

# Package: ksCompare (via r-universe)

July 6, 2026

**Title** Smart Data Frame Comparison in the Spirit of SAS PROC COMPARE

**Version** 0.2.3

**Description** A tidyverse-native engine for comparing two data frames in the spirit of SAS 'PROC COMPARE', with extensions for clinical/pharma workflows: per-column tolerances (absolute, relative, and ULP), smart column matching, intelligent type reconciliation, SAS special-missing awareness, diff-pattern detection, and rich HTML/Excel reports.

**License** MIT + file LICENSE

**Encoding** UTF-8

**Roxygen** list(markdown = TRUE)

**RoxygenNote** 8.0.0

**Depends** R ( $\geq$  4.1)

**Imports** base64enc, cli, digest, dplyr, rlang, stringi, tibble, vctrs

**Suggests** arrow, arsenal, broom, clue, haven, htmltools, knitr, openxlsx2, pointblank, reactable, rmarkdown, stringdist, testthat ( $\geq$  3.0.0), withr

**Config/testthat/edition** 3

**URL** <https://al-garik.github.io/ksCompare/>,  
<https://github.com/al-garik/ksCompare>,  
<https://crow16384.r-universe.dev/ksCompare>

**BugReports** <https://github.com/al-garik/ksCompare/issues>

**VignetteBuilder** knitr

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

**Repository** <https://crow16384.r-universe.dev>

**Date/Publication** 2026-07-06 13:43:16 UTC

**RemoteUrl** <https://github.com/al-garik/ksCompare>

**RemoteRef** HEAD

**RemoteSha** 554841cfa9485bd78809a67cac190b15b6504d09

## Contents

as_ks_comparison . . . . .	2
ks_assert_clean . . . . .	3
ks_cause_summary . . . . .	4
ks_comp_options . . . . .	5
ks_compare . . . . .	7
ks_glance . . . . .	11
ks_match_health . . . . .	12
ks_recommendations . . . . .	13
ks_report_html . . . . .	14
ks_report_xlsx . . . . .	15
ks_row_diff_summary . . . . .	16
ks_set_comp_options . . . . .	17
ks_suggest_key . . . . .	18
ks_sysinfo . . . . .	19
ks_tidy . . . . .	20
ks_tol . . . . .	21
ks_unmatched_rows . . . . .	22
sas-out . . . . .	23
summary.ks_comparison . . . . .	24
<b>Index</b>	<b>25</b>

---

as_ks_comparison	<i>Convert an arsenal::comparedf result into a ks_comparison</i>
------------------	--

---

### Description

Lossy interop: takes the two source frames stored on a `arsenal::comparedf` object and re-runs `ks_compare()` using the same BY keys. Useful for teams migrating from `arsenal::comparedf()` to `ksCompare` without rewriting their inputs.

### Usage

```
as_ks_comparison(x, ...)
```

### Arguments

x	An <code>arsenal::comparedf</code> object.
...	Additional arguments forwarded to <code>ks_compare()</code> (e.g. <code>tolerance</code> , <code>options</code> ).

### Details

Note that `ksCompare`'s diff is computed from scratch — it does not attempt to translate `comparedf`'s internal diff tables row-by-row.

**Value**

A `ks_comparison` object.

---

<code>ks_assert_clean</code>	<i>Use a comparison as a pipeline gate</i>
------------------------------	--

---

**Description**

Treats a `ks_comparison` as a pass/fail check. By default a "passing" comparison has zero schema, key, and value differences; callers can relax those expectations to allow a budget of known acceptable diffs.

**Usage**

```
ks_assert_clean(
  x,
  max_value_diffs = 0L,
  max_schema_diffs = 0L,
  max_unmatched_rows = 0L
)

ks_pointblank_step(
  agent,
  comparison,
  max_value_diffs = 0L,
  max_schema_diffs = 0L,
  max_unmatched_rows = 0L,
  label = "ksCompare clean"
)
```

**Arguments**

<code>x</code>	A <code>ks_comparison</code> object returned by <code>ks_compare()</code> .
<code>max_value_diffs</code>	Integer. Maximum number of differing cells allowed in matched rows. Default 0L (strict).
<code>max_schema_diffs</code>	Integer. Maximum schema differences allowed, counted as the sum of base-only columns, comp-only columns, and matched columns whose <code>kind</code> differs. Default 0L (strict).
<code>max_unmatched_rows</code>	Integer. Maximum number of unmatched rows allowed ( <code>n_base_only_rows</code> + <code>n_comp_only_rows</code> ). Default 0L (strict).
<code>agent</code>	A <code>ptblank_agent</code> (typically piped in from <code>pointblank::create_agent()</code> ) or a data frame / tibble. The comparison gate is added as a step on this object and <code>agent</code> is returned, so it composes with the rest of a pointblank pipeline.

comparison      A `ks_comparison` produced by `ks_compare()` that the step will gate on.  
 label            Optional label for the pointblank step.

## Details

Two flavours are provided:

- `ks_assert_clean()` throws a classed error (`ksCompare_assertion_failed`) and halts a pipeline when expectations are unmet. Returns the comparison invisibly on success, so it can be dropped into a `|>` chain.
- `ks_pointblank_step()` returns a pointblank step that wraps `ks_assert_clean()` for use inside `pointblank::create_agent()` / `pointblank::action_levels()` flows.

## Value

`ks_assert_clean()` returns `x` invisibly when expectations are met. On failure, raises a condition of class `ksCompare_assertion_failed` (catchable in CI with `tryCatch(..., ksCompare_assertion_failed = function(e) ...)`). `ks_pointblank_step()` returns a pointblank step.

## Examples

```
a <- data.frame(id = 1:3, x = 1:3)
ks_compare(a, a, by = "id") |> ks_assert_clean()

# Allow a small budget of known diffs
## Not run:
ks_compare(a, b, by = "id") |>
  ks_assert_clean(max_value_diffs = 5)

## End(Not run)
```

---

`ks_cause_summary`      *Diff causes summary*

---

## Description

Aggregates the `note` column of `cmp$value_diff` into a tidy taxonomy: one row per distinct cause, with counts and the columns it affects. Helps answer "what kinds of discrepancies do we have?" before drilling into individual cells.

## Usage

```
ks_cause_summary(x)
```

## Arguments

`x`                    A `ks_comparison`.

## Details

Compound notes (multiple cues joined by "; ") are split into their elementary cues, so a single cell can contribute to several causes. Cells with `note = NA` (a "plain" value change with no recognised cue) are bucketed as "plain value change".

## Value

A tibble with columns:

- `cause`: short label (e.g. "letter case differs").
- `n_cells`: number of cells contributing this cue.
- `n_columns`: number of distinct base columns affected.
- `columns`: comma-separated sample of the affected columns (capped at 5 names). Sorted by `n_cells` descending. Zero-row tibble when there are no value differences.

## See Also

[ks\\_compare\(\)](#), [ks\\_row\\_diff\\_summary\(\)](#).

## Examples

```
a <- data.frame(id = 1:3, x = c("a", "b ", "c"), y = c(1, 2, 3))
b <- data.frame(id = 1:3, x = c("A", "b", "c"), y = c(1, 2, 4))
cmp <- ks_compare(a, b, by = "id")
ks_cause_summary(cmp)
```

---

<code>ks_comp_options</code>	<i>Comparison options</i>
------------------------------	---------------------------

---

## Description

Builds an options object consumed by [ks\\_compare\(\)](#) that controls missing-value semantics, label / format comparison, string normalisation, and time-zone handling.

## Usage

```
ks_comp_options(
  na_equal = getOption("ksCompare.na_equal", TRUE),
  sas_special_missing = getOption("ksCompare.sas_special_missing", TRUE),
  compare_labels = getOption("ksCompare.compare_labels", TRUE),
  compare_formats = getOption("ksCompare.compare_formats", TRUE),
  str_trim = getOption("ksCompare.str_trim", FALSE),
  str_case = getOption("ksCompare.str_case", "sensitive"),
  str_norm = getOption("ksCompare.str_norm", "none"),
  tz = getOption("ksCompare.tz", "preserve"),
  path = getOption("ksCompare.path", NULL)
)
```

**Arguments**

<code>na_equal</code>	Treat two NAs as equal? Default <code>TRUE</code> . With <code>FALSE</code> , any cell where one side is <code>NA</code> and the other is not is reported as a value diff; cells where <i>both</i> sides are <code>NA</code> are also reported (mirroring SAS PROC COMPARE behaviour with the <code>NOMISS</code> option turned off).
<code>sas_special_missing</code>	Distinguish SAS special missings ( <code>.A-.Z</code> and <code>._</code> ) when comparing numeric columns imported via <code>haven</code> ? Default <code>TRUE</code> . When <code>TRUE</code> , <code>.A</code> and <code>.B</code> compare as different even though both are <code>NA</code> in R; relies on the <code>tagged_na</code> tag attached by <code>haven::read_sas()</code> .
<code>compare_labels</code>	Compare <code>attr(x, "label")</code> between matched columns and surface differences in <code>cmp\$schema_diff</code> ? Default <code>TRUE</code> . Set to <code>FALSE</code> to suppress label-only diffs (e.g. when metadata drift is expected).
<code>compare_formats</code>	Compare <code>attr(x, "format.sas")</code> between matched columns? Default <code>TRUE</code> . Comparison is trailing-dot- and case-tolerant, so <code>DATE9.</code> , <code>DATE9</code> , and <code>date9.</code> all compare as equal.
<code>str_trim</code>	Trim leading/trailing whitespace with <code>stringi::stri_trim_both()</code> before comparing strings? Default <code>FALSE</code> . Useful when one source has been padded by a fixed-width exporter.
<code>str_case</code>	One of <code>"sensitive"</code> (default) or <code>"fold"</code> . When <code>"fold"</code> , strings are compared after Unicode case folding ( <code>stringi::stri_trans_tolower()</code> ).
<code>str_norm</code>	One of <code>"none"</code> (default) or <code>"NFC"</code> . When <code>"NFC"</code> , strings are Unicode-normalised before comparing so that visually-identical sequences with different code-point compositions match.
<code>tz</code>	One of <code>"preserve"</code> (default), <code>"UTC"</code> , or <code>"strip"</code> . Controls how POSIXct columns are reconciled when the two sides carry different <code>tzzone</code> attributes: <ul style="list-style-type: none"> <li>• <code>"preserve"</code>: values are compared as-is (a <code>tz</code> mismatch surfaces as a diff if it changes the underlying instant).</li> <li>• <code>"UTC"</code>: both sides are converted to UTC before compare.</li> <li>• <code>"strip"</code>: <code>tzzone</code> is dropped and values compared as POSIXct.</li> </ul>
<code>path</code>	Optional output folder used by downstream writers ( <code>ks_report_html()</code> , <code>ks_report_xlsx()</code> ) when their own <code>path =</code> argument is left blank. Default <code>NULL</code> means "use the current working directory". The folder is created on first use if it does not already exist. Setting it once on <code>ks_comp_options()</code> and passing the result through <code>ks_compare()</code> keeps every artefact from a single comparison run in the same place.

**Value**

A `ks_comp_options` S3 list.

**Global defaults via options()**

Each argument falls back to a global R option in the `ksCompare.*` namespace, then to the package default. Precedence is: explicit argument > `getOption("ksCompare.<arg>")`

> built-in default. That makes it possible to set project-wide defaults once and then call `ks_compare()` without repeating yourself:

```
# In .Rprofile or at the top of a script
options(ksCompare.path      = "out/qc",
        ksCompare.str_case = "fold")
```

```
ks_compare(base, comp) |> ks_report_html() # honours the globals
```

`ks_set_comp_options()` is a small wrapper around `base::options()` that takes the same argument names without the `ksCompare.` prefix.

Recognised options: `ksCompare.na_equal`, `ksCompare.sas_special_missing`, `ksCompare.compare_labels`, `ksCompare.compare_formats`, `ksCompare.str_trim`, `ksCompare.str_case`, `ksCompare.str_norm`, `ksCompare.tz`, `ksCompare.path`.

## See Also

`ks_set_comp_options()` for setting the `ksCompare.*` globals from R code.

## Examples

```
# Defaults: strict NAs, label & format comparison on
ks_comp_options()

# Loose strings: trim padding and ignore case
ks_comp_options(str_trim = TRUE, str_case = "fold")

# Suppress label/format drift, but keep cell-level strictness
ks_comp_options(compare_labels = FALSE, compare_formats = FALSE)

# Treat any NA as a difference (PROC COMPARE without NOMISS)
ks_comp_options(na_equal = FALSE)

# Pin every report from this run to a dedicated folder
ks_comp_options(path = tempfile("ksCompare_"))

# Pick up a project-wide global path
withr::with_options(
  list(ksCompare.path = tempfile("ksCompare_")),
  ks_comp_options()$path
)
```

---

ks\_compare

*Compare two data frames*

---

## Description

`ks_compare()` is the single entry point for the package. It compares `base` against `comp` and returns a `ks_comparison` S3 object containing schema, key, row, and value differences plus a small QC manifest.

## Usage

```
ks_compare(
  base,
  comp,
  by = NULL,
  mapping = NULL,
  tolerance = ks_tol(),
  coerce = c("safe", "strict", "lossy"),
  dup_keys = c("first", "last", "keep_all", "all_pairs", "error"),
  options = ks_comp_options(),
  base_name = NULL,
  comp_name = NULL,
  find_patterns = FALSE,
  n_first_last = 5L,
  max_unmatched_rows = 100L,
  loglevel = c("verdict", "verbose")
)
```

## Arguments

<code>base, comp</code>	<p>Either an in-memory data frame (or anything coercible via <code>tibble::as_tibble()</code>: tibbles, data.tables, Arrow tables) or a single character file path. When given a path, the file is read automatically based on its extension:</p> <ul style="list-style-type: none"> <li>• <code>.sas7bdat</code>, <code>.xpt</code> (via <code>haven</code>)</li> <li>• <code>.parquet</code>, <code>.feather</code> / <code>.arrow</code> (via <code>arrow</code>)</li> <li>• <code>.rds</code>, <code>.rdata</code> / <code>.rda</code></li> <li>• <code>.csv</code>, <code>.tsv</code> / <code>.txt</code></li> </ul> <p><code>base</code> is treated as the reference; <code>comp</code> is the candidate. When both inputs are file paths that share the same basename (e.g. <code>prod/adsl.sas7bdat</code> vs <code>qc/adsl.sas7bdat</code>), display names are automatically disambiguated with the parent directory.</p>
<code>by</code>	<p>Optional column name(s) to use as the matching key.</p> <ul style="list-style-type: none"> <li>• <code>NULL</code> (default): rows are matched by row position. A warning banner is shown in the HTML report so this is never silent.</li> <li>• <code>character</code>: same column name on both sides. Use a named vector when the key is renamed: <code>c(USUBJID = "SUBJID")</code>.</li> <li>• <code>"auto"</code>: search the smallest combination of shared columns that is unique on <b>both</b> sides (capped at 4 columns). Falls back to position match with a warning if nothing qualifies.</li> </ul>
<code>mapping</code>	<p>Optional explicit column mapping for renamed columns. Either a fully named character vector (<code>c(base_col = "comp_col", ...)</code>) or a 2-column data frame with columns <code>base</code> and <code>comp</code>. Mapping pairs always win over exact-name matching.</p>
<code>tolerance</code>	<p>A <code>ks_tol()</code> specification. Defaults to strict equality. Per-column overrides supported via <code>ks_tol(per_column = list(col = ks_tol(...)))</code>.</p>
<code>coerce</code>	<p>Type-coercion strictness. One of:</p>

- "safe" (default): `vctrs::vec_ptype2()` plus integer double, factor character, `haven_labelled` unwrapping.
  - "strict": only common types `vctrs::vec_ptype2()` agrees on. Type mismatches surface as `kind = "type_mismatch"` rows.
  - "lossy": additionally allows numeric character, date/datetime character (ISO-8601), factor integer (level codes).
- dup\_keys** Strategy for duplicate keys (per-side independently). One of:
- "first" (default): keep the first occurrence per side. Drops rows; emits a `ksCompare_dup_keys_resolved` info message.
  - "last": keep the last occurrence per side. Same message.
  - "keep\_all": pair duplicate rows positionally within each key group (1 1, 2 2, ...). Leftover rows on the longer side become base- or compare-only.
  - "all\_pairs": cartesian-pair every base row with every comp row sharing a key value. Emits a `ksCompare_all_pairs_cardinality` warning when group sizes disagree.
  - "error": raise on any duplicate key.
- options** A `ks_comp_options()` specification controlling NA / SAS special-missing semantics, label / format comparison, string normalisation, and the default output folder for reports.
- base\_name, comp\_name** Optional display names for the two frames (used in `print()` and reports). Default to the names of the supplied expressions.
- find\_patterns** Logical (default `FALSE`). When `TRUE`, run the pattern-detection pass (constant offsets, sign flips, unit rescales, epoch shifts, etc.) and populate `cmp$pattern_summary`. The detectors iterate over every column with a diff and can be noticeably slow on large value-diff tables, so they are off by default. When `FALSE`, `cmp$pattern_summary` is an empty tibble with the usual column layout.
- n\_first\_last** Non-negative integer (default 5). Per matched column with at least one cell diff, capture the first and last `n_first_last` differences in `key_id` order on `cmp$first_last_unequal`. Mirrors the *First / Last N Obs With Some Compared Variables Unequal* tables of SAS PROC COMPARE. Set to 0 to skip.
- max\_unmatched\_rows** Non-negative integer (default 100). Cap on how many full base-only / comp-only rows are stored on `cmp$unmatched_rows` for inspection (and surfaced by `ks_unmatched_rows()` / the HTML / XLSX reports). The cap is apportioned proportionally between the two sides; truncation is recorded on the result's `truncated` and `n_total` attributes. Set to 0 to skip storing the row data entirely (the row counts on `cmp$key_diff` and `summary()` are unaffected).
- loglevel** Controls what `print()` emits after the comparison. One of:
- "verdict" (default): prints a single coloured verdict line (`cli_alert_success` / `_warning` / `_danger`) plus any "critical" recommendations. The

full schema / rows / values breakdown is omitted. Pipeline info messages (duplicate-key resolution, auto-key selection, etc.) are unaffected.

- "verbose": `print()` emits the complete schema / rows / values breakdown. Pipeline info messages fire normally as well.

## Value

A `ks_comparison` object with components:

- `meta`: counts, keys, matching strategy, row-key lookup table.
- `schema_diff`: one row per matched / unmatched column with kind / label / format comparison.
- `key_diff`: counts of matched / base-only / comp-only rows.
- `row_diff`: long table of `key_id`, `base_row`, `comp_row`, `status`.
- `value_diff`: long table of differing cells with `column_base`, `column_comp`, `kind`, `base`, `comp`, `diff`, `note`.
- `pattern_summary`: detected recurring shapes per column (only populated when `find_patterns = TRUE`; otherwise an empty tibble with the standard column layout).
- `unmatched_rows`: full base-only / comp-only rows (capped by `max_unmatched_rows`); see `ks_unmatched_rows()`.
- `first_last_unequal`: per-column first / last `n_first_last` differing observations in `key_id` order (PROC COMPARE-style).
- `verdict`: a list with `headline` (character), `severity` ("ok", "info", "warn", or "critical"), `pct_match` (0–100 numeric), `n_cells`, and `n_diffs`; the same summary shown by `print()`.
- `options`, `tolerance`: the inputs (echoed for the manifest).
- `manifest`: input hashes + run metadata.

Inspect with `print()`, `summary()`, `tibble::as_tibble()`, or render via `ks_report_html()` / `ks_report_xlsx()`.

## Pipeline

1. **Schema alignment** — columns are paired by an explicit `mapping` first, then by exact name match. Unmatched columns become `base_only` / `comp_only` and are reported in `cmp$meta$diff`.
2. **Key resolution** — `by` is taken literally if supplied, auto-inferred from shared columns when `by = "auto"`, or rows are matched by row position when `by = NULL`.
3. **Row matching** — keyed-unique join, duplicate-key resolution (`dup_keys`), or position-`zip`; details land in `cmp$meta$matching` (also displayed by `print()` and in the HTML report).
4. **Cell diff** — per-column equality respecting `ks_tol()` (abs / rel / ULP) for numerics, `ks_comp_options()` for strings / NAs / SAS special missings, and `vctrs::vec_cast()` for compatible types.

5. **Pattern detection** — recurring shapes across diffs (constant offset, sign flip, trim-only, ...) are scored and surfaced on `cmp$pattern_summary`.
6. **Manifest** — xxhash64 digests of inputs / options / tolerance plus a timestamp and package version, for QC traceability.

## Examples

```
a <- data.frame(id = 1:3, x = c(1, 2, 3), y = c("a", "b", "c"))
b <- data.frame(id = 1:3, x = c(1, 2, 4), y = c("a", "B", "c"))

# Strict compare on `id`
cmp <- ks_compare(a, b, by = "id")
cmp
tibble::as_tibble(cmp)

# Tolerant numeric compare
ks_compare(
  data.frame(id = 1, x = 1.0),
  data.frame(id = 1, x = 1.0001),
  by = "id",
  tolerance = ks_tol(abs = 1e-3)
)

# Renamed key column
ks_compare(
  data.frame(USUBJID = "S001", x = 1),
  data.frame(SUBJID = "S001", x = 1),
  by = c(USUBJID = "SUBJID")
)
```

---

ks\_glance

*Glance at a ks\_comparison*

---

## Description

Returns a one-row tibble with the same headline counts as `summary.ks_comparison()`. Convenient for binding many comparisons together in a QC loop.

## Usage

```
ks_glance(x, ...)
```

## Arguments

<code>x</code>	A <code>ks_comparison</code> .
<code>...</code>	Unused.

**Value**

A one-row tibble with columns `n_base_rows`, `n_comp_rows`, `n_matched_rows`, `n_base_only_rows`, `n_comp_only_rows`, `n_matched_columns`, `n_value_diffs`, `n_columns_with_diffs`.

**Examples**

```
cmp <- ks_compare(
  data.frame(id = 1:2, x = c(1, 2)),
  data.frame(id = 1:2, x = c(1, 3)),
  by = "id"
)
ks_glance(cmp)
```

---

<code>ks_match_health</code>	<i>Match-quality health check</i>
------------------------------	-----------------------------------

---

**Description**

Computes a small set of metrics that flag *suspicious* matching: high diff density, many fully-different rows, duplicate keys resolved positionally, and similar patterns that often indicate the wrong key was used.

**Usage**

```
ks_match_health(x)
```

**Arguments**

`x`                    A `ks_comparison`.

**Value**

A list with components:

- `diff_density`:  $n\_value\_diffs / (n\_matched\_rows * n\_matched\_columns)$  (0-1).
- `n_fully_diff_rows`: matched rows where every matched column differs (or, more permissively,  $\geq 80\%$  of columns differ).
- `pct_fully_diff_rows`: as proportion of matched rows.
- `n_value_to_na`, `n_na_to_value`: structural NA transitions.
- `dup_keys`: TRUE if either side had duplicate key values.
- `dup_positional`: TRUE if duplicates were paired positionally (`keep_all`) — the most error-prone path.
- `row_count_delta`: `n_comp_rows - n_base_rows` (signed).
- `flags`: character vector of human-readable warnings (possibly empty).
- `severity`: "ok", "info", "warn", or "critical".

**See Also**

[ks\\_recommendations\(\)](#), [ks\\_compare\(\)](#).

**Examples**

```
a <- data.frame(id = c(1, 1, 2), x = c(1, 2, 3))
b <- data.frame(id = c(1, 1, 2), x = c(9, 9, 3))
cmp <- ks_compare(a, b, by = "id", dup_keys = "keep_all")
ks_match_health(cmp)
```

---

ks_recommendations	<i>Actionable recommendations for a comparison</i>
--------------------	--

---

**Description**

Synthesises [ks\\_match\\_health\(\)](#) flags, [ks\\_suggest\\_key\(\)](#) hints, and a few other heuristics into a short list of user-facing "what to look at first" recommendations.

**Usage**

```
ks_recommendations(x)
```

**Arguments**

x                    A `ks_comparison`.

**Value**

A tibble with one row per recommendation:

- severity ("critical", "warn", "info", "ok")
- title short label
- message longer human-readable text
- action suggested next step (may be NA)

**See Also**

[ks\\_match\\_health\(\)](#), [ks\\_suggest\\_key\(\)](#).

---

ks_report_html	Render a self-contained HTML report
----------------	-------------------------------------

---

## Description

Produces a polished, single-file HTML report summarising a `ks_comparison`. The output uses `htmltools` for layout and `reactable` for interactive tables. No internet access, no Quarto, and no Pandoc are required.

## Usage

```
ks_report_html(
  x,
  path = NULL,
  title = "ksCompare report",
  subtitle = NULL,
  max_rows = 100L,
  theme = c("default", "slate"),
  group_by_key = FALSE,
  max_groups = 200L
)
```

## Arguments

<code>x</code>	A <code>ks_comparison</code> object.
<code>path</code>	File path for the report. Three behaviours: <ul style="list-style-type: none"> <li>• missing / <code>NULL</code> (default): a filename is auto-generated from the comparison's <code>base_name</code> (and, when distinct, <code>comp_name</code>) written into <code>options\$path</code> (if set on <code>ks_comp_options()</code>) or the current working directory, e.g. <code>ksCompare_ads1_vs_ads1_qc.html</code>.</li> <li>• <code>NA</code>: nothing is written; the assembled <code>htmltools::tagList()</code> is returned for embedding in another document.</li> <li>• a character path: written to that path. A <code>.html</code> extension is appended if missing. A bare filename (no directory component) is resolved relative to <code>options\$path</code> when that has been set on <code>ks_comp_options()</code>, otherwise relative to the working directory.</li> </ul>
<code>title</code>	Title shown at the top of the report.
<code>subtitle</code>	Optional subtitle (e.g. study identifier or run id).
<code>max_rows</code>	Per-table row cap (default 100). When the value-diff table exceeds this cap, the report shows a <i>smart stratified sample</i> covering each column and each distinct diff-cause ( <b>note</b> ), prioritising the largest numeric magnitudes (at least one example per column and per cause is always retained). A notice indicates the sample size and recommends <code>as_tibble(cmp)</code> for the full set. <code>as_tibble(x)</code> is never truncated.

<code>theme</code>	One of "default" (steel blue) or "slate" (neutral dark headers on light background).
<code>group_by_key</code>	Logical (default FALSE). When TRUE and a <code>by = key</code> was used in <code>ks_compare()</code> , the value-diff section is rendered as one collapsed <code>&lt;details&gt;</code> block per key value, sorted by number of diffs (most-affected first). For <code>dup_keys = "keep_all" / "all_pairs"</code> , each block also shows a <i>Pair</i> column so you can tell pairs apart. Silently falls back to the flat table when no <code>by =</code> was supplied (position match).
<code>max_groups</code>	Maximum number of key-value blocks to render when <code>group_by_key = TRUE</code> (default 200). Excess groups are summarised in a closing notice.

## Details

Layout: a left-hand sticky table of contents, a header with dataset names, a status pill, KPI cards, then sections for schema, row diff, value diff, patterns, and the run manifest. Tables use friendly column labels, column groups, and `tick/cross` markers instead of raw TRUE/FALSE.

## Value

Invisibly returns `path` (or the assembled tag list when `path = NULL`).

## Examples

```
## Not run:
  cmp <- ks_compare(iris, iris, by = "Species")

  # Explicit path
  ks_report_html(cmp, "report.html", title = "ADSL QC")

  # Auto-generated filename in working directory
  ks_report_html(cmp)

  # In-memory tagList for embedding
  tags <- ks_report_html(cmp, path = NA)

## End(Not run)
```

---

`ks_report_xlsx`

*Render an Excel workbook report*

---

## Description

Writes a multi-sheet workbook summarising a `ks_comparison`. Sheets: `Summary`, `Schema`, `KeyDiff`, `Values`, `Patterns`, `OUT_BASE`, `OUT_COMP`, `OUT_DIF`, `OUT_NOEQUAL`, `Manifest`. Cells with value differences are highlighted on the wide `OUT_DIF` sheet.

**Usage**

```
ks_report_xlsx(x, path = NULL, highlight = TRUE, threshold = 0)
```

**Arguments**

<b>x</b>	A <code>ks_comparison</code> object.
<b>path</b>	File path for the workbook. When <code>NULL</code> (default), the filename is auto-generated from the comparison's <code>base_name</code> / <code>comp_name</code> and placed in <code>options\$path</code> (if set on <code>ks_comp_options()</code> ) or the current working directory, e.g. <code>ksCompare_adsl_vs_adsl_qc.xlsx</code> . A <code>.xlsx</code> extension is appended if missing. A bare filename is resolved against <code>options\$path</code> when that has been set.
<b>highlight</b>	If <code>TRUE</code> (default), apply conditional formatting to highlight numeric <code>OUT_DIF</code> cells whose magnitude is above <code>threshold</code> .
<b>threshold</b>	Numeric threshold for highlighting (default 0, so any non-zero diff is highlighted).

**Value**

Invisibly returns `path`.

**Examples**

```
## Not run:
cmp <- ks_compare(iris, iris, by = "Species")
ks_report_xlsx(cmp, "report.xlsx")

# Auto-named output in a pinned folder
cmp2 <- ks_compare(
  iris, iris, by = "Species",
  options = ks_comp_options(path = tempfile("ksCompare_"))
)
ks_report_xlsx(cmp2)

## End(Not run)
```

---

`ks_row_diff_summary` *Per-row diff summary*

---

**Description**

Tallies the number of cell differences per matched observation, sorted descending. Useful for spotting rows where most/all columns disagree (often a sign of a wrong row match, a duplicated key resolved differently, or a systemic shift on a single record).

**Usage**

```
ks_row_diff_summary(x, n = NULL)
```

**Arguments**

<code>x</code>	A <code>ks_comparison</code> .
<code>n</code>	Optional cap on the number of rows returned (default <code>NULL = all</code> ). When set, rows beyond <code>n</code> are dropped.

**Details**

Only matched rows are reported; unmatched rows are available via `ks_unmatched_rows()`.

**Value**

A tibble with one row per matched observation that has at least one cell diff:

- `key_id`, `base_row`, `comp_row`, `key_label`
- `n_diffs`: count of differing cells on this observation.
- `columns`: comma-separated sample of the affected columns (capped at 8 names). Sorted by `n_diffs` descending. Zero-row tibble when there are no value differences.

**See Also**

`ks_compare()`, `ks_cause_summary()`.

**Examples**

```
a <- data.frame(id = 1:3, x = 1:3, y = 1:3)
b <- data.frame(id = 1:3, x = c(1, 9, 9), y = c(1, 9, 3))
cmp <- ks_compare(a, b, by = "id")
ks_row_diff_summary(cmp)
```

---

`ks_set_comp_options` *Set ksCompare global options*

---

**Description**

Thin wrapper around `base::options()` for the `ksCompare.*` namespace. Accepts the same argument names as `ks_comp_options()` (without the `ksCompare.` prefix) and forwards them into `base::options()`, returning the previous values invisibly so the call can be unwound with `base::options()` or `withr::with_options()`.

**Usage**

```
ks_set_comp_options(...)
```

**Arguments**

... Named arguments. Recognised names match the formals of `ks_comp_options()`: `na_equal`, `sas_special_missing`, `compare_labels`, `compare_formats`, `str_trim`, `str_case`, `str_norm`, `tz`, `path`. A single unnamed `ks_comp_options` object may also be passed in, in which case all of its fields are pushed.

## Details

Globals set this way become the new defaults of `ks_comp_options()` (and therefore of `ks_compare()` when no `options =` argument is supplied).

## Value

Invisibly, a named list of the previous values (suitable for `do.call(options, prev)` to restore).

## Examples

```
old <- ks_set_comp_options(path = tempfile("ksCompare_"), str_case = "fold")
ks_comp_options()$path
do.call(options, old) # restore
```

---

<code>ks_suggest_key</code>	<i>Suggest additional key columns when duplicates exist</i>
-----------------------------	---

---

## Description

Scans matched columns to find candidates whose combination with the current by key would make the join key unique on at least one side (and ideally both). Returns columns ranked by how much they reduce duplicate-key cardinality. Useful when `ks_compare()` used `dup_keys = "keep_all"` and the result looks suspicious.

## Usage

```
ks_suggest_key(x, base = NULL, comp = NULL, top_n = 10L)
```

## Arguments

<code>x</code>	A <code>ks_comparison</code> .
<code>base, comp</code>	Optional: the original input data frames used to build <code>x</code> . Required to inspect candidate column uniqueness.
<code>top_n</code>	Maximum number of candidate columns to return.

## Details

Because `ks_compare()` does not snapshot the input data, the original `base` and `comp` frames must be passed in to inspect candidate columns. When `base / comp` are not supplied (or do not contain the original columns) an empty tibble is returned.

## Value

A tibble with columns `column`, `pct_unique_base`, `pct_unique_comp`, `would_make_unique`. Zero rows when no key is in use, no duplicates exist, or no improvement is possible.

**See Also**

[ks\\_match\\_health\(\)](#), [ks\\_compare\(\)](#).

---

ks_sysinfo	<i>SAS SYSINFO-compatible bitmask for a comparison</i>
------------	--

---

**Description**

SAS's `&SYSINFO` macro variable encodes the result of `PROC COMPARE` as a bitmask. `ks_sysinfo()` returns a compatible-style integer summarizing the same kinds of facts, plus a named integer vector showing which bits are set. Bit positions follow SAS documentation:

**Usage**

```
ks_sysinfo(x)
```

**Arguments**

`x`                    A `ks_comparison` object.

**Details**

Bit	Mask	Meaning
1	1	Data set labels differ
2	2	Data set types differ
3	4	Variable has different informat
4	8	Variable has different format
5	16	Variable has different length
6	32	Variable has different label
7	64	Base data set has observation not in compare
8	128	Compare data set has observation not in base
9	256	Base data set has BY group not in compare
10	512	Compare data set has BY group not in base
11	1024	Base data set has variable not in compare
12	2048	Compare data set has variable not in base
13	4096	A value comparison was unequal
14	8192	Conflicting variable types
15	16384	BY variables do not match
16	32768	Fatal error - comparison not done

Length and informat are not currently tracked; their bits are always clear. `haven` does not surface SAS informats, and column "length" is not a meaningful R-side concept for numerics.

**Value**

An integer with class `ks_sysinfo`. Use `as.integer()` for the raw value or print to see the named bits.

**Examples**

```
a <- data.frame(id = 1:3, x = c(1, 2, 3))
b <- data.frame(id = 1:3, x = c(1, 2, 4), z = 1:3)
cmp <- ks_compare(a, b, by = "id")
ks_sysinfo(cmp)
as.integer(ks_sysinfo(cmp))
```

---

`ks_tidy`
*Tidy a ks\_comparison*


---

**Description**

Returns the long-format cell-level diff table. Equivalent to `as_tibble(cmp)` and identical to `cmp$value_diff` when `include_unmatched = FALSE`.

**Usage**

```
ks_tidy(x, ...)

## S3 method for class 'ks_comparison'
ks_tidy(x, include_unmatched = FALSE, ...)
```

**Arguments**

```
x          A ks_comparison.
...        Unused.
include_unmatched Logical (default FALSE). When TRUE, append base-only / comp-only rows
           from cmp$unmatched_rows.
```

**Details**

When `include_unmatched = TRUE`, rows describing observations present on only one side of the comparison are appended to the result with `kind = "base_only" / "comp_only"` and `column_base / column_comp` set to NA. One row is added per unmatched observation, capped at `ks_compare(max_unmatched_rows = ...)`.

**Value**

A tibble of value differences, optionally with appended unmatched-row markers. Has columns `key_id`, `base_row`, `comp_row`, `column_base`, `column_comp`, `kind`, `base`, `comp`, `diff`, `na_flow`, `note`.

## Examples

```
cmp <- ks_compare(
  data.frame(id = 1:2, x = c(1, 2)),
  data.frame(id = 1:3, x = c(1, 3, 4)),
  by = "id"
)
ks_tidy(cmp)
ks_tidy(cmp, include_unmatched = TRUE)
```

---

 ks\_tol

*Tolerance specification for numeric comparisons*


---

## Description

Builds a tolerance object consumed by `ks_compare()`. Two numeric values `a` and `b` are considered equal when **any** of the active rules pass:

## Usage

```
ks_tol(abs = 0, rel = 0, ulp = 0, per_column = NULL)
```

## Arguments

<code>abs</code>	Non-negative absolute tolerance (default 0 — strict equality).
<code>rel</code>	Non-negative relative tolerance, scaled by <code>max(abs(a), abs(b))</code> (default 0).
<code>ulp</code>	Non-negative integer ULP tolerance. Typical values are 4–16; defaults to 0.
<code>per_column</code>	Optional named list of per-column <code>ks_tol()</code> overrides keyed by base-side column name.

## Details

- $\text{abs}(a - b) \leq \text{abs}$
- $\text{abs}(a - b) \leq \text{rel} * \max(\text{abs}(a), \text{abs}(b))$
- $\text{ulp\_distance}(a, b) \leq \text{ulp}$  — IEEE-754 units in the last place, useful for catching floating-point round-trip noise without false positives on genuine differences.

Per-column overrides may be supplied via `per_column`, a named list whose names are **base-side** column names and whose values are themselves the result of `ks_tol()`. Columns not listed fall back to the top-level `abs` / `rel` / `ulp`.

## Value

A `ks_tol` S3 list.

**Examples**

```
ks_tol(abs = 1e-9)
ks_tol(rel = 1e-6, per_column = list(price = ks_tol(abs = 0.005)))
ks_tol(ulp = 4)
```

---

ks_unmatched_rows	<i>Unmatched rows from a comparison</i>
-------------------	---

---

**Description**

Returns the rows that are present on only one side of a comparison ("base\_only" or "comp\_only"), with their original data columns. Mirrors the "Observations in only" tables produced by SAS PROC COMPARE.

**Usage**

```
ks_unmatched_rows(x, side = c("both", "base_only", "comp_only"))
```

**Arguments**

x	A ks_comparison.
side	One of "both" (default), "base_only", or "comp_only".

**Details**

The result is precomputed at `ks_compare()` time and capped by the `max_unmatched_rows` argument (default 100). When the cap kicks in, the `truncated` attribute is set and the full counts are recorded on `n_total`.

**Value**

A tibble with columns `side`, `key_id`, `key_label`, `base_row`, `comp_row`, followed by the data columns from the side that holds each row (`base_row` is NA for `comp_only` rows and vice versa, making it trivial to look an observation up in the original `base / comp` frame). Empty (zero-row) tibble when there are no unmatched rows.

**See Also**

[ks\\_compare\(\)](#), [ks\\_tidy\(\)](#).

**Examples**

```
a <- data.frame(id = 1:3, x = c(1, 2, 3))
b <- data.frame(id = 2:4, x = c(2, 3, 4))
cmp <- ks_compare(a, b, by = "id")
ks_unmatched_rows(cmp)
```

---

 sas-out

*SAS PROC COMPARE-style output datasets*


---

## Description

These helpers reshape the long-format `value_diff` table inside a `ks_comparison` into the four classic SAS PROC COMPARE output datasets:

## Usage

```
as_outbase(x)
```

```
as_outcomp(x)
```

```
as_outdif(x)
```

```
as_outnoequal(x)
```

## Arguments

`x`                    A `ks_comparison` object.

## Details

- `as_outbase()` — values from the base frame, only for matched rows that contain at least one differing column; one row per matched key, one column per compared variable.
- `as_outcomp()` — same shape as `as_outbase()` but values from the compare frame.
- `as_outdif()` — numeric difference (`base - comp`) for every matched row that contains at least one differing column; non-numeric cells are reported as `NA` with a `.note` column.
- `as_outnoequal()` — only the rows where at least one matched cell differs, in long format (one row per differing cell).

The shapes mirror SAS's `OUTBASE=`, `OUTCOMP=`, `OUTDIF=`, and `OUTNOEQUAL=` datasets in spirit; we do not replicate the exact metadata columns (`_TYPE_`, `_OBS_`, ...) but include `key_id` and the original key columns so the rows can be related back to the source frames.

## Value

A tibble.

## Examples

```
a <- data.frame(id = 1:3, x = c(1, 2, 3), y = c("a", "b", "c"))
b <- data.frame(id = 1:3, x = c(1, 2, 4), y = c("a", "B", "c"))
cmp <- ks_compare(a, b, by = "id")
as_outbase(cmp)
```

```
as_outcomp(cmp)
as_outdif(cmp)
as_outnoequal(cmp)
```

---

```
summary.ks_comparison
```

*Summary statistics for a ks\_comparison*

---

## Description

Computes a small list of headline counts that drive the printed summary, the HTML report KPI cards, and CI gates.

## Usage

```
## S3 method for class 'ks_comparison'
summary(object, ...)
```

## Arguments

`object`            A `ks_comparison` returned by `ks_compare()`.  
`...`              Unused; present for S3 conformance.

## Value

A `ks_comparison_summary` list (also pretty-printed via `print()`) with components:

- `n_base_rows`, `n_comp_rows`: input row counts.
- `n_matched_rows`, `n_base_only_rows`, `n_comp_only_rows`: row counts after matching.
- `n_matched_columns`, `n_base_only_columns`, `n_comp_only_columns`: schema-side counts.
- `n_value_diffs`: number of differing cells in matched rows.
- `n_columns_with_diffs`: how many distinct columns hold those differences.

## See Also

[ks\\_glance\(\)](#) for a tibble version, [ks\\_tidy\(\)](#) for the long cell-level table.

## Examples

```
cmp <- ks_compare(
  data.frame(id = 1:3, x = c(1, 2, 3)),
  data.frame(id = 1:3, x = c(1, 2, 4)),
  by = "id"
)
summary(cmp)
```

# Index

`as.integer()`, 20  
`as_ks_comparison`, 2  
`as_outbase` (*sas-out*), 23  
`as_outcomp` (*sas-out*), 23  
`as_outdif` (*sas-out*), 23  
`as_outnoequal` (*sas-out*), 23  
  
`base::options()`, 7, 17  
  
`ks_assert_clean`, 3  
`ks_assert_clean()`, 4  
`ks_cause_summary`, 4  
`ks_cause_summary()`, 17  
`ks_comp_options`, 5  
`ks_comp_options()`, 9, 10, 14, 16–18  
`ks_compare`, 7  
`ks_compare()`, 2–5, 7, 13, 15, 17–19, 21, 22, 24  
`ks_glance`, 11  
`ks_glance()`, 24  
`ks_match_health`, 12  
`ks_match_health()`, 13, 19  
`ks_pointblank_step` (*ks\_assert\_clean*), 3  
`ks_pointblank_step()`, 4  
`ks_recommendations`, 13  
`ks_recommendations()`, 13  
`ks_report_html`, 14  
`ks_report_html()`, 6, 10  
`ks_report_xlsx`, 15  
`ks_report_xlsx()`, 6, 10  
`ks_row_diff_summary`, 16  
`ks_row_diff_summary()`, 5  
`ks_set_comp_options`, 17  
`ks_set_comp_options()`, 7  
`ks_suggest_key`, 18  
`ks_suggest_key()`, 13  
`ks_sysinfo`, 19  
`ks_tidy`, 20  
`ks_tidy()`, 22, 24  
  
`ks_tol`, 21  
`ks_tol()`, 8, 10, 21  
`ks_unmatched_rows`, 22  
`ks_unmatched_rows()`, 9, 10, 17  
  
`pointblank::action_levels()`, 4  
`pointblank::create_agent()`, 3, 4  
`print()`, 9, 10  
  
*sas-out*, 23  
`stringi::stri_trim_both()`, 6  
`summary()`, 10  
`summary.ks_comparison`, 24  
`summary.ks_comparison()`, 11  
  
`tibble::as_tibble()`, 8, 10  
  
`withr::with_options()`, 17