: to_code maps each decode to its
# submission code, and the spec dataType makes the result integer.
vars <- rbind(
  cdisc_sdtm_variables,
  data.frame(
    dataset = "DM", variable = "SEXN", label = "Sex (N)",
    data_type = "integer", length = 8L, order = NA_integer_,
    codelist_id = "SEXN"
  )
)
cls <- rbind(
  cdisc_codelists,
  data.frame(
    codelist_id = "SEXN", term = c("1", "2"),
    decode = c("F", "M"), order = 1:2
  )
)
spec_n <- artoo_spec(cdisc_sdtm_datasets, vars, codelists = cls)
dm_n <- decode_column(cdisc_dm, spec_n, "DM",
  from = "SEX", to = "SEXN", direction = "to_code"
)
str(dm_n$SEXN)

```

get_meta

Read the metadata a dataset carries

Description

Pull the `artoo_meta` off a data frame produced by `apply_spec()` or read back by any `read_*`() codec. The metadata travels as a single Dataset-JSON string in the frame's `metadata_json` attribute; `get_meta()` parses it to the S7 object, the form every codec writes from. This is the read half of the lossless round-trip.

Usage

```
get_meta(x)
```

Arguments

`x` *A data frame carrying artoo metadata.* `<data.frame>`: required. Typically the output of `apply_spec()` or a `read_*`() codec.

Requirement: `x` must carry a `metadata_json` attribute (set by `set_meta()`, `apply_spec()`, or a reader); a bare frame aborts with `artoo_error_input`.

Value

A `<artoo_meta>` with two properties. `@dataset` is a named list of dataset-level attributes: `itemGroupOID`, `name`, `label`, `records`, `studyOID`, `metaDataVersionOID`, `encoding`, and `keys`. `@columns` is a named list with one entry per variable, each carrying `itemOID`, `name`, `label`, `dataType`, `targetDataType`, `length`, `displayFormat`, `informat`, `keySequence`, `codelist`, `significantDigits`, and `origin` (absent values are `NULL`). Pass it to `set_meta()` to re-attach, or index it directly (`meta@columnsAGElabel`).

See Also

`set_meta()` for the write half; `apply_spec()` which stamps it.

Examples

```
# ---- Example 1: read metadata off a conformed dataset ----
#
# apply_spec() stamps the metadata; get_meta() reads it back as the S7
# object whose @columns holds one CDISC attribute set per variable.
spec <- artoo_spec(cdisc_adam_datasets, cdisc_adam_variables, codelists = cdisc_codelists)
adsl <- apply_spec(cdisc_adsl, spec, "ADSL")
meta <- get_meta(adsl)
meta@columns$STUDYID

# ---- Example 2: round-trip metadata across two frames ----
#
# The metadata is a portable object: read it off one frame and stamp it
# onto another with set_meta().
bare <- as.data.frame(adsl)
attr(bare, "metadata_json") <- NULL
restamped <- set_meta(bare, meta)
identical(get_meta(restamped)@columns, meta@columns)
```

is_artoo_checks

Test for a artoo_checks control

Description

Report whether an object is a `artoo_checks` control built by `artoo_checks()`. Use it to guard a `checks` argument before threading it into `check_spec()` or `apply_spec()`.

Usage

```
is_artoo_checks(x)
```

Arguments

`x` *Object to test. <any>*.

Value

A `<logical(1)>`: TRUE when x is a `artoo_checks`.

See Also

[artoo_checks\(\)](#) to build one.

Examples

```
# ---- Example 1: confirm a control before reusing it ----
#
# is_artoo_checks() distinguishes a real control from a bare list of flags.
is_artoo_checks(artoo_checks())
is_artoo_checks(list(missing_variable = TRUE))
```

is_artoo_meta

Test for a artoo_meta object

Description

Report whether an object is a `artoo_meta` — the CDISC-shaped metadata a conformed dataset carries through the artoo workflow (`spec -> apply_spec -> read_/write_`). [get_meta\(\)](#) returns one; this is the type guard before you inspect its `@dataset` and `@columns` slots.

Usage

```
is_artoo_meta(x)
```

Arguments

x *Object to test. <any>*.

Value

A `<logical(1)>`: TRUE when x is a `artoo_meta`, else FALSE.

See Also

[get_meta\(\)](#) and [set_meta\(\)](#) to read and attach metadata.

Examples

```
# ---- Example 1: guard before inspecting metadata ----
#
# get_meta() yields a artoo_meta; is_artoo_meta() confirms the type before
# you reach into its slots.
spec <- artoo_spec(cdisc_adam_datasets, cdisc_adam_variables, codelists = cdisc_codelists)
adsl <- apply_spec(cdisc_adsl, spec, "ADSL")
meta <- get_meta(adsl)
is_artoo_meta(meta)

# ---- Example 2: a bare data frame carries no meta object ----
#
# The raw frame itself is not a artoo_meta – only the object get_meta()
# returns is.
is_artoo_meta(cdisc_adsl)
```

is_artoo_spec

Test for a artoo_spec object

Description

Report whether an object is a `artoo_spec` — the validated CDISC specification that drives the artoo workflow (`spec` -> `apply_spec` -> `read_/write_`). `artoo_spec()` builds one; this is the type guard before you pass it to `apply_spec()` or reach into it with the spec accessors.

Usage

```
is_artoo_spec(x)
```

Arguments

x *Object to test. <any>.*

Value

A `<logical(1)>`: TRUE when x is a `artoo_spec`, else FALSE.

See Also

`artoo_spec()` to build one; `is_artoo_meta()` for the metadata guard.

Examples

```
# ---- Example 1: guard a built specification ----
#
# artoo_spec() assembles and validates a spec; is_artoo_spec() confirms the
# type before you drive apply_spec() with it.
spec <- artoo_spec(cdisc_sdtm_datasets, cdisc_sdtm_variables, codelists = cdisc_codelists)
```

```
is_artoo_spec(spec)

# ---- Example 2: an ordinary object is not a spec ----
#
# Any non-artoo_spec value – a bare data frame, say – returns FALSE.
is_artoo_spec(cdisc_dm)
```

members

List the datasets in a file or directory

Description

Inventory the dataset(s) a path contains, one row per dataset, dispatched by extension through the same codec registry as `read_dataset()`. A SAS XPORT library lists every member; a single-dataset file (`.json`, `.ndjson`, `.parquet`, `.rds`) reports one row; a directory inventories each dataset file it holds. The format-neutral companion to the xpt-specific `xpt_members()`.

Usage

```
members(path)
```

Arguments

`path` *A dataset file or a directory.* `<character(1)>`: required. A path to a dataset file (`.xpt`, `.json`, `.ndjson`, `.parquet`, `.rds`) or to a directory holding such files. A path that does not exist, or a file whose extension no codec claims, aborts.

Details

One dataset per file, except XPORT. XPORT is the only multi-dataset container artoo handles, so only an `.xpt` path can return more than one row. Every other format is one dataset per file.

A directory is inventoried, not descended. Only the files directly in the directory are listed (no recursion); files whose extension no codec claims are skipped, and a directory with no dataset files returns an empty inventory rather than aborting. A dataset file that fails to read aborts with its codec's error, naming the file.

Note: counting records reads the file through its codec (the one lossless reader), so `members()` is an honest count, not a header guess; for a large directory it reads every dataset.

Value

A `<artoo_members>` *data frame*, one row per dataset, with columns `file` (source basename), `member` (dataset name), `label`, `records` (row count), `variables` (column count), and `format` (the codec format). Empty when a directory holds no dataset files. It is an ordinary data frame underneath.

See Also

Members of one XPORT file: [xpt_members\(\)](#).

Per-variable attributes: [columns\(\)](#) for one dataset's variable pane.

Examples

```
dm <- apply_spec(cdisc_dm, sdtm_spec, "DM", conformance = "off")

# ---- Example 1: one dataset in a file ----
#
# A single-dataset format reports exactly one member.
p <- tempfile(fileext = ".json")
write_json(dm, p)
members(p)

# ---- Example 2: every dataset in a directory ----
#
# Point members() at a folder to inventory each dataset file it holds, one
# row per dataset, dispatched by extension.
dir <- tempfile("datasets")
dir.create(dir)
write_json(dm, file.path(dir, "dm.json"))
write_rds(dm, file.path(dir, "dm.rds"))
members(dir)
```

read_dataset

Read a dataset from any supported format

Description

Read a clinical file back to a data frame, restoring its `artoo_meta`. The codec is chosen from the file extension (or an explicit format), and the metadata the file carries is re-attached, so a value written by [write_dataset\(\)](#) round-trips losslessly. This is the ingest end of the I/O layer; the per-format wrappers like [read_rds\(\)](#) call it.

Usage

```
read_dataset(path, format = NULL, col_select = NULL, n_max = Inf, ...)
```

Arguments

path	<i>Source file path.</i> <character(1)>: required. Its extension selects the codec unless format is given.
format	<i>Force a codec instead of inferring from the extension.</i> <character(1)> NULL. One of the registered formats (see artoo_formats()).

col_select	<i>Variables to read.</i> <character> NULL. NULL (default) reads every column; otherwise a vector of variable names. Columns return in file order (not the requested order) and the artoo_meta is filtered to match. Works on every format: parquet narrows columns natively, the rest filter after decode. Note: an unknown name is a artoo_error_input, never a silent drop.
n_max	<i>Maximum records to read.</i> <numeric(1)>: default Inf. Caps the row count; the returned artoo_meta reports the rows actually read. xpt v8 bounds the disk read; the other formats cap after decode.
...	<i>Codec-specific arguments</i> passed through to the decoder (see the per-format wrappers, e.g. read_xpt()). An argument the codec does not know is an error, never silently ignored.

Value

A <data.frame> carrying artoo_meta when the file recorded it (read it with [get_meta\(\)](#)). A file whose payload is not a data frame is a artoo_error_codec.

See Also

[write_dataset\(\)](#) for the inverse; [read_rds\(\)](#) for the per-format wrapper.

Examples

```
spec <- artoo_spec(cdisc_adam_datasets, cdisc_adam_variables, codelists = cdisc_codelists)

# ---- Example 1: round-trip a dataset through rds ----
#
# Write a conformed dataset, then read it back; the metadata survives.
adsl <- apply_spec(cdisc_adsl, spec, "ADSL", conformance = "off")
path <- tempfile(fileext = ".rds")
write_dataset(adsl, path)
back <- read_dataset(path)
identical(get_meta(back)@columns, get_meta(adsl)@columns)

# ---- Example 2: the metadata names the dataset and row count ----
#
# The restored artoo_meta exposes the dataset-level attributes.
get_meta(back)@dataset$records
```

read_json

Read a dataset from CDISC Dataset-JSON

Description

Read a CDISC Dataset-JSON v1.1 (.json) file back to a data frame, restoring the full artoo_meta it carries and realizing SAS date/datetime/time variables to R Date / POSIXct / hms: :hms. Column types are reconstructed from the recorded metadata, not guessed from the JSON tokens, so the round-trip is lossless. The ingest end of the I/O layer; a thin wrapper over [read_dataset\(\)](#) with format = "json".

Usage

```
read_json(path, col_select = NULL, n_max = Inf, encoding = NULL)
```

Arguments

path	<i>Source .json path.</i> <character(1)>: required. A JSON file that is not Dataset-JSON v1.1 aborts with <code>artoo_error_codec</code> .
col_select	<i>Variables to read.</i> <character> NULL. NULL (default) reads every column; otherwise a vector of variable names. Columns return in file order (not the requested order) and the <code>artoo_meta</code> is filtered to match. Works on every format: parquet narrows columns natively, the rest filter after decode. Note: an unknown name is a <code>artoo_error_input</code> , never a silent drop.
n_max	<i>Maximum records to read.</i> <numeric(1)>: default <code>Inf</code> . Caps the row count; the returned <code>artoo_meta</code> reports the rows actually read. <code>xpt v8</code> bounds the disk read; the other formats cap after decode.
encoding	<i>Source charset of the file bytes.</i> <character(1)> NULL. NULL (default) reads UTF-8, as Dataset-JSON requires. Pass an IANA or SAS charset name (e.g. "windows-1252") only to read a non-conformant file a producer wrote in that charset; the bytes are transcoded to UTF-8 on read. Tip: any SAS or IANA spelling listed by <code>artoo_encodings()</code> is accepted.

Value

A <data.frame> carrying `artoo_meta` (read it with `get_meta()`).

See Also

`write_json()` for the inverse; `read_dataset()` for the generic dispatcher.

Examples

```
spec <- artoo_spec(cdisc_adam_datasets, cdisc_adam_variables, codelists = cdisc_codelists)

# ---- Example 1: round-trip a conformed dataset through Dataset-JSON ----
#
# The variable labels, types, and keys survive the round-trip.
adsl <- apply_spec(cdisc_adsl, spec, "ADSL", conformance = "off")
path <- tempfile(fileext = ".json")
write_json(adsl, path)
back <- read_json(path)
identical(get_meta(back)@columns, get_meta(adsl)@columns)

# ---- Example 2: the metadata names the dataset and row count ----
#
# The restored artoo_meta exposes the dataset-level attributes.
get_meta(back)@dataset$records
```

read_ndjson

Read a dataset from CDISC Dataset-JSON NDJSON

Description

Read a newline-delimited CDISC Dataset-JSON v1.1 (.ndjson) file back to a data frame, restoring the full `artoo_meta` from its metadata line and realizing SAS date/datetime/time variables to R Date / POSIXct / hms: :hms. Rows are parsed in bounded slabs, and `n_max` stops the line loop early, so a partial read of a huge file never parses the tail. A thin wrapper over `read_dataset()` with `format = "ndjson"`.

Usage

```
read_ndjson(path, col_select = NULL, n_max = Inf, encoding = NULL)
```

Arguments

path	<i>Source .ndjson path.</i> <character(1)>: required. A gzip stream (.ndjson.gz) is inflated transparently. A file whose first line is not the Dataset-JSON metadata object aborts with <code>artoo_error_codec</code> .
col_select	<i>Variables to read.</i> <character> NULL. NULL (default) reads every column; otherwise a vector of variable names. Columns return in file order (not the requested order) and the <code>artoo_meta</code> is filtered to match. Works on every format: parquet narrows columns natively, the rest filter after decode. Note: an unknown name is a <code>artoo_error_input</code> , never a silent drop.
n_max	<i>Maximum records to read.</i> <numeric(1)>: default Inf. Caps the row count; the returned <code>artoo_meta</code> reports the rows actually read. xpt v8 bounds the disk read; the other formats cap after decode.
encoding	<i>Source charset of the file bytes.</i> <character(1)> NULL. NULL (default) reads UTF-8, as Dataset-JSON requires. Pass an IANA or SAS charset name (e.g. "windows-1252") only to read a non-conformant file a producer wrote in that charset; each line is transcoded to UTF-8 on read, preserving the bounded <code>n_max</code> streaming.

Value

A <data.frame> carrying `artoo_meta` (read it with `get_meta()`).

See Also

`write_ndjson()` for the inverse; `read_json()` for the array-form file; `read_dataset()` for the generic dispatcher.

Examples

```
spec <- artoo_spec(cdisc_adam_datasets, cdisc_adam_variables, codelists = cdisc_codelists)

# ---- Example 1: round-trip a conformed dataset through NDJSON ----
#
# The variable labels, types, and keys survive the round-trip.
adsl <- apply_spec(cdisc_adsl, spec, "ADSL", conformance = "off")
path <- tempfile(fileext = ".ndjson")
write_ndjson(adsl, path)
back <- read_ndjson(path)
identical(get_meta(back)@columns, get_meta(adsl)@columns)

# ---- Example 2: a bounded partial read of the first rows ----
#
# n_max stops the line loop as soon as enough rows are in.
head_rows <- read_ndjson(path, n_max = 5)
get_meta(head_rows)@dataset$records
```

read_parquet

*Read a dataset from Apache Parquet***Description**

Read an Apache Parquet (.parquet) file back to a data frame, restoring the `artoo_meta` from its `metadata_json` sidecar and realizing SAS date/datetime/time variables to R Date / POSIXct / hms : : hms. A parquet written by another tool (with no `artoo` sidecar) reads back as a bare frame. A thin wrapper over `read_dataset()` with `format = "parquet"`. Requires the lightweight `nanoparquet` package.

Usage

```
read_parquet(path, col_select = NULL, n_max = Inf, encoding = NULL)
```

Arguments

<code>path</code>	<i>Source .parquet path.</i> <character(1)>: required.
<code>col_select</code>	<i>Variables to read.</i> <character> NULL. NULL (default) reads every column; otherwise a vector of variable names. Columns return in file order (not the requested order) and the <code>artoo_meta</code> is filtered to match. Works on every format: <code>parquet</code> narrows columns natively, the rest filter after decode. Note: an unknown name is a <code>artoo_error_input</code> , never a silent drop.
<code>n_max</code>	<i>Maximum records to read.</i> <numeric(1)>: default <code>Inf</code> . Caps the row count; the returned <code>artoo_meta</code> reports the rows actually read. <code>xpt v8</code> bounds the disk read; the other formats cap after decode.

encoding *Source charset of the string columns.* <character(1)> | NULL. NULL (default) reads the UTF-8 bytes parquet stores. Pass a charset name only to read a foreign file whose string columns hold that charset's bytes; they are transcoded to UTF-8 on read.

Tip: any SAS or IANA spelling listed by [artoo_encodings\(\)](#) is accepted.

Value

A <data.frame> carrying `artoo_meta` when the file recorded it (read it with [get_meta\(\)](#)); otherwise a plain data frame.

See Also

[write_parquet\(\)](#) for the inverse; [read_dataset\(\)](#) for the generic dispatcher.

Examples

```
spec <- artoo_spec(cdisc_adam_datasets, cdisc_adam_variables, codelists = cdisc_codelists)

# ---- Example 1: round-trip a conformed dataset through Parquet ----
#
# The variable labels, types, and keys survive the round-trip.
adsl <- apply_spec(cdisc_adsl, spec, "ADSL", conformance = "off")
path <- tempfile(fileext = ".parquet")
write_parquet(adsl, path)
back <- read_parquet(path)
get_meta(back)@columns$STUDYID$label

# ---- Example 2: the metadata names the dataset and row count ----
#
# The restored artoo_meta exposes the dataset-level attributes.
get_meta(back)@dataset$records
```

read_rds

Read a dataset from rds

Description

Read an R .rds file written by [write_rds\(\)](#) (or any rds carrying a `metadata_json` attribute) back to a data frame with its `artoo_meta` restored. A thin wrapper over [read_dataset\(\)](#) with `format = "rds"`.

Usage

```
read_rds(path, col_select = NULL, n_max = Inf, encoding = NULL)
```

Arguments

path	<i>Source .rds path.</i> <character(1)>: required.
col_select	<i>Variables to read.</i> <character> NULL. NULL (default) reads every column; otherwise a vector of variable names. Columns return in file order (not the requested order) and the artoo_meta is filtered to match. Works on every format: parquet narrows columns natively, the rest filter after decode. Note: an unknown name is a artoo_error_input, never a silent drop.
n_max	<i>Maximum records to read.</i> <numeric(1)>: default Inf. Caps the row count; the returned artoo_meta reports the rows actually read. xpt v8 bounds the disk read; the other formats cap after decode.
encoding	<i>Source charset of the string columns.</i> <character(1)> NULL. NULL (default) returns the strings exactly as saved (faithful R round-trip). Pass a charset name only to transcode a foreign rds whose string columns hold that charset's bytes. Tip: any SAS or IANA spelling listed by artoo_encodings() is accepted.

Value

A <data.frame> carrying artoo_meta when the file recorded it. An rds holding anything other than a data frame is a artoo_error_codec; use readRDS() for arbitrary objects.

See Also

[write_rds\(\)](#) for the inverse; [read_dataset\(\)](#) for the generic dispatcher.

Examples

```
spec <- artoo_spec(
  cdisc_adam_datasets, cdisc_adam_variables,
  codelists = cdisc_codelists
)

# ---- Example 1: read a dataset written by write_rds() ----
#
# The restored frame carries the same metadata it was written with.
adsl <- apply_spec(cdisc_adsl, spec, "ADSL", conformance = "off")
path <- tempfile(fileext = ".rds")
write_rds(adsl, path)
back <- read_rds(path)
get_meta(back)@dataset$records

# ---- Example 2: a plain rds still reads as a data frame ----
#
# An rds without artoo metadata reads back as an ordinary frame.
bare <- tempfile(fileext = ".rds")
saveRDS(cdisc_dm, bare)
nrow(read_rds(bare))
```

read_spec

*Read a specification from JSON, Excel, or Define-XML***Description**

Read a clinical-dataset specification into a validated `artoo_spec`, dispatching on the file extension: `artoo`'s native JSON (the inverse of `write_spec()`), a Pinnacle 21 (P21) Excel workbook, or a native Define-XML 2.0/2.1 document. The returned spec is the lingua franca the rest of `artoo` applies and serialises.

Usage

```
read_spec(path, datasets = NULL, on_duplicate = c("error", "first", "warn"))
```

Arguments

path	<i>The specification file to read.</i> <character(1)>: required. A <code>.json</code> (native) or <code>.xlsx / .xls</code> (P21) file. Requirement: reading a P21 workbook needs the <code>readxl</code> package.
datasets	<i>Read only these datasets.</i> <character> NULL. NULL (default) reads the whole spec. Otherwise the spec is scoped to the named datasets before validation, so one broken sheet elsewhere in a workbook cannot block the dataset you are working on. An unknown name aborts listing what the file defines.
on_duplicate	<i>Policy for a variable defined more than once.</i> <character(1)>. A workbook row duplicated within one dataset makes the spec ambiguous; the finding is reported with its source location (sheet and row numbers for Excel). One of: <ul style="list-style-type: none"> • "error" (default) abort, naming each duplicate's rows. • "first" keep the first definition of each, dropping the rest with a message. • "warn" keep the first definition and warn (<code>artoo_warning_spec</code>).

Details

Three formats, one validator. A `.json` file is read as `artoo` native JSON; a `.xlsx / .xls` file is read as a P21 workbook; a `.xml` file is read as Define-XML 2.x. Either way the result is built through `artoo_spec()`, so type canonicalisation and cross-slot integrity checks are identical regardless of source.

Define-XML ingestion (needs the `xml2` package). `ItemGroupDefs` become datasets (keys derived from the `ItemRef KeySequence`), `ItemRef + ItemDef` pairs become variables, `CodeLists` become codelists (`def:ExtendedValue = "Yes"` marks an extended term), `MethodDefs / CommentDefs / leaves` become the supporting slots, and `ValueListDefs` land in the value-level slot with their where-clauses rendered as readable text.

Note: an `ExternalCodeList` (MedDRA, ISO-3166) names a dictionary, not an enumerable membership list; it is dropped, and variables that referenced it carry no codelist. Define-XML v1.0 (the 2005 model) is refused with guidance.

P21 ingestion. Sheets are located by a tolerant alias match (case-, space-, and spelling-variant insensitive). Datasets and Variables are required; Codelists and ValueLevel are optional (the latter becomes the spec's value-level slot). Every cell is read as text, then the dataset and codelist foreign keys are forward-filled to recover merged cells (which the Excel reader returns as NA on continuation rows). A key that cannot be resolved aborts with `artoo_error_spec` rather than being silently dropped.

Value

A *validated* `artoo_spec`. Inspect it with `spec_datasets()` / `spec_variables()`, check it with `validate_spec()`, or persist it with `write_spec()`.

See Also

Inverse: `write_spec()` serialises a spec to native JSON.

Build / inspect: `artoo_spec()`, `spec_datasets()`, `spec_variables()`, `validate_spec()`.

Examples

```
# ---- Example 1: round-trip a spec through native JSON ----
#
# write_spec() and read_spec() are inverses on the JSON path: the spec
# that comes back is identical to the one written.
spec <- artoo_spec(cdisc_sdtm_datasets, cdisc_sdtm_variables, codelists = cdisc_codelists)
path <- tempfile(fileext = ".json")
write_spec(spec, path)
back <- read_spec(path)
identical(back, spec)

# ---- Example 2: scope the read to one dataset ----
#
# `datasets =` reads just the domain you are working on – validation is
# scoped with it, so a problem elsewhere in the workbook cannot block
# this dataset.
dm_spec <- read_spec(path, datasets = "DM")
spec_datasets(dm_spec)
head(spec_variables(dm_spec, "DM")[, c("variable", "label", "data_type")])
```

read_xpt

Read a dataset from SAS XPORT

Description

Read a SAS Transport (.xpt) file (v5 or v8) back to a data frame, restoring the `artoo_meta` its NAMESTR records carry and realizing SAS date/datetime/time variables to R Date / POSIXct / hms: :hms. The ingest end of the I/O layer; a thin wrapper over `read_dataset()` with `format = "xpt"`.

Usage

```
read_xpt(path, encoding = NULL, col_select = NULL, n_max = Inf, member = NULL)
```

Arguments

path	<i>Source .xpt path.</i> <character(1)>: required.
encoding	<i>Force a source charset.</i> <character(1)> NULL. NULL (default) auto-detects (UTF-8 when every character value and label is valid UTF-8, else Windows-1252). IANA and SAS names both work. Tip: any SAS or IANA spelling listed by artoo_encodings() is accepted.
col_select	<i>Variables to read.</i> <character> NULL. NULL (default) reads every column; otherwise a vector of variable names (matching the names as stored, uppercase for v5). Columns return in file order, and the <code>artoo_meta</code> is filtered to match. Note: an unknown name is a <code>artoo_error_input</code> , never a silent drop.
n_max	<i>Maximum records to read.</i> <numeric(1)>: default Inf. Caps the row count; the returned <code>artoo_meta</code> reports the rows actually read.
member	<i>Which member of a multi-member transport file to read.</i> <character(1) numeric(1)> NULL. A transport file can hold several datasets; pass a member name (case-insensitive) or 1-based index to pick one. NULL (default) reads a single-member file directly and aborts on a multi-member file, pointing at xpt_members() . Tip: <code>xpt_members(path)</code> lists what a file holds before you choose.

Details

The character encoding is auto-detected (UTF-8 if every character value is valid UTF-8, else Windows-1252) and recorded on the returned `artoo_meta`, so a later [write_xpt\(\)](#) reproduces it; pass `encoding` to override. XPORT cannot record its own encoding, so this detection is a heuristic. See [write_xpt\(\)](#) for what XPORT can and cannot preserve.

Value

A <data.frame> carrying `artoo_meta` (read it with [get_meta\(\)](#)).

See Also

[xpt_members\(\)](#) to list a file's members; [write_xpt\(\)](#) for the inverse; [read_dataset\(\)](#) for the generic dispatcher.

Examples

```
spec <- artoo_spec(
  cdisc_adam_datasets, cdisc_adam_variables,
  codelists = cdisc_codelists
)

# ---- Example 1: round-trip a conformed dataset through xpt ----
#
# Write ADSL, read it back; the variable labels and lengths survive.
```

```

adsl <- apply_spec(cdisc_adsl, spec, "ADSL", conformance = "off")
path <- tempfile(fileext = ".xpt")
write_xpt(adsl, path)
back <- read_xpt(path)
get_meta(back)@columns$STUDYID$label

# ---- Example 2: pick one member of a multi-member transport file ----
#
# Build a two-member file by concatenating two single-member files (every
# member section is 80-byte padded), then read one dataset out of it.
dm <- apply_spec(cdisc_dm, sdtm_spec, "DM", conformance = "off")
p_dm <- tempfile(fileext = ".xpt")
write_xpt(dm, p_dm)
multi <- tempfile(fileext = ".xpt")
writeBin(
  c(
    readBin(path, "raw", file.size(path)),
    readBin(p_dm, "raw", file.size(p_dm))[-(1:240)]
  ),
  multi
)
xpt_members(multi)$name
nrow(read_xpt(multi, member = "DM"))

```

 repair_spec

Repair a spec from its conformance findings

Description

Take the `integer_fraction` and `integer_overflow` findings a `check_spec()` or `check_study()` run reports and return a new spec with every offending variable retyped to "float", so a frame that the original spec would refuse to coerce now conforms. This closes the loop on the spec-side fix: inspect the findings, then apply them all at once instead of editing the source workbook variable by variable. Persist the result with `write_spec()`.

Usage

```
repair_spec(spec, findings)
```

Arguments

spec	<i>The specification to repair.</i> <artoo_spec>: required.
findings	<i>A findings data frame.</i> <data.frame>: required. The result of <code>check_spec()</code> or <code>check_study()</code> ; must carry the check, dataset, and variable columns.

Details

Scope. Only the two lossy-integer findings are repaired (`integer_fraction`, `integer_overflow`) — both mean "the spec says integer but the data is not", and "float" is the loss-free fix. Other findings are ignored; this is not a general spec rewriter. When no repairable finding is present the spec is returned unchanged, with a note.

Built on `set_type()`. Each (dataset, variable) pair is applied through the same validated override `set_type()` uses, so the result is a fully re-validated `artoo_spec`, never a hand-edited internal.

Value

A new `<artoo_spec>` with the flagged variables retyped to "float", or spec unchanged when there is nothing to repair. The input is never mutated.

See Also

Primitive: `set_type()` to retype a chosen variable directly.

Findings: `check_spec()` for one dataset, `check_study()` across a study. **Persist:** `write_spec()`.

Examples

```
# ---- Example 1: auto-repair an integer/fractional mismatch ----
#
# adam_spec types ADSL.AGE as integer. Give it fractional ages and
# check_spec() raises an integer_fraction error; repair_spec() flips AGE
# (and only AGE) to float, and the corrected spec then applies cleanly.
dat <- cdisc_adsl
dat$AGE <- dat$AGE + 0.5
findings <- check_spec(dat, adam_spec, "ADSL")
fixed <- repair_spec(adam_spec, findings)
spec_variables(fixed, "ADSL")$data_type[
  spec_variables(fixed, "ADSL")$variable == "AGE"
]

# ---- Example 2: nothing to repair is a no-op ----
#
# The bundled data conforms, so its findings carry no integer_fraction or
# integer_overflow rows and the spec is returned unchanged.
clean <- check_spec(cdisc_adsl, adam_spec, "ADSL")
identical(repair_spec(adam_spec, clean), adam_spec)
```

Description

Stamp a `artoo_meta` onto a data frame as a single Dataset-JSON string in its `metadata_json` attribute. Every `write_*()` codec reads that string back with `get_meta()` and embeds it verbatim, so the metadata survives the trip to any format. Use it to attach metadata to a bare frame before a write, or to re-stamp after a tidyverse verb has dropped attributes.

Usage

```
set_meta(x, meta)
```

Arguments

<code>x</code>	<i>The data frame to stamp.</i> <code><data.frame></code> : required.
<code>meta</code>	<i>The metadata to attach.</i> <code><artoo_meta></code> : required. Usually from <code>get_meta()</code> or built by <code>apply_spec()</code> .

Value

The data frame `x`, with its `metadata_json` attribute set. Pass it on to a `write_*()` codec or back through `get_meta()`.

See Also

`get_meta()` for the read half; `apply_spec()` which stamps it.

Examples

```
# ---- Example 1: re-stamp metadata a dplyr verb would drop ----
#
# Conform a dataset, capture its metadata, then re-attach after an
# attribute-dropping transform so the write stays lossless.
spec <- artoo_spec(
  cdisc_adam_datasets, cdisc_adam_variables,
  codelists = cdisc_codelists
)
adsl <- apply_spec(cdisc_adsl, spec, "ADSL")
meta <- get_meta(adsl)
trimmed <- head(as.data.frame(adsl), 5)
attr(trimmed, "metadata_json") <- NULL
set_meta(trimmed, meta)

# ---- Example 2: borrow metadata from a conformed dataset ----
#
# A writer with a raw frame can lift metadata off a conformed dataset and
# stamp it onto the bare frame (DM is SDTM, so its spec is sdtm-shaped).
sdtm <- artoo_spec(
  cdisc_sdtm_datasets, cdisc_sdtm_variables,
  codelists = cdisc_codelists
)
meta_dm <- get_meta(apply_spec(cdisc_dm, sdtm, "DM"))
dm <- set_meta(cdisc_dm, meta_dm)
```

```
is_artoo_meta(get_meta(dm))
```

set_type	<i>Override a variable's dataType in a spec</i>
----------	---

Description

Return a new `artoo_spec` with one or more variables retyped. This is the supported, in-R way to correct a spec when the data disagrees with its declared `dataType` (e.g. a variable typed integer whose extract holds fractional values): fix it here rather than editing the source workbook, then drive `apply_spec()` with the corrected spec. The spec is immutable, so the original is never changed.

Usage

```
set_type(spec, dataset, ...)
```

Arguments

<code>spec</code>	<i>The specification to amend.</i> <code><artoo_spec></code> : required.
<code>dataset</code>	<i>The dataset whose variables to retype.</i> <code><character(1)></code> : required. Must name a dataset in spec.
<code>...</code>	<i>Named variable = type pairs.</i> Each name is a variable in dataset; each value is a CDISC <code>dataType</code> ("string", "integer", "decimal", "float", "double", "boolean", "date", "datetime", "time", "URI") or a recognised spelling of one. At least one pair is required and every argument must be named. Tip: to undo an integer <code>dataType</code> that the data does not satisfy, set "float" (IEEE double) or "decimal" (exact, exchanged as text).

Details

Per-dataset scope. A type is set only on the named dataset's row. A variable that appears in several datasets keeps its other rows' types; call `set_type()` once per dataset to change them all. Spec-wide consequences (a variable typed inconsistently across datasets) are a `validate_spec()` concern, not a construction error here.

Canonicalised, then validated. Each supplied type is mapped through the closed CDISC `dataType` vocabulary, so "Float", "decimal", and "text" all resolve; an unrecognised token aborts with `artoo_error_type`. The rebuilt spec is re-validated, so an override that would break the spec aborts with `artoo_error_spec`.

Value

A new `<artoo_spec>` with the named variables retyped, ready for `apply_spec()` or `write_spec()`. The input spec is unchanged.

See Also

Auto-repair: `repair_spec()` to apply every `integer_fraction / integer_overflow` fix from a findings frame at once.

Workflow: `apply_spec()` to conform with the corrected spec; `write_spec()` to persist it; `check_spec()` to find the mismatches.

Examples

```
# ---- Example 1: retype one variable the data disagrees with ----
#
# The bundled adam_spec types ADSL.AGE as integer. If an extract stored it
# with fractional values, retype it to float so apply_spec() coerces
# without loss. set_type() returns a new spec; the original is untouched.
fixed <- set_type(adam_spec, "ADSL", AGE = "float")
v <- spec_variables(fixed, "ADSL")
v[v$variable == "AGE", c("variable", "data_type")]

# ---- Example 2: retype several at once, original left intact ----
#
# Pass any number of variable = type pairs; canonical dataTypes and common
# spellings both resolve. The source spec is immutable, so adam_spec still
# reports AGE as its original type.
patched <- set_type(adam_spec, "ADSL", AGE = "decimal", TRTSDT = "date")
spec_variables(adam_spec, "ADSL")$data_type[
  spec_variables(adam_spec, "ADSL")$variable == "AGE"
]
```

spec_codelists

*Codelist terms***Description**

Return the controlled-terminology terms and decodes a spec carries: one codelist's terms when `codelist_id` names it, or the full `codelists` slot when `codelist_id` is `NULL`. Use it to inspect the values a coded variable is allowed to take before applying the spec. Mirrors the `spec_variables()` filter pattern.

Usage

```
spec_codelists(spec, codelist_id = NULL)
```

Arguments

`spec` *The specification to read.* `<artoo_spec>`: required.

`codelist_id` *The codelist to return.* `<character(1)>` | `NULL`. When `NULL` (default) the whole `codelists` table is returned.

Restriction: a non-`NULL` id must name a codelist present in the spec's `codelists` slot; an unknown id aborts with `artoo_error_input`.

Value

A data frame of *codelist terms*, one row per term: every term when `codelist_id` is NULL, else the named codelist's terms. Columns:

- `codelist_id` — the codelist identifier variables reference.
- `term` — the submission value (what conformed data carries).
- `decode` — the human-readable decoded value.
- `order` — display order within the codelist.
- `extended` — TRUE marks an extensible codelist (sponsor terms allowed; non-members down-grade to notes in `check_spec()`).
- `comment_id` — reference into the comments slot.

See Also

`spec_variables()` for which variables reference a codelist.

Examples

```
# ---- Example 1: the terms behind a coded variable ----
#
# SEX is coded against C66731; spec_codelists() returns the terms and their
# decodes that apply_spec() will enforce or decode.
spec <- artoo_spec(
  cdisc_adam_datasets, cdisc_adam_variables,
  codelists = cdisc_codelists
)
spec_codelists(spec, "C66731")

# ---- Example 2: the whole codelists table ----
#
# Called with no id, it returns every term across every codelist.
head(spec_codelists(spec))
```

spec_comments

Comment definitions in a spec

Description

Return the comment definitions a specification carries. Datasets, variables, value-level rows, and codelists reference these by `comment_id`; `validate_spec()` checks the references resolve and each referenced comment has a body.

Usage

```
spec_comments(spec)
```

Arguments

spec *The specification to read.* <artoo_spec>: required.

Value

A data frame of comment metadata, one row per comment, with all four columns: comment_id, description, document_id, pages. Empty when the spec defines no comments.

See Also

[spec_methods\(\)](#), [spec_documents\(\)](#), [validate_spec\(\)](#).

Examples

```
# ---- Example 1: the comments a spec defines ----
#
# Build a spec with one comment and read it back.
spec <- artoo_spec(
  data.frame(dataset = "ADSL"),
  data.frame(dataset = "ADSL", variable = "AGE", data_type = "integer"),
  comments = data.frame(
    comment_id = "C.AGE",
    description = "Age in years at informed consent.",
    stringsAsFactors = FALSE
  )
)
spec_comments(spec)
```

spec_datasets

Dataset names in a spec

Description

List the datasets a specification defines. The result is the set of names you pass as the dataset argument to the other accessors and to `apply_spec()`.

Usage

```
spec_datasets(spec)
```

Arguments

spec *The specification to read.* <artoo_spec>: required.

Value

A character vector of dataset names, de-duplicated and with NAs dropped. Empty when the spec has no datasets.

See Also

[spec_variables\(\)](#) for one dataset's variables; [spec_keys\(\)](#) for its sort keys.

Examples

```
# ---- Example 1: the datasets the pilot ADaM spec defines ----
#
# Build the spec from the bundled CDISC-pilot tables and list its
# datasets – the names you pass to the other accessors.
spec <- artoo_spec(
  cdisc_adam_datasets, cdisc_adam_variables,
  codelists = cdisc_codelists
)
spec_datasets(spec)
```

spec_documents

Document references in a spec

Description

Return the document references a specification carries. Methods and comments point to these by `document_id`.

Usage

```
spec_documents(spec)
```

Arguments

`spec` *The specification to read.* `<artoo_spec>`: required.

Value

A data frame of document metadata (`document_id`, `title`, `href`), one row per document. Empty when the spec defines none.

See Also

[spec_methods\(\)](#), [spec_comments\(\)](#), [validate_spec\(\)](#).

Examples

```
# ---- Example 1: the documents a spec defines ----
#
# Build a spec with one document reference and read it back.
spec <- artoo_spec(
  data.frame(dataset = "ADSL"),
  data.frame(dataset = "ADSL", variable = "AGE", data_type = "integer"),
```

```

documents = data.frame(
  document_id = "SAP",
  title = "Statistical Analysis Plan",
  stringsAsFactors = FALSE
)
)
spec_documents(spec)

```

spec_keys

Sort keys for a dataset

Description

Parse a dataset's sort keys into a character vector of variable names. These keys drive the sort step of `apply_spec()` and the keySequence written to each output format.

Usage

```
spec_keys(spec, dataset)
```

Arguments

`spec` *The specification to read.* <artoo_spec>: required.
`dataset` *The dataset whose keys to parse.* <character(1)>: required.
Restriction: must name a dataset in the spec.

Value

A character vector of key variable names, split from the dataset's keys cell (whitespace- or comma-separated). Empty when no keys are declared.

See Also

[spec_datasets\(\)](#) for the dataset names; [spec_variables\(\)](#) for the variables a key must reference.

Examples

```

# ---- Example 1: parse a dataset's sort keys ----
#
# Declare DM's keys, then read them back as the ordered vector apply_spec()
# sorts by. (STUDYID and USUBJID are real DM variables in the demo data.)
ds <- cdisc_sdtm_datasets
ds$keys[ds$dataset == "DM"] <- "STUDYID USUBJID"
spec <- artoo_spec(ds, cdisc_sdtm_variables, codelists = cdisc_codelists)
spec_keys(spec, "DM")

```

spec_methods

Derivation methods in a spec

Description

Return the method definitions a specification carries. Variables and value-level rows reference these by `method_id`; `validate_spec()` checks that every reference resolves and that each referenced method is complete (has a description).

Usage

```
spec_methods(spec)
```

Arguments

`spec` *The specification to read.* `<artoo_spec>`: required.

Value

A data frame of method metadata, one row per method, with all eight columns: `method_id`, `description`, `name`, `type`, `expression_context`, `expression_code`, `document_id`, `pages`. Empty when the spec defines no methods.

See Also

[spec_comments\(\)](#), [spec_documents\(\)](#), [validate_spec\(\)](#).

Examples

```
# ---- Example 1: the methods a spec defines ----
#
# Build a spec with one derivation method and read it back.
spec <- artoo_spec(
  data.frame(dataset = "ADSL"),
  data.frame(dataset = "ADSL", variable = "AGEGR1", data_type = "string"),
  methods = data.frame(
    method_id = "MT.AGEGR1",
    description = "Age group from AGE.",
    stringsAsFactors = FALSE
  )
)
spec_methods(spec)
```

spec_standard	<i>The CDISC standard a spec implements</i>
---------------	---

Description

Return the one CDISC standard the specification carries (e.g. "ADaMIG 1.1", "SDTMIG 3.2"). A `artoo_spec` is single-standard by construction — `artoo_spec()` aborts when its sources mix standards — so this is always a scalar; NA when no source named one.

Usage

```
spec_standard(spec)
```

Arguments

`spec` *The specification to read.* `<artoo_spec>`: required.

Value

A `<character(1)>`: the standard, or NA when unspecified.

See Also

[spec_study\(\)](#) for the rest of the study-level metadata; [artoo_spec\(\)](#) for how the standard is resolved.

Examples

```
# ---- Example 1: the standard set at construction ----
#
# Pass the standard explicitly (or let it resolve from a P21 workbook's
# Standard column / a Define-XML study block) and read it back.
spec <- artoo_spec(
  cdisc_adam_datasets, cdisc_adam_variables,
  codelists = cdisc_codelists,
  standard = "ADaMIG 1.1"
)
spec_standard(spec)

# ---- Example 2: unspecified resolves to NA ----
#
# A spec built without any standard source carries NA.
bare <- artoo_spec(
  cdisc_adam_datasets, cdisc_adam_variables,
  codelists = cdisc_codelists
)
spec_standard(bare)
```

spec_study	<i>Study-level metadata</i>
------------	-----------------------------

Description

Return the study-level metadata row, or a single field from it. Holds the canonical study fields (`study_name`, `study_description`, `protocol_name` — every source spelling is canonicalised by `artoo_spec()`) plus any other study-scoped fields a source provides (the CDISC standard lives on its own property — see `spec_standard()`).

Usage

```
spec_study(spec, field = NULL)
```

Arguments

<code>spec</code>	<i>The specification to read.</i> <code><artoo_spec></code> : required.
<code>field</code>	<i>Return one field instead of the row.</i> <code><character(1)></code> <code>NULL</code> . When <code>NULL</code> (default) the whole study data frame is returned. Restriction: a non- <code>NULL</code> field must be a column of the study table; an unknown field aborts with <code>artoo_error_input</code> .

Value

The study data frame (one row), or the value of one field. The canonical fields are `study_name`, `study_description`, and `protocol_name` — every source spelling is canonicalised to these by `artoo_spec()` — plus any other field the source carried verbatim (e.g. `define_version` from a Define-XML read).

See Also

`spec_datasets()` for the datasets the study scopes; `spec_standard()` for the spec's CDISC standard.

Examples

```
# ---- Example 1: the whole study row, then one field ----
#
# spec_study() with no field returns the study-level table; pass a field
# name to pull a single value such as the study name.
spec <- artoo_spec(
  cdisc_adam_datasets, cdisc_adam_variables,
  codelists = cdisc_codelists,
  study = data.frame(study_name = "CDISCPIL0T01")
)
spec_study(spec)
spec_study(spec, "study_name")
```

spec_variables	<i>Variables in a spec</i>
----------------	----------------------------

Description

Return the variable-metadata table for one dataset, or for the whole spec. Each row carries the variable's CDISC data_type, label, length, display format, key sequence, and codelist reference.

Usage

```
spec_variables(spec, dataset = NULL)
```

Arguments

spec	<i>The specification to read.</i> <artoo_spec>: required.
dataset	<i>Restrict to one dataset.</i> <character(1)> NULL. When NULL (default) every dataset's variables are returned; otherwise only the named dataset's rows. Restriction: a non-NULL dataset must name a dataset in the spec (see spec_datasets()); an unknown name aborts with <code>artoo_error_input</code> .

Value

A data frame of variable metadata, one row per variable, with 22 columns (absent ones are filled with typed NA at construction):

- dataset, variable — the identifying pair (unique within a spec).
- itemoid — the Define-XML / Dataset-JSON item OID, when recorded.
- label — the variable label (<= 40 bytes for XPORT v5).
- data_type — canonical CDISC dataType (string, integer, decimal, float, double, boolean, date, datetime, time, URI).
- target_data_type — integer/decimal when a temporal variable stores as a SAS-epoch numeric; NA means ISO 8601 text (--DTC).
- length — declared storage length (bytes for character).
- display_format, informat — SAS format / informat strings.
- key_sequence — 1-based position in the dataset sort key.
- order — column position in the dataset.
- codelist_id, method_id, comment_id — references into the codelists / methods / comments slots.
- mandatory — logical obligation flag (NA is treated as mandatory by [check_spec\(\)](#)).
- significant_digits — for decimal variables.
- origin, source, predecessor, assigned_value, pages, role — Define-XML provenance fields, carried as-is.

Filter or arrange it with ordinary base / dplyr verbs.

See Also

[spec_datasets\(\)](#) for the dataset names; [spec_codelists\(\)](#) for a variable's controlled terminology.

Examples

```
spec <- artoo_spec(cdisc_sdtm_datasets, cdisc_sdtm_variables, codelists = cdisc_codelists)

# ---- Example 1: one dataset's variables ----
#
# Pass a dataset name to get just that domain's variables, already
# canonicalised to CDISC dataTypes.
head(spec_variables(spec, "DM")[, c("variable", "label", "data_type")])

# ---- Example 2: every variable across the spec ----
#
# Omit `dataset` to get the full table, e.g. to count variables per domain.
table(spec_variables(spec)$dataset)
```

sync_meta

Re-align metadata with a transformed data frame

Description

Re-attach and reconcile a `artoo_meta` after a transformation that dropped or reshaped it: the metadata's columns are narrowed and reordered to the frame's current columns, the record count is refreshed, the keys are recomputed, and a column the metadata does not describe gets an entry synthesized from its class and attributes. The one-liner to run after a `dplyr` (or `base`) pipeline, before handing the frame to a `write_*()` codec.

Usage

```
sync_meta(x, meta = NULL)
```

Arguments

`x` *The transformed data frame.* `<data.frame>`: required.

`meta` *The metadata to reconcile against.* `<artoo_meta>` | `NULL`. `NULL` (default) uses the frame's own `metadata_json` attribute.

Requirement: when the transform dropped the attribute (`base` [`subsetting` does), capture `get_meta()` before the pipeline and pass it here; a bare frame with no meta aborts with `artoo_error_input`.

Details

Why it exists. Base row subsetting (`x[i,]`) drops the frame's `metadata_json` attribute, and many tidyverse verbs rebuild the frame. `sync_meta()` takes the last-known metadata (the frame's own attribute when it survived, or an explicit meta) and makes it agree with the data again, so the round trip stays lossless without hand-editing.

Value

A `data.frame`: `x` re-stamped with the reconciled `artoo_meta`. Hand it to any `write_*`(`)` codec.

See Also

Read / attach: `get_meta()`, `set_meta()`.

Produce conformed frames: `apply_spec()`.

Examples

```
spec <- artoo_spec(cdisc_adam_datasets, cdisc_adam_variables, codelists = cdisc_codelists)

# ---- Example 1: re-attach after an attribute-dropping subset ----
#
# Base subsetting drops the metadata; capture it first, transform, then
# sync. The metadata narrows to the kept columns and the new row count.
adsl <- apply_spec(cdisc_adsl, spec, "ADSL", conformance = "off")
meta <- get_meta(adsl)
elderly <- adsl[adsl$AGE > 65, c("STUDYID", "USUBJID", "AGE")]
synced <- sync_meta(elderly, meta)
get_meta(synced)@dataset$records

# ---- Example 2: a derived column gains a synthesized entry ----
#
# A new column the metadata does not describe is profiled from its class,
# so the frame still writes losslessly.
adsl$AGEGR9 <- ifelse(adsl$AGE > 65, ">65", "<=65")
synced2 <- sync_meta(adsl)
get_meta(synced2)@columns$AGEGR9$dataType
```

 validate_spec

 Validate a specification for submission-readiness

Description

Run artoo's bundled, self-contained checks over a `artoo_spec`, **scoped to the dataset(s) you are working on**, and return a `artoo_check` that prints a sectioned report. Every finding is keyed to an open rule in the shipped catalog (see `spec_rules.json`); the result object keeps the findings as a plain data frame in `@findings` for programmatic use.

Usage

```
validate_spec(
  spec,
  data = NULL,
  dataset = NULL,
  on_error = c("off", "warn", "abort")
)
```

Arguments

spec	<i>The specification to validate.</i> <artoo_spec>: required.
data	<i>Optional input data for controlled-terminology checks.</i> <data.frame> named list NULL. When supplied, data values are cross-checked against the spec codelists. A single data frame requires a length-1 dataset; pass a named list to validate several at once.
dataset	<i>Restrict to one or more datasets.</i> <character> NULL. NULL (default) validates every dataset. Restriction: each name must be a dataset in the spec.
on_error	<i>What to do with an error-severity finding.</i> <character(1)>. One of: <ul style="list-style-type: none"> • "off" (default) collect and return every finding; never signal. • "warn" additionally cli_warn(artoo_warning_validation) with the error count. • "abort" additionally abort with artoo_error_validation. All findings are collected and returned in every case.

Details

Dataset-scoped. A spec workbook carries many datasets. Pass dataset to validate only the one(s) you are working on — the methods, comments, and codelists those datasets reference are checked for completeness, but unrelated datasets are not. dataset = NULL validates the whole spec.

Collect, do not stop. Every finding is collected and returned; validate_spec() does not abort on an error-severity finding unless on_error = "abort".

Value

A artoo_check *object*. Its @findings data frame has columns check, dimension, severity, dataset, variable, message. Print it for the sectioned report.

See Also

[artoo_spec\(\)](#) to build a spec; [spec_methods\(\)](#) / [spec_comments\(\)](#) for the metadata checked.

Examples

```
# ---- Example 1: validate one dataset ----
#
# Build a spec from the bundled ADaM tables and validate it; the
```

```

# result prints a sectioned report and keeps the findings table.
spec <- artoo_spec(
  cdisc_adam_datasets, cdisc_adam_variables,
  codelists = cdisc_codelists
)
chk <- validate_spec(spec, dataset = "ADSL")
chk@findings

# ---- Example 2: gate on errors with on_error = "abort" ----
#
# Point a key at a missing variable, then validate with on_error = "abort"
# and catch the resulting error.
bad_ds <- cdisc_sdtm_datasets
bad_ds$keys <- "NOTAVAR"
bad <- artoo_spec(bad_ds, cdisc_sdtm_variables, codelists = cdisc_codelists)
tryCatch(
  validate_spec(bad, dataset = "DM", on_error = "abort"),
  artoo_error_validation = function(e) conditionMessage(e)[1]
)

```

write_dataset

Write a dataset to any supported format

Description

Serialize a data frame to a clinical file format, preserving its `artoo_meta` losslessly. The codec is chosen from the file extension (or an explicit format), so one call covers `xpt`, Dataset-JSON, Parquet, and `rds`. This is the emit end of the `artoo` workflow; the per-format wrappers like `write_rds()` are thin sugar over it.

Usage

```
write_dataset(x, path, format = NULL, ...)
```

Arguments

<code>x</code>	<i>The dataset to write.</i> <data.frame>: required. Typically the output of <code>apply_spec()</code> , carrying <code>artoo_meta</code> .
<code>path</code>	<i>Destination file path.</i> <character(1)>: required. Its extension selects the codec unless format is given.
<code>format</code>	<i>Force a codec instead of inferring from the extension.</i> <character(1)> NULL. One of the registered formats (see <code>artoo_formats()</code>).
<code>...</code>	<i>Codec-specific arguments</i> passed through to the encoder (see the per-format wrappers, e.g. <code>write_xpt()</code> , for what each codec accepts). An argument the codec does not know is an error, never silently ignored.

Value

The input `x`, invisibly, so a write can sit mid-pipeline. Called for the side effect of writing path.

See Also

[read_dataset\(\)](#) for the inverse; [write_rds\(\)](#) for the per-format wrapper; [artoo_formats\(\)](#) for what is available.

Examples

```
spec <- artoo_spec(cdisc_adam_datasets, cdisc_adam_variables, codelists = cdisc_codelists)

# ---- Example 1: write a conformed dataset, inferring rds from the path ----
#
# apply_spec() attaches the metadata; write_dataset() carries it into the
# file so a later read is lossless.
adsl <- apply_spec(cdisc_adsl, spec, "ADSL", conformance = "off")
path <- tempfile(fileext = ".rds")
write_dataset(adsl, path)

# ---- Example 2: force the format for an unconventional extension ----
#
# When the extension does not name the format, pass it explicitly.
alt <- tempfile(fileext = ".data")
write_dataset(adsl, alt, format = "rds")
```

write_json

Write a dataset to CDISC Dataset-JSON

Description

Serialize a data frame to a CDISC Dataset-JSON v1.1 (.json) file, Dataset-JSON being the native home of the `artoo_meta` shape: the file is the metadata block plus a flat rows array. The emit end of the artoo workflow (`spec -> apply_spec -> write_json`); a thin wrapper over [write_dataset\(\)](#) with `format = "json"`.

Usage

```
write_json(
  x,
  path,
  on_invalid = c("error", "replace", "ignore"),
  created = NULL,
  strict = FALSE
)
```

Arguments

x	<i>The dataset to write.</i> <data.frame>: required. Typically the output of apply_spec() , carrying artoo_meta.
path	<i>Destination .json path.</i> <character(1)>: required.
on_invalid	<i>Policy for values that are not valid UTF-8.</i> <character(1)>: default "error". One of "error" (abort with artoo_error_codec, naming the offenders with their invalid bytes hex-escaped), "replace" (substitute ? and warn with artoo_warning_encoding), or "ignore" (drop the invalid bytes). The same policy vocabulary as write_xpt() ; text correctly read through artoo is always valid UTF-8, so this only fires on bytes that entered the frame through a mis-declared source encoding.
created	<i>Creation timestamp.</i> <POSIXct(1)> NULL. NULL (default) stamps the current time into datasetJSONCreationDateTime; freeze it for byte-stable output.
strict	<i>Suppress the _artoo extension block.</i> <logical(1)>: default FALSE. By default the file carries a single namespaced _artoo object when (and only when) there is content strict CDISC cannot express: SAS special-missing tags (.A-.Z, ._), the recorded source encoding, and informats. Data values stay plain nulls either way, so a foreign reader degrades gracefully. Note: strict = TRUE writes a pure closed-vocabulary file and warns (artoo_warning_codec) naming exactly what was dropped; those attributes will not survive a read-back.

Details

Full metadata, no loss. Unlike .xpt, a .json file records the complete artoo_meta: keySequence, codelist, origin, targetDataType, and significantDigits all survive. Dates, datetimes, and times are exchanged as ISO 8601 strings, or as SAS-epoch numbers when their targetDataType is "integer" (the ADaM numeric-date convention); decimal rides as a string so exact precision is preserved. The file is always UTF-8 (RFC 8259 / CDISC v1.1). NaN and infinite values are not valid CDISC numerics and abort the write.

Streaming write, whole-file read. The writer streams the rows array in bounded slabs (a .json.gz path gzips the stream transparently), but [read_json\(\)](#) must parse the whole array at once. For multi-million-row datasets prefer the NDJSON variant ([write_ndjson\(\)](#) / [read_ndjson\(\)](#)), which bounds memory in both directions.

Value

The input x, invisibly, so a write can sit mid-pipeline.

See Also

[read_json\(\)](#) for the inverse; [write_dataset\(\)](#) for the generic dispatcher.

Examples

```
# ---- Example 1: write a conformed dataset as Dataset-JSON ----
#
# apply_spec() attaches the metadata; write_json() serializes the full
# itemGroup plus the data rows.
adsl <- apply_spec(cdisc_adsl, adam_spec, "ADSL", conformance = "off")
```

```

path <- tempfile(fileext = ".json")
write_json(adsl, path)

# ---- Example 2: a frozen timestamp for reproducible bytes ----
#
# Fixing `created` makes two writes byte-identical; the columns() pane on
# the written file shows the full metadata the file carries (DM is SDTM,
# so it conforms against the bundled sdtm_spec).
dm <- apply_spec(cdisc_dm, sdtm_spec, "DM", conformance = "off")
path2 <- tempfile(fileext = ".json")
write_json(dm, path2, created = as.POSIXct("2020-01-01", tz = "UTC"))
columns(path2)

```

write_ndjson

Write a dataset to CDISC Dataset-JSON NDJSON

Description

Serialize a data frame to the newline-delimited variant of CDISC Dataset-JSON v1.1 (.ndjson): line 1 carries the complete metadata block, every following line one row array. The streaming end of the artoo workflow (spec -> apply_spec -> write_ndjson) for datasets too large for the array-form .json file; a thin wrapper over [write_dataset\(\)](#) with format = "ndjson".

Usage

```

write_ndjson(
  x,
  path,
  on_invalid = c("error", "replace", "ignore"),
  created = NULL,
  strict = FALSE
)

```

Arguments

x	<i>The dataset to write.</i> <data.frame>: required. Typically the output of apply_spec() , carrying artoo_meta.
path	<i>Destination .ndjson path.</i> <character(1)>: required. A .ndjson.gz path writes gzip-compressed bytes.
on_invalid	<i>Policy for values that are not valid UTF-8.</i> <character(1)>: default "error". One of "error" (abort with artoo_error_codec), "replace" (substitute ? and warn with artoo_warning_encoding), or "ignore" (drop the invalid bytes). See write_json() for when this fires.
created	<i>Creation timestamp.</i> <POSIXct(1)> NULL. NULL (default) stamps the current time into datasetJSONCreationDateTime; freeze it for byte-stable output.
strict	<i>Suppress the _artoo extension block.</i> <logical(1)>: default FALSE. See write_json() : the same extension semantics apply to the metadata line.

Details

Bounded memory, both directions. The writer streams slabs of per-column JSON literals and `read_ndjson()` parses slab-sized line batches, so a multi-million-row dataset never materializes a whole rows array the way the `.json` codec must. A `.ndjson.gz` path gzips the stream transparently.

Value

The input `x`, invisibly, so a write can sit mid-pipeline.

See Also

`read_ndjson()` for the inverse; `write_json()` for the array-form file; `write_dataset()` for the generic dispatcher.

Examples

```
spec <- artoo_spec(cdisc_adam_datasets, cdisc_adam_variables, codelists = cdisc_codelists)

# ---- Example 1: write a conformed dataset as NDJSON ----
#
# apply_spec() attaches the metadata; write_ndjson() streams the metadata
# line and one row per line.
adsl <- apply_spec(cdisc_adsl, spec, "ADSL", conformance = "off")
path <- tempfile(fileext = ".ndjson")
write_ndjson(adsl, path)
readLines(path, n = 2)[2]

# ---- Example 2: gzip the stream via the file extension ----
#
# A .ndjson.gz path compresses transparently; read_ndjson() inflates it.
gz <- tempfile(fileext = ".ndjson.gz")
write_ndjson(adsl, gz)
nrow(read_ndjson(gz))
```

write_parquet

Write a dataset to Apache Parquet

Description

Serialize a data frame to an Apache Parquet (`.parquet`) file, storing the data natively while preserving the full `artoo_meta` as a CDISC-shaped sidecar in the file's key-value metadata. The emit end of the `artoo` workflow (`spec -> apply_spec -> write_parquet`); a thin wrapper over `write_dataset()` with `format = "parquet"`. Requires the lightweight `nanoparquet` package.

Usage

```
write_parquet(
  x,
  path,
  encoding = NULL,
  on_invalid = c("error", "replace", "ignore"),
  compression = "snappy"
)
```

Arguments

x	<i>The dataset to write.</i> <data.frame>: required. Typically the output of apply_spec() , carrying <code>artoo_meta</code> .
path	<i>Destination .parquet path.</i> <character(1)>: required.
encoding	<i>Source charset to record.</i> <character(1)> NULL. The parquet bytes are always written as UTF-8 (the format's STRING type is UTF-8 by spec); encoding only records the data's original charset in the <code>artoo_meta</code> , so a later write_xpt() can reproduce the source bytes. NULL (default) leaves the recorded encoding untouched. Tip: any SAS or IANA spelling listed by artoo_encodings() is accepted.
on_invalid	<i>Policy for values that are not valid UTF-8.</i> <character(1)>: default "error". One of "error" (abort with <code>artoo_error_codec</code>), "replace" (substitute ? and warn with <code>artoo_warning_encoding</code>), or "ignore" (drop the invalid bytes). See write_json() for when this fires; parquet STRING bytes are UTF-8 by spec, exactly like Dataset-JSON.
compression	<i>Column compression codec.</i> <character(1)>: default "snappy". One of: <ul style="list-style-type: none"> "snappy" (default) — fast, the parquet ecosystem default. "gzip" — smaller files, slower. "zstd" — the best size/speed trade-off where supported. "uncompressed" — raw pages.

Details

Metadata where plain Parquet has none. A bare nanoparquet/arrow file drops labels, formats, and codelists; `write_parquet()` embeds the complete `artoo_meta` as a single Dataset-JSON-shaped string under the `metadata_json` key, so [read_parquet\(\)](#) restores every CDISC attribute. The same string is what a `.json` file or an `rds` carries, so conversion between any two formats stays lossless. A reader without `artoo` still opens the data and can see the `metadata_json` block.

Value

The input x, invisibly, so a write can sit mid-pipeline.

See Also

[read_parquet\(\)](#) for the inverse; [write_dataset\(\)](#) for the generic dispatcher.

Examples

```
spec <- artoo_spec(cdisc_adam_datasets, cdisc_adam_variables, codelists = cdisc_codelists)

# ---- Example 1: write a conformed dataset to Parquet ----
#
# apply_spec() attaches the metadata; write_parquet() stores the data
# natively and the metadata as a CDISC-shaped sidecar.
adsl <- apply_spec(cdisc_adsl, spec, "ADSL", conformance = "off")
path <- tempfile(fileext = ".parquet")
write_parquet(adsl, path)

# ---- Example 2: round-trip and confirm the metadata survived ----
#
# Reading it back yields an identical artoo_meta.
back <- read_parquet(path)
identical(get_meta(back)@columns, get_meta(adsl)@columns)
```

write_rds

Write a dataset to rds

Description

Write a data frame to an R `.rds` file, preserving its `artoo_meta`. A thin wrapper over `write_dataset()` with `format = "rds"`; the `rds` carries the metadata both as live R attributes and as the language-agnostic `metadata_json` string, so `read_rds()` restores it exactly.

Usage

```
write_rds(x, path, encoding = NULL)
```

Arguments

<code>x</code>	<i>The dataset to write.</i> <code><data.frame></code> : required.
<code>path</code>	<i>Destination .rds path.</i> <code><character(1)></code> : required.
<code>encoding</code>	<i>Source charset to record.</i> <code><character(1)></code> <code>NULL</code> . <code>rds</code> is R-native and faithful: strings are saved as-is, never transcoded. <code>encoding</code> only records the data's original charset in the <code>artoo_meta</code> , so a later <code>write_xpt()</code> can reproduce the source bytes. <code>NULL</code> (default) leaves the recorded encoding untouched. Tip: any SAS or IANA spelling listed by <code>artoo_encodings()</code> is accepted.

Value

The input x, invisibly, so a write can sit mid-pipeline.

See Also

`read_rds()` for the inverse; `write_dataset()` for the generic dispatcher.

Examples

```
spec <- artoo_spec(cdisc_adam_datasets, cdisc_adam_variables, codelists = cdisc_codelists)

# ---- Example 1: write a conformed dataset to rds ----
#
# apply_spec() attaches the metadata; write_rds() carries it into the file.
adsl <- apply_spec(cdisc_adsl, spec, "ADSL", conformance = "off")
path <- tempfile(fileext = ".rds")
write_rds(adsl, path)

# ---- Example 2: round-trip and confirm the metadata survived ----
#
# Reading it back yields an identical artoo_meta.
back <- read_rds(path)
identical(get_meta(back)@columns, get_meta(adsl)@columns)
```

write_spec

Write a specification to native JSON or a P21 Excel workbook

Description

Serialise a `artoo_spec`, dispatching on the file extension: a `.json` path writes artoo's native, lossless JSON; a `.xlsx` path writes a Pinnacle 21 (P21) style Excel workbook. Both are inverses of `read_spec()` on their format, which makes the spec converters free compositions: `read_spec("define.xml") |> write_spec("spec.xlsx")` is a Define-XML to P21 bridge in one line.

Usage

```
write_spec(spec, path)
```

Arguments

spec	<i>The specification to serialise.</i> <code><artoo_spec></code> : required. Build one with <code>artoo_spec()</code> or <code>read_spec()</code> .
path	<i>Destination file.</i> <code><character(1)></code> : required. The extension picks the format: <code>.json</code> (native, lossless) or <code>.xlsx</code> (P21 interchange; needs the <code>writexl</code> package). Any other extension aborts with <code>artoo_error_input</code> .

Details

Native JSON is the lossless format. Each slot is written as an array of row objects, with NA encoded as JSON null and numbers at full precision, so `read_spec()` rebuilds an identical `artoo_spec` through `artoo_spec()`. Object keys are emitted in a fixed order, so writing the same spec twice yields byte-identical output.

P21 xlsx is the interchange format. Sheets are emitted with the headers the P21 reader recognises (Define, Datasets, Variables, ValueLevel, Codelists, Methods, Comments, Documents; empty

optional sheets are omitted), foreign keys repeated on every row (no merged cells), and the spec's `spec_standard()` as the Datasets sheet's Standard column. The study row writes back as the Define sheet's Attribute/Value pairs (StudyName, StudyDescription, ProtocolName). The Data Type column is written in the Define-XML / ODM vocabulary the workbook expects: a character variable is text (not the Dataset-JSON string), and decimal / double collapse to float, boolean / URI to text.

Columns the P21 vocabulary does not model are not lost: a foreign column carried on a slot is re-emitted verbatim under its own header, so an xlsx round-trip keeps user columns.

Note: fields with no P21 column (itemoid, target_data_type, per-variable key_sequence) do not survive an xlsx round-trip; persist to JSON when you need the spec back exactly. The Data Type re-encoding is also non-injective: decimal, double, boolean, and URI fold to float or text on a read-back. A Define-XML partialDate / partialDatetime (and the other partial / incomplete subtypes) is read as the base date / datetime – CDISC Dataset-JSON v1.1 has no partial dataType – so it is written back as the base type.

Value

The output path, invisibly. Read it back with `read_spec()`.

See Also

Inverse: `read_spec()` reads native JSON, a P21 Excel workbook, or Define-XML back into a `artoo_spec`.

Build / inspect: `artoo_spec()`, `spec_datasets()`, `spec_variables()`, `spec_standard()`.

Examples

```
# ---- Example 1: persist a spec to JSON, then read it back ----
#
# Build a spec from the bundled CDISC-pilot tables, write it to a temp
# JSON file, and confirm read_spec() reconstructs it intact.
spec <- artoo_spec(
  cdisc_adam_datasets, cdisc_adam_variables,
  codelists = cdisc_codelists
)
path <- tempfile(fileext = ".json")
write_spec(spec, path)
identical(read_spec(path), spec)

# ---- Example 2: the same spec as a P21 workbook ----
#
# The .xlsx path emits P21-shaped sheets; reading the workbook back
# recovers the P21-representable surface (here: the dataset names).
if (requireNamespace("writexl", quietly = TRUE)) {
  xlsx <- tempfile(fileext = ".xlsx")
  write_spec(spec, xlsx)
  spec_datasets(read_spec(xlsx))
}
```

write_xpt

*Write a dataset to SAS XPORT***Description**

Serialize a data frame to a SAS Transport (.xpt) file in v5 (the FDA submission standard) or v8 (extended names and labels), preserving the `artoo_meta` a column can hold. The emit end of the artoo workflow (`spec -> apply_spec -> write_xpt`); a thin wrapper over `write_dataset()` with `format = "xpt"`.

Usage

```
write_xpt(
  x,
  path,
  version = 5,
  encoding = NULL,
  on_invalid = c("error", "replace", "ignore"),
  created = NULL
)
```

Arguments

<code>x</code>	<i>The dataset to write.</i> <data.frame>: required. Typically the output of <code>apply_spec()</code> , carrying <code>artoo_meta</code> .
<code>path</code>	<i>Destination .xpt path.</i> <character(1)>: required.
<code>version</code>	<i>XPORT transport version.</i> <integer(1)>: default 5. 5 (the FDA standard: names <= 8 characters, labels <= 40 bytes) or 8 (names <= 32, long labels).
<code>encoding</code>	<i>Target charset.</i> <character(1)> NULL. NULL (default) inherits the source encoding recorded in <code>artoo_meta</code> , else UTF-8. IANA and SAS names ("US-ASCII", "wlatin1") both work. Tip: any SAS or IANA spelling listed by <code>artoo_encodings()</code> is accepted.
<code>on_invalid</code>	<i>Policy for values not representable in encoding.</i> <character(1)>: default "error". One of "error" (abort with <code>artoo_error_codec</code> , naming the offenders), "replace" (substitute ? and warn with <code>artoo_warning_encoding</code>), or "ignore" (drop them). The same policy vocabulary as the UTF-8 writers (<code>write_json()</code> , <code>write_ndjson()</code> , <code>write_parquet()</code>).
<code>created</code>	<i>Header timestamp.</i> <POSIXct(1)> NULL. NULL (default) stamps the current time; freeze it for byte-stable output.

Details

What XPORT can carry. An .xpt file's NAMESTR stores only variable name, label, length, and SAS format. CDISC metadata beyond that (`keySequence`, `codelist`, `origin`, `targetDataType`, ...) and the source encoding are not representable in the bytes; they ride the in-session `artoo_meta` and the

sidecar in self-describing formats (Dataset-JSON, Parquet, rds). XPORT also cannot distinguish an empty string from NA (both store as blanks) and drops trailing spaces.

Character ISO dates (-DTC) write as text. A character column whose `dataType` is `date/datetime/time` with no numeric `targetDataType` is the CDISC ISO 8601 text form — the SDTM --DTC convention — and stores as a character variable, partial dates ("1951", "1951-12") included, byte for byte. The SAS-numeric encoding (with DATE9.-style formats) is used for columns that are R Date/POSIXct/hms or whose metadata records `targetDataType = "integer"` (the ADaM numeric-date convention). A character column *under* `targetDataType = "integer"` aborts loudly — a partial date can never become a SAS numeric silently.

Value

The input x, invisibly, so a write can sit mid-pipeline.

See Also

[read_xpt\(\)](#) for the inverse; [write_dataset\(\)](#) for the generic dispatcher.

Examples

```
spec <- artoo_spec(
  cdisc_adam_datasets, cdisc_adam_variables,
  codelists = cdisc_codelists
)

# ---- Example 1: write a conformed dataset as v5 (FDA standard) ----
#
# apply_spec() attaches the metadata; write_xpt() carries the label, length,
# and SAS format for each variable into the transport file.
adsl <- apply_spec(cdisc_adsl, spec, "ADSL", conformance = "off")
path <- tempfile(fileext = ".xpt")
write_xpt(adsl, path)

# ---- Example 2: v8 for long names, with a frozen timestamp ----
#
# Version 8 keeps names over 8 characters; a fixed `created` makes the bytes
# reproducible. Reading it back shows the labels, types, and record count
# survived the transport. DM is SDTM, so it conforms against the bundled
# sdtm_spec.
dm <- apply_spec(cdisc_dm, sdtm_spec, "DM", conformance = "off")
path8 <- tempfile(fileext = ".xpt")
write_xpt(dm, path8, version = 8, created = as.POSIXct("2020-01-01", tz = "UTC"))
get_meta(read_xpt(path8))@dataset$records
```

xpt_members

List the members of a SAS XPORT transport file

Description

Report every dataset (member) a SAS Transport (.xpt) file holds, with its label, variable count, and row count — the survey step before `read_xpt()` with `member =` picks one. A single-member file (the FDA submission convention) returns one row.

Usage

```
xpt_members(path)
```

Arguments

`path` *Source .xpt path.* <character(1)>: required. A file that is not a valid XPORT library aborts with `artoo_error_codec`.

Details

v5 has no recorded row count. A v8 member records its rows; a v5 member's count is derived from the byte span up to the next member (or end of file) minus trailing padding, so an all-character v5 member whose last row is entirely blank reports one row fewer (the documented v5 ambiguity, see `write_xpt()`).

Value

A <data.frame> with one row per member and columns `member` (1-based index), `name`, `label`, `nvars`, and `nobs`. Pass `member` or `name` to `read_xpt()`.

See Also

`read_xpt()` with `member =` to read one of them.

Examples

```
spec <- artoo_spec(
  cdisc_adam_datasets, cdisc_adam_variables,
  codelists = cdisc_codelists
)

# ---- Example 1: a single-member file reports one row ----
#
# The FDA convention is one dataset per transport file.
dm <- apply_spec(cdisc_dm, sdtm_spec, "DM", conformance = "off")
p <- tempfile(fileext = ".xpt")
write_xpt(dm, p)
xpt_members(p)
```

```
# ---- Example 2: survey a multi-member file, then read one member ----  
#  
# Concatenate two single-member files into one library and list it.  
adsl <- apply_spec(cdisc_adsl, spec, "ADSL", conformance = "off")  
p2 <- tempfile(fileext = ".xpt")  
write_xpt(adsl, p2)  
multi <- tempfile(fileext = ".xpt")  
writeBin(  
  c(  
    readBin(p, "raw", file.size(p)),  
    readBin(p2, "raw", file.size(p2))[-(1:240)]  
  ),  
  multi  
)  
xpt_members(multi)
```

Index

* datasets

- cdisc_adae, 13
 - cdisc_adsl, 13
 - cdisc_dm, 14
 - cdisc_spec, 14
 - cdisc_specs, 15
 - cdisc_suppdm, 16
 - cdisc_ts, 17
 - cdisc_vs, 17
- adam_spec (cdisc_specs), 15
- apply_spec, 3
- apply_spec(), 8, 18, 19, 21–27, 29, 43–45, 55, 57, 59, 60, 62, 66
- artoo_checks, 6
- artoo_checks(), 18–20, 23, 27, 28
- artoo_encodings, 8
- artoo_encodings(), 33, 36, 37, 40, 62, 63, 66
- artoo_formats, 10
- artoo_formats(), 9, 31, 57, 58
- artoo_spec, 10
- artoo_spec(), 14, 16, 29, 38, 39, 51, 52, 56, 64, 65
- cdisc_adae, 13
- cdisc_adam_datasets (cdisc_spec), 14
- cdisc_adam_variables (cdisc_spec), 14
- cdisc_adsl, 13
- cdisc_codelists (cdisc_spec), 14
- cdisc_dm, 14
- cdisc_sdtm_datasets (cdisc_spec), 14
- cdisc_sdtm_variables (cdisc_spec), 14
- cdisc_spec, 14
- cdisc_specs, 15
- cdisc_suppdm, 16
- cdisc_ts, 17
- cdisc_vs, 17
- check_spec, 18
- check_spec(), 4–6, 8, 19, 20, 23, 25, 27, 41, 42, 45, 46, 53
- check_study, 19
- check_study(), 19, 41, 42
- columns, 21
- columns(), 31
- conformance, 22
- conformance(), 4, 5
- decode_column, 24
- decode_column(), 5
- get_meta, 26
- get_meta(), 5, 22, 28, 32–34, 36, 40, 43, 55
- is_artoo_checks, 27
- is_artoo_meta, 28
- is_artoo_meta(), 29
- is_artoo_spec, 29
- is_artoo_spec(), 12
- members, 30
- read_dataset, 31
- read_dataset(), 10, 22, 30, 32–37, 39, 40, 58
- read_json, 32
- read_json(), 9, 34, 59
- read_ndjson, 34
- read_ndjson(), 59, 61
- read_parquet, 35
- read_parquet(), 62
- read_rds, 36
- read_rds(), 31, 32, 63
- read_spec, 38
- read_spec(), 64, 65
- read_xpt, 39
- read_xpt(), 9, 32, 67, 68
- repair_spec, 41
- repair_spec(), 4, 5, 19, 20, 45
- sdtm_spec (cdisc_specs), 15
- set_meta, 42
- set_meta(), 5, 26–28, 55

set_type, 44
set_type(), 4, 5, 42
spec_codelists, 45
spec_codelists(), 12, 25, 54
spec_comments, 46
spec_comments(), 48, 50, 56
spec_datasets, 47
spec_datasets(), 12, 16, 18, 39, 49, 52–54,
65
spec_documents, 48
spec_documents(), 47, 50
spec_keys, 49
spec_keys(), 12, 48
spec_methods, 50
spec_methods(), 47, 48, 56
spec_standard, 51
spec_standard(), 15, 16, 52, 65
spec_study, 52
spec_study(), 12, 51
spec_variables, 53
spec_variables(), 12, 16, 39, 45, 46, 48, 49,
65
sync_meta, 54

validate_spec, 55
validate_spec(), 11, 12, 18–20, 39, 44,
46–48, 50

write_dataset, 57
write_dataset(), 10, 31, 32, 58–63, 66, 67
write_json, 58
write_json(), 33, 60–62, 66
write_ndjson, 60
write_ndjson(), 34, 59, 66
write_parquet, 61
write_parquet(), 36, 66
write_rds, 63
write_rds(), 36, 37, 57, 58
write_spec, 64
write_spec(), 15, 38, 39, 41, 42, 44, 45
write_xpt, 66
write_xpt(), 9, 40, 57, 59, 62, 63, 68

xpt_members, 68
xpt_members(), 22, 30, 31, 40