

# Package: dtplyr (via r-universe)

July 20, 2024

**Title** Data Table Back-End for 'dplyr'

**Version** 1.3.1.9000

**Description** Provides a data.table backend for 'dplyr'. The goal of 'dtplyr' is to allow you to write 'dplyr' code that is automatically translated to the equivalent, but usually much faster, data.table code.

**License** MIT + file LICENSE

**URL** <https://dtplyr.tidyverse.org>, <https://github.com/tidyverse/dtplyr>

**BugReports** <https://github.com/tidyverse/dtplyr/issues>

**Depends** R (>= 3.6)

**Imports** cli (>= 3.4.0), data.table (>= 1.13.0), dplyr (>= 1.1.0), glue, lifecycle, rlang (>= 1.0.4), tibble, tidyselect (>= 1.2.0), vctrs (>= 0.4.1)

**Suggests** bench, covr, knitr, rmarkdown, testthat (>= 3.1.2), tidyr (>= 1.1.0), waldo (>= 0.3.1)

**VignetteBuilder** knitr

**Config/Needs/website** tidyverse/tidytemplate

**Config/testthat/edition** 3

**Encoding** UTF-8

**Roxygen** {library(tidyr); list(markdown = TRUE)}

**RoxygenNote** 7.3.1

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

**RemoteUrl** <https://github.com/tidyverse/dtplyr>

**RemoteRef** HEAD

**RemoteSha** 82aee6bb6d4000482a1360204252038d02542228

## Contents

|                                   |           |
|-----------------------------------|-----------|
| arrange.dplyr_step . . . . .      | 2         |
| collect.dplyr_step . . . . .      | 3         |
| complete.dplyr_step . . . . .     | 4         |
| count.dplyr_step . . . . .        | 5         |
| distinct.dplyr_step . . . . .     | 6         |
| drop_na.dplyr_step . . . . .      | 7         |
| expand.dplyr_step . . . . .       | 8         |
| fill.dplyr_step . . . . .         | 9         |
| filter.dplyr_step . . . . .       | 11        |
| group_by.dplyr_step . . . . .     | 12        |
| group_modify.dplyr_step . . . . . | 13        |
| head.dplyr_step . . . . .         | 14        |
| intersect.dplyr_step . . . . .    | 14        |
| lazy_dt . . . . .                 | 15        |
| left_join.dplyr_step . . . . .    | 16        |
| mutate.dplyr_step . . . . .       | 18        |
| nest.dplyr_step . . . . .         | 19        |
| pivot_longer.dplyr_step . . . . . | 20        |
| pivot_wider.dplyr_step . . . . .  | 22        |
| reframe.dplyr_step . . . . .      | 24        |
| relocate.dplyr_step . . . . .     | 25        |
| rename.dplyr_step . . . . .       | 26        |
| replace_na.dplyr_step . . . . .   | 27        |
| select.dplyr_step . . . . .       | 27        |
| separate.dplyr_step . . . . .     | 28        |
| slice.dplyr_step . . . . .        | 29        |
| summarise.dplyr_step . . . . .    | 31        |
| transmute.dplyr_step . . . . .    | 32        |
| unite.dplyr_step . . . . .        | 33        |
| <b>Index</b>                      | <b>34</b> |

---

arrange.dplyr\_step    *Arrange rows by column values*

---

### Description

This is a method for dplyr generic `arrange()`. It is translated to an `order()` call in the `i` argument of `[.data.table]`.

### Usage

```
## S3 method for class 'dplyr_step'
arrange(.data, ..., .by_group = FALSE)
```

**Arguments**

|           |  |
|-----------|--|
| .data     | A <code>lazy_dt()</code> .   |
| ...       | <data-masking> Variables, or functions of variables. Use <code>desc()</code> to sort a variable in descending order. |
| .by_group | If TRUE, will sort first by grouping variable. Applies to grouped data frames only.                                  |

**Examples**

```
library(dplyr, warn.conflicts = FALSE)

dt <- lazy_dt(mtcars)
dt %>% arrange(vs, cyl)
dt %>% arrange(desc(vs), cyl)
dt %>% arrange(across(mpg:disp))
```

---

collect.dplyr\_step    *Force computation of a lazy data.table*

---

**Description**

- `collect()` returns a tibble, grouped if needed.
- `compute()` generates an intermediate assignment in the translation.
- `as.data.table()` returns a data.table.
- `as.data.frame()` returns a data frame.
- `as_tibble()` returns a tibble.

**Usage**

```
## S3 method for class 'dplyr_step'
collect(x, ...)

## S3 method for class 'dplyr_step'
compute(x, name = unique_name(), ...)

## S3 method for class 'dplyr_step'
as.data.table(x, keep.rownames = FALSE, ...)

## S3 method for class 'dplyr_step'
as.data.frame(x, ...)

## S3 method for class 'dplyr_step'
as_tibble(x, ..., .name_repair = "check_unique")
```

**Arguments**

|               |  |
|---------------|--|
| x             | A <a href="#">lazy_dt</a>                  |
| ...           | Arguments used by other methods.           |
| name          | Name of intermediate data.table.           |
| keep.rownames | Ignored as dplyr never preserves rownames. |
| .name_repair  | Treatment of problematic column names      |

**Examples**

```
library(dplyr, warn.conflicts = FALSE)

dt <- lazy_dt(mtcars)

# Generate translation
avg_mpg <- dt %>%
  filter(am == 1) %>%
  group_by(cyl) %>%
  summarise(mpg = mean(mpg))

# Show translation and temporarily compute result
avg_mpg

# compute and return tibble
avg_mpg_tb <- as_tibble(avg_mpg)
avg_mpg_tb

# compute and return data.table
avg_mpg_dt <- data.table::as.data.table(avg_mpg)
avg_mpg_dt

# modify translation to use intermediate assignment
compute(avg_mpg)
```

---

complete.dplyr\_step *Complete a data frame with missing combinations of data*

---

**Description**

This is a method for the tidyr `complete()` generic. This is a wrapper around dplyr translations for `expand()`, `full_join()`, and `replace_na()` that's useful for completing missing combinations of data.

**Usage**

```
## S3 method for class 'dplyr_step'
complete(data, ..., fill = list())
```

**Arguments**

|      |   |
|------|---|
| data | A <code>lazy_dt()</code> .  |
| ...  | <p>&lt;data-masking&gt; Specification of columns to expand or complete. Columns can be atomic vectors or lists.</p> <ul style="list-style-type: none"> <li>• To find all unique combinations of x, y and z, including those not present in the data, supply each variable as a separate argument: <code>expand(df, x, y, z)</code> or <code>complete(df, x, y, z)</code>.</li> <li>• To find only the combinations that occur in the data, use nesting: <code>expand(df, nesting(x, y, z))</code>.</li> <li>• You can combine the two forms. For example, <code>expand(df, nesting(school_id, student_id), date)</code> would produce a row for each present school-student combination for all possible dates.</li> </ul> <p>When used with factors, <code>expand()</code> and <code>complete()</code> use the full set of levels, not just those that appear in the data. If you want to use only the values seen in the data, use <code>forcats::fct_drop()</code>.</p> <p>When used with continuous variables, you may need to fill in values that do not appear in the data: to do so use expressions like <code>year = 2010:2020</code> or <code>year = full_seq(year, 1)</code>.</p> |
| fill | A named list that for each variable supplies a single value to use instead of NA for missing combinations.  |

**Examples**

```
library(tidyr)
tbl <- tibble(x = 1:2, y = 1:2, z = 3:4)
dt <- lazy_dt(tbl)

dt %>%
  complete(x, y)

dt %>%
  complete(x, y, fill = list(z = 10L))
```

---

|                  |                                    |
|------------------|------------------------------------|
| count.dplyr_step | <i>Count observations by group</i> |
|------------------|------------------------------------|

---

**Description**

This is a method for the dplyr `count()` generic. It is translated using `.N` in the `j` argument, and supplying groups to `keyby` as appropriate.

**Usage**

```
## S3 method for class 'dplyr_step'
count(x, ..., wt = NULL, sort = FALSE, name = NULL)
```

**Arguments**

|      |   |
|------|---|
| x    | A <code>lazy_dt()</code>  |
| ...  | <data-masking> Variables to group by.   |
| wt   | <data-masking> Frequency weights. Can be NULL or a variable: <ul style="list-style-type: none"> <li>• If NULL (the default), counts the number of rows in each group.</li> <li>• If a variable, computes <code>sum(wt)</code> for each group.</li> </ul>  |
| sort | If TRUE, will show the largest groups at the top.   |
| name | The name of the new column in the output.<br>If omitted, it will default to <code>n</code> . If there's already a column called <code>n</code> , it will use <code>nn</code> . If there's a column called <code>n</code> and <code>nn</code> , it'll use <code>nnn</code> , and so on, adding <code>ns</code> until it gets a new name. |

**Examples**

```
library(dplyr, warn.conflicts = FALSE)

dt <- lazy_dt(dplyr::starwars)
dt %>% count(species)
dt %>% count(species, sort = TRUE)
dt %>% count(species, wt = mass, sort = TRUE)
```

---

`distinct.dplyr_step` *Subset distinct/unique rows*

---

**Description**

This is a method for the dplyr `distinct()` generic. It is translated to `data.table::unique.data.table()`.

**Usage**

```
## S3 method for class 'dplyr_step'
distinct(.data, ..., .keep_all = FALSE)
```

**Arguments**

|           |   |
|-----------|---|
| .data     | A <code>lazy_dt()</code>  |
| ...       | <data-masking> Optional variables to use when determining uniqueness. If there are multiple rows for a given combination of inputs, only the first row will be preserved. If omitted, will use all variables in the data frame. |
| .keep_all | If TRUE, keep all variables in <code>.data</code> . If a combination of <code>...</code> is not distinct, this keeps the first row of values.   |

## Examples

```
library(dplyr, warn.conflicts = FALSE)
df <- lazy_dt(data.frame(
  x = sample(10, 100, replace = TRUE),
  y = sample(10, 100, replace = TRUE)
))

df %>% distinct(x)
df %>% distinct(x, y)
df %>% distinct(x, .keep_all = TRUE)
```

---

drop\_na.dtplyr\_step    *Drop rows containing missing values*

---

## Description

This is a method for the tidyr `drop_na()` generic. It is translated to `data.table::na.omit()`

## Usage

```
## S3 method for class 'dtplyr_step'
drop_na(data, ...)
```

## Arguments

|                   |   |
|-------------------|---|
| <code>data</code> | A <code>lazy_dt()</code> .  |
| <code>...</code>  | <code>&lt;tidy-select&gt;</code> Columns to inspect for missing values. If empty, all columns are used. |

## Examples

```
library(dplyr)
library(tidyr)

dt <- lazy_dt(tibble(x = c(1, 2, NA), y = c("a", NA, "b")))
dt %>% drop_na()
dt %>% drop_na(x)

vars <- "y"
dt %>% drop_na(x, any_of(vars))
```

---

expand.dtplyr\_step      *Expand data frame to include all possible combinations of values.*

---

## Description

This is a method for the tidy `expand()` generic. It is translated to `data.table::CJ()`.

## Usage

```
## S3 method for class 'dtplyr_step'
expand(data, ..., .name_repair = "check_unique")
```

## Arguments

`data`      A `lazy_dt()`.

`...`      Specification of columns to expand. Columns can be atomic vectors or lists.

- To find all unique combinations of `x`, `y` and `z`, including those not present in the data, supply each variable as a separate argument: `expand(df, x, y, z)`.
- To find only the combinations that occur in the data, use `nesting`: `expand(df, nesting(x, y, z))`.
- You can combine the two forms. For example, `expand(df, nesting(school_id, student_id), date)` would produce a row for each present school-student combination for all possible dates.

Unlike the `data.frame` method, this method does not use the full set of levels, just those that appear in the data.

When used with continuous variables, you may need to fill in values that do not appear in the data: to do so use expressions like `year = 2010:2020` or `year = full_seq(year, 1)`.

`.name_repair`      Treatment of problematic column names:

- "minimal": No name repair or checks, beyond basic existence,
- "unique": Make sure names are unique and not empty,
- "check\_unique": (default value), no name repair, but check they are unique,
- "universal": Make the names unique and syntactic
- a function: apply custom name repair (e.g., `.name_repair = make.names` for names in the style of base R).
- A purrr-style anonymous function, see `rlang::as_function()`

This argument is passed on as `repair` to `vctrs::vec_as_names()`. See there for more details on these terms and the strategies used to enforce them.



**Examples**

```

library(tidyr)

fruits <- lazy_dt(tibble(
  type = c("apple", "orange", "apple", "orange", "orange", "orange"),
  year = c(2010, 2010, 2012, 2010, 2010, 2012),
  size = factor(
    c("XS", "S", "M", "S", "S", "M"),
    levels = c("XS", "S", "M", "L")
  ),
  weights = rnorm(6, as.numeric(size) + 2)
))

# All possible combinations -----
# Note that only present levels of the factor variable `size` are retained.
fruits %>% expand(type)
fruits %>% expand(type, size)

# This is different from the data frame behaviour:
fruits %>% dplyr::collect() %>% expand(type, size)

# Other uses -----
fruits %>% expand(type, size, 2010:2012)

# Use `anti_join()` to determine which observations are missing
all <- fruits %>% expand(type, size, year)
all
all %>% dplyr::anti_join(fruits)

# Use with `right_join()` to fill in missing rows
fruits %>% dplyr::right_join(all)

```

---

|                 |   |
|-----------------|---|
| fill.dplyr_step | <i>Fill in missing values with previous or next value</i> |
|-----------------|---|

---

**Description**

This is a method for the tidyr fill() generic. It is translated to `data.table::nafill()`. Note that `data.table::nafill()` currently only works for integer and double columns.

**Usage**

```

## S3 method for class 'dplyr_step'
fill(data, ..., .direction = c("down", "up", "downup", "updown"))

```

**Arguments**

|      |   |
|------|---|
| data | A data frame.                                     |
| ...  | <code>&lt;tidy-select&gt;</code> Columns to fill. |

`.direction` Direction in which to fill missing values. Currently either "down" (the default), "up", "downup" (i.e. first down and then up) or "updown" (first up and then down).

### Examples

```
library(tidyr)

# Value (year) is recorded only when it changes
sales <- lazy_dt(tibble::tribble(
  ~quarter, ~year, ~sales,
  "Q1", 2000, 66013,
  "Q2", NA, 69182,
  "Q3", NA, 53175,
  "Q4", NA, 21001,
  "Q1", 2001, 46036,
  "Q2", NA, 58842,
  "Q3", NA, 44568,
  "Q4", NA, 50197,
  "Q1", 2002, 39113,
  "Q2", NA, 41668,
  "Q3", NA, 30144,
  "Q4", NA, 52897,
  "Q1", 2004, 32129,
  "Q2", NA, 67686,
  "Q3", NA, 31768,
  "Q4", NA, 49094
))

# `fill()` defaults to replacing missing data from top to bottom
sales %>% fill(year)

# Value (n_squirrels) is missing above and below within a group
squirrels <- lazy_dt(tibble::tribble(
  ~group, ~name, ~role, ~n_squirrels,
  1, "Sam", "Observer", NA,
  1, "Mara", "Scorekeeper", 8,
  1, "Jesse", "Observer", NA,
  1, "Tom", "Observer", NA,
  2, "Mike", "Observer", NA,
  2, "Rachael", "Observer", NA,
  2, "Sydekea", "Scorekeeper", 14,
  2, "Gabriela", "Observer", NA,
  3, "Derrick", "Observer", NA,
  3, "Kara", "Scorekeeper", 9,
  3, "Emily", "Observer", NA,
  3, "Danielle", "Observer", NA
))

# The values are inconsistently missing by position within the group
# Use .direction = "downup" to fill missing values in both directions
squirrels %>%
  dplyr::group_by(group) %>%
```

```
fill(n_squirrels, .direction = "downup") %>%
  dplyr::ungroup()

# Using `direction = "updown"` accomplishes the same goal in this example
```

---

filter.dtplyr\_step      *Subset rows using column values*

---

### Description

This is a method for the dplyr `arrange()` generic. It is translated to the `i` argument of `[.data.table]`

### Usage

```
## S3 method for class 'dtplyr_step'
filter(.data, ..., .by = NULL, .preserve = FALSE)
```

### Arguments

|                        |   |
|------------------------|---|
| <code>.data</code>     | A <code>lazy_dt()</code> .  |
| <code>...</code>       | <code>&lt;data-masking&gt;</code> Expressions that return a logical value, and are defined in terms of the variables in <code>.data</code> . If multiple expressions are included, they are combined with the <code>&amp;</code> operator. Only rows for which all conditions evaluate to <code>TRUE</code> are kept. |
| <code>.by</code>       | <b>[Experimental]</b><br><code>&lt;tidy-select&gt;</code> Optionally, a selection of columns to group by for just this operation, functioning as an alternative to <code>group_by()</code> . For details and examples, see <code>?dplyr_by</code> .   |
| <code>.preserve</code> | Ignored   |

### Examples

```
library(dplyr, warn.conflicts = FALSE)

dt <- lazy_dt(mtcars)
dt %>% filter(cyl == 4)
dt %>% filter(vs, am)

dt %>%
  group_by(cyl) %>%
  filter(mpg > mean(mpg))
```

---

group\_by.dtplyr\_step *Group and ungroup*

---

### Description

These are methods for dplyr's `group_by()` and `ungroup()` generics. Grouping is translated to the either `keyby` and `by` argument of `[.data.table]` depending on the value of the `arrange` argument.

### Usage

```
## S3 method for class 'dtplyr_step'
group_by(.data, ..., .add = FALSE, arrange = TRUE)

## S3 method for class 'dtplyr_step'
ungroup(x, ...)
```

### Arguments

|                        |   |
|------------------------|---|
| <code>.data</code>     | A <code>lazy_dt()</code>  |
| <code>...</code>       | In <code>group_by()</code> , variables or computations to group by. Computations are always done on the ungrouped data frame. To perform computations on the grouped data, you need to use a separate <code>mutate()</code> step before the <code>group_by()</code> . Computations are not allowed in <code>nest_by()</code> . In <code>ungroup()</code> , variables to remove from the grouping. |
| <code>.add, add</code> | When <code>FALSE</code> , the default, <code>group_by()</code> will override existing groups. To add to the existing groups, use <code>.add = TRUE</code> .<br>This argument was previously called <code>add</code> , but that prevented creating a new grouping variable called <code>add</code> , and conflicts with our naming conventions.  |
| <code>arrange</code>   | If <code>TRUE</code> , will automatically arrange the output of subsequent grouped operations by group. If <code>FALSE</code> , output order will be left unchanged. In the generated <code>data.table</code> code this switches between using the <code>keyby</code> ( <code>TRUE</code> ) and <code>by</code> ( <code>FALSE</code> ) arguments.   |
| <code>x</code>         | A <code>tbl()</code>  |

### Examples

```
library(dplyr, warn.conflicts = FALSE)
dt <- lazy_dt(mtcars)

# group_by() is usually translated to `keyby` so that the groups
# are ordered in the output
dt %>%
  group_by(cyl) %>%
  summarise(mpg = mean(mpg))

# use `arrange = FALSE` to instead use `by` so the original order
# or groups is preserved
```

```
dt %>%
  group_by(cyl, arrange = FALSE) %>%
  summarise(mpg = mean(mpg))
```

---

```
group_modify.dtplyr_step
```

*Apply a function to each group*

---

## Description

These are methods for the dplyr `group_map()` and `group_modify()` generics. They are both translated to `[.data.table]`.

## Usage

```
## S3 method for class 'dtplyr_step'
group_modify(.data, .f, ..., keep = FALSE)

## S3 method for class 'dtplyr_step'
group_map(.data, .f, ..., keep = FALSE)
```

## Arguments

|                    |  |
|--------------------|--|
| <code>.data</code> | A <code>lazy_dt()</code>   |
| <code>.f</code>    | The name of a two argument function. The first argument is passed <code>.SD</code> , the <code>data.table</code> representing the current group; the second argument is passed <code>.BY</code> , a list giving the current values of the grouping variables. The function should return a list or <code>data.table</code> . |
| <code>...</code>   | Additional arguments passed to <code>.f</code>   |
| <code>keep</code>  | Not supported for <code>lazy_dt</code> .   |

## Value

`group_map()` applies `.f` to each group, returning a list. `group_modify()` replaces each group with the results of `.f`, returning a modified `lazy_dt()`.

## Examples

```
library(dplyr)

dt <- lazy_dt(mtcars)

dt %>%
  group_by(cyl) %>%
  group_modify(head, n = 2L)

dt %>%
  group_by(cyl) %>%
  group_map(head, n = 2L)
```

---

head.dtplyr\_step      *Subset first or last rows*

---

### Description

These are methods for the base generics `head()` and `tail()`. They are not translated.

### Usage

```
## S3 method for class 'dtplyr_step'
head(x, n = 6L, ...)
```

```
## S3 method for class 'dtplyr_step'
tail(x, n = 6L, ...)
```

### Arguments

|     |  |
|-----|--|
| x   | A <code>lazy_dt()</code>   |
| n   | Number of rows to select. Can use a negative number to instead drop rows from the other end. |
| ... | Passed on to <code>head()/tail()</code> .  |

### Examples

```
library(dplyr, warn.conflicts = FALSE)
dt <- lazy_dt(data.frame(x = 1:10))

# first three rows
head(dt, 3)
# last three rows
tail(dt, 3)

# drop first three rows
tail(dt, -3)
```

---

intersect.dtplyr\_step      *Set operations*

---

### Description

These are methods for the dplyr generics `intersect()`, `union()`, `union_all()`, and `setdiff()`. They are translated to `data.table::fintersect()`, `data.table::funion()`, and `data.table::fsetdiff()`.

**Usage**

```
## S3 method for class 'dtplyr_step'
intersect(x, y, ...)

## S3 method for class 'dtplyr_step'
union(x, y, ...)

## S3 method for class 'dtplyr_step'
union_all(x, y, ...)

## S3 method for class 'dtplyr_step'
setdiff(x, y, ...)
```

**Arguments**

|      |                                     |
|------|-------------------------------------|
| x, y | A pair of <code>lazy_dt()</code> s. |
| ...  | Ignored                             |

**Examples**

```
dt1 <- lazy_dt(data.frame(x = 1:4))
dt2 <- lazy_dt(data.frame(x = c(2, 4, 6)))

intersect(dt1, dt2)
union(dt1, dt2)
setdiff(dt1, dt2)
```

---

lazy\_dt

---

*Create a "lazy" data.table for use with dplyr verbs*


---

**Description**

A lazy `data.table` captures the intent of dplyr verbs, only actually performing computation when requested (with `collect()`, `pull()`, `as.data.frame()`, `data.table::as.data.table()`, or `tibble::as_tibble()`). This allows dtplyr to convert dplyr verbs into as few `data.table` expressions as possible, which leads to a high performance translation.

See vignette("translation") for the details of the translation.

**Usage**

```
lazy_dt(x, name = NULL, immutable = TRUE, key_by = NULL)
```

**Arguments**

|           |   |
|-----------|---|
| x         | A data table (or something that can be coerced to a data table).  |
| name      | Optionally, supply a name to be used in generated expressions. For expert use only.   |
| immutable | If TRUE, x is treated as immutable and will never be modified by any code generated by dtplyr. Alternatively, you can set <code>immutable = FALSE</code> to allow dtplyr to modify the input object.  |
| key_by    | Set keys for data frame, using <code>select()</code> semantics (e.g. <code>key_by = c(key1, key2)</code> ).<br>This uses <code>data.table::setkey()</code> to sort the table and build an index. This will considerably improve performance for subsets, summaries, and joins that use the keys.<br>See <code>vignette("datatable-keys-fast-subset")</code> for more details. |

**Examples**

```
library(dplyr, warn.conflicts = FALSE)

mtcars2 <- lazy_dt(mtcars)
mtcars2
mtcars2 %>% select(mpg:cyl)
mtcars2 %>% select(x = mpg, y = cyl)
mtcars2 %>% filter(cyl == 4) %>% select(mpg)
mtcars2 %>% select(mpg, cyl) %>% filter(cyl == 4)
mtcars2 %>% mutate(cyl2 = cyl * 2, cyl4 = cyl2 * 2)
mtcars2 %>% transmute(cyl2 = cyl * 2, vs2 = vs * 2)
mtcars2 %>% filter(cyl == 8) %>% mutate(cyl2 = cyl * 2)

# Learn more about translation in vignette("translation")
by_cyl <- mtcars2 %>% group_by(cyl)
by_cyl %>% summarise(mpg = mean(mpg))
by_cyl %>% mutate(mpg = mean(mpg))
by_cyl %>%
  filter(mpg < mean(mpg)) %>%
  summarise(hp = mean(hp))
```

---

left\_join.dtplyr\_step *Join data tables*

---

**Description**

These are methods for the dplyr generics `left_join()`, `right_join()`, `inner_join()`, `full_join()`, `anti_join()`, and `semi_join()`. Left, right, inner, and anti join are translated to the `[.data.table]` equivalent, full joins to `data.table::merge.data.table()`. Left, right, and full joins are in some cases followed by calls to `data.table::setcolorder()` and `data.table::setnames()` to ensure that column order and names match dplyr conventions. Semi-joins don't have a direct `data.table` equivalent.



**Usage**

```
## S3 method for class 'dplyr_step'
left_join(x, y, ..., by = NULL, copy = FALSE, suffix = c(".x", ".y"))
```

**Arguments**

|        |   |
|--------|---|
| x, y   | A pair of <code>lazy_dt()</code> s.   |
| ...    | Other parameters passed onto methods.   |
| by     | A join specification created with <code>join_by()</code> , or a character vector of variables to join by.<br><br>If NULL, the default, <code>*_join()</code> will perform a natural join, using all variables in common across x and y. A message lists the variables so that you can check they're correct; suppress the message by supplying <code>by</code> explicitly.<br><br>To join on different variables between x and y, use a <code>join_by()</code> specification. For example, <code>join_by(a == b)</code> will match <code>x\$a</code> to <code>y\$b</code> .<br><br>To join by multiple variables, use a <code>join_by()</code> specification with multiple expressions. For example, <code>join_by(a == b, c == d)</code> will match <code>x\$a</code> to <code>y\$b</code> and <code>x\$c</code> to <code>y\$d</code> . If the column names are the same between x and y, you can shorten this by listing only the variable names, like <code>join_by(a, c)</code> .<br><br><code>join_by()</code> can also be used to perform inequality, rolling, and overlap joins. See the documentation at <code>?join_by</code> for details on these types of joins.<br><br>For simple equality joins, you can alternatively specify a character vector of variable names to join by. For example, <code>by = c("a", "b")</code> joins <code>x\$a</code> to <code>y\$a</code> and <code>x\$b</code> to <code>y\$b</code> . If variable names differ between x and y, use a named character vector like <code>by = c("x_a" = "y_a", "x_b" = "y_b")</code> .<br><br>To perform a cross-join, generating all combinations of x and y, see <code>cross_join()</code> . |
| copy   | If x and y are not from the same data source, and <code>copy</code> is TRUE, then y will be copied into the same src as x. This allows you to join tables across srcs, but it is a potentially expensive operation so you must opt into it.   |
| suffix | If there are non-joined duplicate variables in x and y, these suffixes will be added to the output to disambiguate them. Should be a character vector of length 2.  |

**Examples**

```
library(dplyr, warn.conflicts = FALSE)

band_dt <- lazy_dt(dplyr::band_members)
instrument_dt <- lazy_dt(dplyr::band_instruments)

band_dt %>% left_join(instrument_dt)
band_dt %>% right_join(instrument_dt)
band_dt %>% inner_join(instrument_dt)
band_dt %>% full_join(instrument_dt)

band_dt %>% semi_join(instrument_dt)
band_dt %>% anti_join(instrument_dt)
```

---

mutate.dtplyr\_step      *Create and modify columns*

---

## Description

This is a method for the dplyr `mutate()` generic. It is translated to the `j` argument of `[.data.table]`, using `:=` to modify "in place". If `.before` or `.after` is provided, the new columns are relocated with a call to `data.table::setcolorder()`.

## Usage

```
## S3 method for class 'dtplyr_step'
mutate(
  .data,
  ...,
  .by = NULL,
  .keep = c("all", "used", "unused", "none"),
  .before = NULL,
  .after = NULL
)
```

## Arguments

|                    |  |
|--------------------|--|
| <code>.data</code> | A <code>lazy_dt()</code> .   |
| <code>...</code>   | <p>&lt;data-masking&gt; Name-value pairs. The name gives the name of the column in the output.</p> <p>The value can be:</p> <ul style="list-style-type: none"> <li>• A vector of length 1, which will be recycled to the correct length.</li> <li>• A vector the same length as the current group (or the whole data frame if ungrouped).</li> <li>• NULL, to remove the column.</li> <li>• A data frame or tibble, to create multiple columns in the output.</li> </ul>   |
| <code>.by</code>   | <p><b>[Experimental]</b></p> <p>&lt;tidy-select&gt; Optionally, a selection of columns to group by for just this operation, functioning as an alternative to <code>group_by()</code>. For details and examples, see <code>?dplyr_by</code>.</p>  |
| <code>.keep</code> | <p>Control which columns from <code>.data</code> are retained in the output. Grouping columns and columns created by <code>...</code> are always kept.</p> <ul style="list-style-type: none"> <li>• "all" retains all columns from <code>.data</code>. This is the default.</li> <li>• "used" retains only the columns used in <code>...</code> to create new columns. This is useful for checking your work, as it displays inputs and outputs side-by-side.</li> <li>• "unused" retains only the columns <i>not</i> used in <code>...</code> to create new columns. This is useful if you generate new columns, but no longer need the columns used to generate them.</li> </ul> |

- "none" doesn't retain any extra columns from .data. Only the grouping variables and columns created by ... are kept.

Note: With dplyr .keep will only work with column names passed as symbols, and won't work with other workflows (e.g. eval(parse(text = "x + 1")))

.before, .after **<tidy-select>** Optionally, control where new columns should appear (the default is to add to the right hand side). See [relocate\(\)](#) for more details.

## Examples

```
library(dplyr, warn.conflicts = FALSE)

dt <- lazy_dt(data.frame(x = 1:5, y = 5:1))
dt %>%
  mutate(a = (x + y) / 2, b = sqrt(x^2 + y^2))

# It uses a more sophisticated translation when newly created variables
# are used in the same expression
dt %>%
  mutate(x1 = x + 1, x2 = x1 + 1)
```

---

|                 |             |
|-----------------|-------------|
| nest.dplyr_step | <i>Nest</i> |
|-----------------|-------------|

---

## Description

This is a method for the tidyr [tidyr::nest\(\)](#) generic. It is translated using the non-nested variables in the by argument and .SD in the j argument.

## Usage

```
## S3 method for class 'dplyr_step'
nest(.data, ..., .names_sep = NULL, .key = deprecated())
```

## Arguments

|            |   |
|------------|---|
| .data      | A data frame.   |
| ...        | <b>&lt;tidy-select&gt;</b> Columns to nest, specified using name-variable pairs of the form new_col = c(col1, col2, col3). The right hand side can be any valid tidy select expression.   |
| .names_sep | If NULL, the default, the inner names will come from the former outer names. If a string, the new inner names will use the outer names with names_sep automatically stripped. This makes names_sep roughly symmetric between nesting and unnesting. |
| .key       | Not supported.  |
| data       | A <a href="#">lazy_dt()</a> .   |

## Examples

```
if (require("tidyr", quietly = TRUE)) {  
  dt <- lazy_dt(tibble(x = c(1, 2, 1), y = c("a", "a", "b")))  
  dt %>% nest(data = y)  
  
  dt %>% dplyr::group_by(x) %>% nest()  
}
```

---

`pivot_longer.dtplyr_step`

*Pivot data from wide to long*

---

## Description

This is a method for the tidyr `pivot_longer()` generic. It is translated to `data.table::melt()`

## Usage

```
## S3 method for class 'dtplyr_step'  
pivot_longer(  
  data,  
  cols,  
  names_to = "name",  
  names_prefix = NULL,  
  names_sep = NULL,  
  names_pattern = NULL,  
  names_ptypes = NULL,  
  names_transform = NULL,  
  names_repair = "check_unique",  
  values_to = "value",  
  values_drop_na = FALSE,  
  values_ptypes = NULL,  
  values_transform = NULL,  
  ...  
)
```

## Arguments

|                       |  |
|-----------------------|--|
| <code>data</code>     | A <code>lazy_dt()</code> .   |
| <code>cols</code>     | <code>&lt;tidy-select&gt;</code> Columns to pivot into longer format.  |
| <code>names_to</code> | A character vector specifying the new column or columns to create from the information stored in the column names of data specified by <code>cols</code> . <ul style="list-style-type: none"><li>• If length 0, or if NULL is supplied, no columns will be created.</li><li>• If length 1, a single column will be created which will contain the column names specified by <code>cols</code>.</li></ul> |

- If `length > 1`, multiple columns will be created. In this case, one of `names_sep` or `names_pattern` must be supplied to specify how the column names should be split. There are also two additional character values you can take advantage of:
    - `NA` will discard the corresponding component of the column name.
    - `".value"` indicates that the corresponding component of the column name defines the name of the output column containing the cell values, overriding `values_to` entirely.
- `names_prefix` A regular expression used to remove matching text from the start of each variable name.
- `names_sep, names_pattern`  
 If `names_to` contains multiple values, these arguments control how the column name is broken up.  
`names_sep` takes the same specification as `separate()`, and can either be a numeric vector (specifying positions to break on), or a single string (specifying a regular expression to split on).  
`names_pattern` takes the same specification as `extract()`, a regular expression containing matching groups (`()`).  
 If these arguments do not give you enough control, use `pivot_longer_spec()` to create a `spec` object and process manually as needed.
- `names_ptypes, names_transform, values_ptypes, values_transform`  
 Not currently supported by `dplyr`.
- `names_repair` What happens if the output has invalid column names? The default, `"check_unique"` is to error if the columns are duplicated. Use `"minimal"` to allow duplicates in the output, or `"unique"` to de-duplicated by adding numeric suffixes. See `vctrs::vec_as_names()` for more options.
- `values_to` A string specifying the name of the column to create from the data stored in cell values. If `names_to` is a character containing the special `.value` sentinel, this value will be ignored, and the name of the value column will be derived from part of the existing column names.
- `values_drop_na` If `TRUE`, will drop rows that contain only `NA`s in the `value_to` column. This effectively converts explicit missing values to implicit missing values, and should generally be used only when missing values in data were created by its structure.
- `...` Additional arguments passed on to methods.

## Examples

```
library(tidyr)

# Simplest case where column names are character data
relig_income_dt <- lazy_dt(relig_income)
relig_income_dt %>%
  pivot_longer(!religion, names_to = "income", values_to = "count")

# Slightly more complex case where columns have common prefix,
# and missing missings are structural so should be dropped.
```

```
billboard_dt <- lazy_dt(billboard)
billboard %>%
  pivot_longer(
    cols = starts_with("wk"),
    names_to = "week",
    names_prefix = "wk",
    values_to = "rank",
    values_drop_na = TRUE
  )

# Multiple variables stored in column names
lazy_dt(who) %>%
  pivot_longer(
    cols = new_sp_m014:newrel_f65,
    names_to = c("diagnosis", "gender", "age"),
    names_pattern = "new_?(.*)_(.)(.*)",
    values_to = "count"
  )

# Multiple observations per row
anscombe_dt <- lazy_dt(anscombe)
anscombe_dt %>%
  pivot_longer(
    everything(),
    names_to = c(".value", "set"),
    names_pattern = "(.)(.)(.)(.)(.)(.)"
  )
```

---

pivot\_wider.dplyr\_step

*Pivot data from long to wide*

---

## Description

This is a method for the `tidyr::pivot_wider()` generic. It is translated to `data.table::dcast()`

## Usage

```
## S3 method for class 'dplyr_step'
pivot_wider(
  data,
  id_cols = NULL,
  names_from = name,
  names_prefix = "",
  names_sep = "_",
  names_glue = NULL,
  names_sort = FALSE,
  names_repair = "check_unique",
  values_from = value,
```

```

  values_fill = NULL,
  values_fn = NULL,
  ...
)

```

## Arguments

|                         |  |
|-------------------------|--|
| data                    | A <code>lazy_dt()</code> .   |
| id_cols                 | <code>&lt;tidy-select&gt;</code> A set of columns that uniquely identify each observation. Typically used when you have redundant variables, i.e. variables whose values are perfectly correlated with existing variables.<br>Defaults to all columns in data except for the columns specified through <code>names_from</code> and <code>values_from</code> . If a tidyselect expression is supplied, it will be evaluated on data after removing the columns specified through <code>names_from</code> and <code>values_from</code> . |
| names_from, values_from | <code>&lt;tidy-select&gt;</code> A pair of arguments describing which column (or columns) to get the name of the output column ( <code>names_from</code> ), and which column (or columns) to get the cell values from ( <code>values_from</code> ).<br>If <code>values_from</code> contains multiple values, the value will be added to the front of the output column.  |
| names_prefix            | String added to the start of every variable name. This is particularly useful if <code>names_from</code> is a numeric vector and you want to create syntactic variable names.  |
| names_sep               | If <code>names_from</code> or <code>values_from</code> contains multiple variables, this will be used to join their values together into a single string to use as a column name.  |
| names_glue              | Instead of <code>names_sep</code> and <code>names_prefix</code> , you can supply a glue specification that uses the <code>names_from</code> columns (and special <code>.value</code> ) to create custom column names.  |
| names_sort              | Should the column names be sorted? If <code>FALSE</code> , the default, column names are ordered by first appearance.  |
| names_repair            | What happens if the output has invalid column names? The default, <code>"check_unique"</code> is to error if the columns are duplicated. Use <code>"minimal"</code> to allow duplicates in the output, or <code>"unique"</code> to de-duplicated by adding numeric suffixes. See <code>vctrs::vec_as_names()</code> for more options.  |
| values_fill             | Optionally, a (scalar) value that specifies what each value should be filled in with when missing.<br>This can be a named list if you want to apply different fill values to different value columns.  |
| values_fn               | A function, the default is <code>length()</code> . Note this is different behavior than <code>tidyr::pivot_wider()</code> , which returns a list column by default.  |
| ...                     | Additional arguments passed on to methods.   |

## Examples

```
library(tidyr)
```

```

fish_encounters_dt <- lazy_dt(fish_encounters)
fish_encounters_dt
fish_encounters_dt %>%
  pivot_wider(names_from = station, values_from = seen)
# Fill in missing values
fish_encounters_dt %>%
  pivot_wider(names_from = station, values_from = seen, values_fill = 0)

# Generate column names from multiple variables
us_rent_income_dt <- lazy_dt(us_rent_income)
us_rent_income_dt
us_rent_income_dt %>%
  pivot_wider(names_from = variable, values_from = c(estimate, moe))

# When there are multiple `names_from` or `values_from`, you can use
# use `names_sep` or `names_glue` to control the output variable names
us_rent_income_dt %>%
  pivot_wider(
    names_from = variable,
    names_sep = ".",
    values_from = c(estimate, moe)
  )

# Can perform aggregation with values_fn
warpbreaks_dt <- lazy_dt(as_tibble(warpbreaks[c("wool", "tension", "breaks")]))
warpbreaks_dt
warpbreaks_dt %>%
  pivot_wider(
    names_from = wool,
    values_from = breaks,
    values_fn = mean
  )

```

---

reframe.dtplyr\_step    *Summarise each group to one row*

---

## Description

This is a method for the dplyr `reframe()` generic. It is translated to the `j` argument of `[.data.table]`.

## Usage

```
## S3 method for class 'dtplyr_step'
reframe(.data, ..., .by = NULL)
```

## Arguments

`.data`            A `lazy_dt()`.



... [<data-masking>](#)  
 Name-value pairs of functions. The name will be the name of the variable in the result. The value can be a vector of any length.  
 Unnamed data frame values add multiple columns from a single expression.

.by **[Experimental]**  
[<tidy-select>](#) Optionally, a selection of columns to group by for just this operation, functioning as an alternative to `group_by()`. For details and examples, see `?dplyr_by`.

## Examples

```
library(dplyr, warn.conflicts = FALSE)

dt <- lazy_dt(mtcars)

dt %>%
  reframe(qs = quantile(displacement, c(0.25, 0.75)),
          prob = c(0.25, 0.75),
          .by = cyl)

dt %>%
  group_by(cyl) %>%
  reframe(qs = quantile(displacement, c(0.25, 0.75)),
          prob = c(0.25, 0.75))
```

---

relocate.dtplyr\_step *Relocate variables using their names*

---

## Description

This is a method for the dplyr `relocate()` generic. It is translated to the `j` argument of `[.data.table]`.

## Usage

```
## S3 method for class 'dtplyr_step'
relocate(.data, ..., .before = NULL, .after = NULL)
```

## Arguments

.data A `lazy_dt()`.

... [<tidy-select>](#) Columns to move.

.before, .after [<tidy-select>](#) Destination of columns selected by ... Supplying neither will move columns to the left-hand side; specifying both is an error.

**Examples**

```
library(dplyr, warn.conflicts = FALSE)

dt <- lazy_dt(data.frame(x = 1, y = 2, z = 3))

dt %>% relocate(z)
dt %>% relocate(y, .before = x)
dt %>% relocate(y, .after = y)
```

---

|                    |   |
|--------------------|---|
| rename.dtplyr_step | <i>Rename columns using their names</i> |
|--------------------|---|

---

**Description**

These are methods for the dplyr generics `rename()` and `rename_with()`. They are both translated to `data.table::setnames()`.

**Usage**

```
## S3 method for class 'dtplyr_step'
rename(.data, ...)

## S3 method for class 'dtplyr_step'
rename_with(.data, .fn, .cols = everything(), ...)
```

**Arguments**

|                    |   |
|--------------------|---|
| <code>.data</code> | A <code>lazy_dt()</code>  |
| <code>...</code>   | For <code>rename()</code> : <code>&lt;tidy-select&gt;</code> Use <code>new_name = old_name</code> to rename selected variables.<br>For <code>rename_with()</code> : additional arguments passed onto <code>.fn</code> . |
| <code>.fn</code>   | A function used to transform the selected <code>.cols</code> . Should return a character vector the same length as the input.   |
| <code>.cols</code> | <code>&lt;tidy-select&gt;</code> Columns to rename; defaults to all columns.  |

**Examples**

```
library(dplyr, warn.conflicts = FALSE)
dt <- lazy_dt(data.frame(x = 1, y = 2, z = 3))
dt %>% rename(new_x = x, new_y = y)
dt %>% rename_with(toupper)
```

---

```
replace_na.dtplyr_step
```

*Replace NAs with specified values*

---

### Description

This is a method for the tidyr `replace_na()` generic. It is translated to `data.table::fcoalesce()`.

Note that unlike `tidyr::replace_na()`, `data.table::fcoalesce()` cannot replace NULL values in lists.

### Usage

```
## S3 method for class 'dtplyr_step'
replace_na(data, replace = list())
```

### Arguments

|                      |  |
|----------------------|--|
| <code>data</code>    | A <code>lazy_dt()</code> .   |
| <code>replace</code> | If <code>data</code> is a data frame, <code>replace</code> takes a named list of values, with one value for each column that has missing values to be replaced. Each value in <code>replace</code> will be cast to the type of the column in <code>data</code> that it being used as a replacement in. If <code>data</code> is a vector, <code>replace</code> takes a single value. This single value replaces all of the missing values in the vector. <code>replace</code> will be cast to the type of <code>data</code> . |

### Examples

```
library(tidyr)

# Replace NAs in a data frame
dt <- lazy_dt(tibble(x = c(1, 2, NA), y = c("a", NA, "b")))
dt %>% replace_na(list(x = 0, y = "unknown"))

# Replace NAs using `dplyr::mutate()`
dt %>% dplyr::mutate(x = replace_na(x, 0))
```

---

```
select.dtplyr_step
```

*Subset columns using their names*

---

### Description

This is a method for the dplyr `select()` generic. It is translated to the `j` argument of `[.data.table]`.

### Usage

```
## S3 method for class 'dtplyr_step'
select(.data, ...)
```

**Arguments**

`.data` A `lazy_dt()`.

`...` `<tidy-select>` One or more unquoted expressions separated by commas. Variable names can be used as if they were positions in the data frame, so expressions like `x:y` can be used to select a range of variables.

**Examples**

```
library(dplyr, warn.conflicts = FALSE)

dt <- lazy_dt(data.frame(x1 = 1, x2 = 2, y1 = 3, y2 = 4))

dt %>% select(starts_with("x"))
dt %>% select(ends_with("2"))
dt %>% select(z1 = x1, z2 = x2)
```

---

`separate.dplyr_step` *Separate a character column into multiple columns with a regular expression or numeric locations*

---

**Description**

This is a method for the `tidyr::separate()` generic. It is translated to `data.table::tstrsplit()` in the `j` argument of `[.data.table]`.

**Usage**

```
## S3 method for class 'dplyr_step'
separate(
  data,
  col,
  into,
  sep = "[^[:alnum:]]+",
  remove = TRUE,
  convert = FALSE,
  ...
)
```

**Arguments**

`data` A `lazy_dt()`.

`col` Column name or position.  
This argument is passed by expression and supports quasiquotation (you can unquote column names or column positions).

`into` Names of new variables to create as character vector. Use `NA` to omit the variable in the output.

|         |   |
|---------|---|
| sep     | Separator between columns. The default value is a regular expression that matches any sequence of non-alphanumeric values.  |
| remove  | If TRUE, remove the input column from the output data frame.  |
| convert | If TRUE, will run <code>type.convert()</code> with <code>as.is = TRUE</code> on new columns. This is useful if the component columns are integer, numeric or logical.<br>NB: this will cause string "NA"s to be converted to NAs. |
| ...     | Arguments passed on to methods  |

### Examples

```
library(tidyr)
# If you want to split by any non-alphanumeric value (the default):
df <- lazy_dt(data.frame(x = c(NA, "x.y", "x.z", "y.z")), "DT")
df %>% separate(x, c("A", "B"))

# If you just want the second variable:
df %>% separate(x, c(NA, "B"))

# Use regular expressions to separate on multiple characters:
df <- lazy_dt(data.frame(x = c(NA, "x?y", "x.z", "y:z")), "DT")
df %>% separate(x, c("A", "B"), sep = "[.:?]")

# convert = TRUE detects column classes:
df <- lazy_dt(data.frame(x = c("x:1", "x:2", "y:4", "z", NA)), "DT")
df %>% separate(x, c("key", "value"), ":") %>% str
df %>% separate(x, c("key", "value"), ":", convert = TRUE) %>% str
```

---

slice.dplyr\_step      *Subset rows using their positions*

---

### Description

These are methods for the dplyr `slice()`, `slice_head()`, `slice_tail()`, `slice_min()`, `slice_max()` and `slice_sample()` generics. They are translated to the `i` argument of `[.data.table]`.

Unlike dplyr, `slice()` (and `slice()` alone) returns the same number of rows per group, regardless of whether or not the indices appear in each group.

### Usage

```
## S3 method for class 'dplyr_step'
slice(.data, ..., .by = NULL)

## S3 method for class 'dplyr_step'
slice_head(.data, ..., n, prop, by = NULL)

## S3 method for class 'dplyr_step'
slice_tail(.data, ..., n, prop, by = NULL)
```

```
## S3 method for class 'dtplyr_step'
slice_min(.data, order_by, ..., n, prop, by = NULL, with_ties = TRUE)

## S3 method for class 'dtplyr_step'
slice_max(.data, order_by, ..., n, prop, by = NULL, with_ties = TRUE)
```

## Arguments

|                        |   |
|------------------------|---|
| <code>.data</code>     | A <code>lazy_dt()</code> .  |
| <code>...</code>       | For <code>slice()</code> : <code>&lt;data-masking&gt;</code> Integer row values.<br>Provide either positive values to keep, or negative values to drop. The values provided must be either all positive or all negative. Indices beyond the number of rows in the input are silently ignored.<br>For <code>slice_*</code> (), these arguments are passed on to methods.   |
| <code>.by, by</code>   | <b>[Experimental]</b><br><code>&lt;tidy-select&gt;</code> Optionally, a selection of columns to group by for just this operation, functioning as an alternative to <code>group_by()</code> . For details and examples, see <code>?dplyr_by</code> .   |
| <code>n, prop</code>   | Provide either <code>n</code> , the number of rows, or <code>prop</code> , the proportion of rows to select. If neither are supplied, <code>n = 1</code> will be used. If <code>n</code> is greater than the number of rows in the group (or <code>prop &gt; 1</code> ), the result will be silently truncated to the group size. <code>prop</code> will be rounded towards zero to generate an integer number of rows.<br>A negative value of <code>n</code> or <code>prop</code> will be subtracted from the group size. For example, <code>n = -2</code> with a group of 5 rows will select $5 - 2 = 3$ rows; <code>prop = -0.25</code> with 8 rows will select $8 * (1 - 0.25) = 6$ rows. |
| <code>order_by</code>  | <code>&lt;data-masking&gt;</code> Variable or function of variables to order by. To order by multiple variables, wrap them in a data frame or tibble.   |
| <code>with_ties</code> | Should ties be kept together? The default, <code>TRUE</code> , may return more rows than you request. Use <code>FALSE</code> to ignore ties, and return the first <code>n</code> rows.  |

## Examples

```
library(dplyr, warn.conflicts = FALSE)

dt <- lazy_dt(mtcars)
dt %>% slice(1, 5, 10)
dt %>% slice(-(1:4))

# First and last rows based on existing order
dt %>% slice_head(n = 5)
dt %>% slice_tail(n = 5)

# Rows with minimum and maximum values of a variable
dt %>% slice_min(mpg, n = 5)
dt %>% slice_max(mpg, n = 5)

# slice_min() and slice_max() may return more rows than requested
```

```

# in the presence of ties. Use with_ties = FALSE to suppress
dt %>% slice_min(cyl, n = 1)
dt %>% slice_min(cyl, n = 1, with_ties = FALSE)

# slice_sample() allows you to random select with or without replacement
dt %>% slice_sample(n = 5)
dt %>% slice_sample(n = 5, replace = TRUE)

# you can optionally weight by a variable - this code weights by the
# physical weight of the cars, so heavy cars are more likely to get
# selected
dt %>% slice_sample(weight_by = wt, n = 5)

```

---

summarise.dtplyr\_step *Summarise each group to one row*

---

## Description

This is a method for the dplyr `summarise()` generic. It is translated to the `j` argument of `[.data.table]`.

## Usage

```

## S3 method for class 'dtplyr_step'
summarise(.data, ..., .by = NULL, .groups = NULL)

```

## Arguments

|                      |   |
|----------------------|---|
| <code>.data</code>   | A <code>lazy_dt()</code> .  |
| <code>...</code>     | <p>&lt;data-masking&gt; Name-value pairs of summary functions. The name will be the name of the variable in the result.</p> <p>The value can be:</p> <ul style="list-style-type: none"> <li>• A vector of length 1, e.g. <code>min(x)</code>, <code>n()</code>, or <code>sum(is.na(y))</code>.</li> <li>• A data frame, to add multiple columns from a single expression.</li> </ul> <p><b>[Deprecated]</b> Returning values with size 0 or &gt;1 was deprecated as of 1.1.0. Please use <code>reframe()</code> for this instead.</p> |
| <code>.by</code>     | <p><b>[Experimental]</b></p> <p>&lt;tidy-select&gt; Optionally, a selection of columns to group by for just this operation, functioning as an alternative to <code>group_by()</code>. For details and examples, see <code>?dplyr_by</code>.</p>   |
| <code>.groups</code> | <p><b>[Experimental]</b> Grouping structure of the result.</p> <ul style="list-style-type: none"> <li>• "drop_last": dropping the last level of grouping. This was the only supported option before version 1.0.0.</li> <li>• "drop": All levels of grouping are dropped.</li> <li>• "keep": Same grouping structure as <code>.data</code>.</li> <li>• "rowwise": Each row is its own group.</li> </ul>   |

When `.groups` is not specified, it is chosen based on the number of rows of the results:

- If all the results have 1 row, you get "drop\_last".
- If the number of rows varies, you get "keep" (note that returning a variable number of rows was deprecated in favor of `reframe()`, which also unconditionally drops all levels of grouping).

In addition, a message informs you of that choice, unless the result is ungrouped, the option "dplyr.summarise.inform" is set to FALSE, or when `summarise()` is called from a function in a package.

### Examples

```
library(dplyr, warn.conflicts = FALSE)

dt <- lazy_dt(mtcars)

dt %>%
  group_by(cyl) %>%
  summarise(vs = mean(vs))

dt %>%
  group_by(cyl) %>%
  summarise(across(dispatch, mean))
```

---

transmute.dtplyr\_step *Create new columns, dropping old*

---

### Description

This is a method for the dplyr `transmute()` generic. It is translated to the `j` argument of `[.data.table]`.

### Usage

```
## S3 method for class 'dtplyr_step'
transmute(.data, ...)
```

### Arguments

|                    |   |
|--------------------|---|
| <code>.data</code> | A <code>lazy_dt()</code> .  |
| <code>...</code>   | <data-masking> Name-value pairs. The name gives the name of the column in the output.<br>The value can be: <ul style="list-style-type: none"> <li>• A vector of length 1, which will be recycled to the correct length.</li> <li>• A vector the same length as the current group (or the whole data frame if ungrouped).</li> <li>• NULL, to remove the column.</li> <li>• A data frame or tibble, to create multiple columns in the output.</li> </ul> |



**Examples**

```
library(dplyr, warn.conflicts = FALSE)

dt <- lazy_dt(dplyr::starwars)
dt %>% transmute(name, sh = paste0(species, "/", homeworld))
```

---

unite.dplyr\_step      *Unite multiple columns into one by pasting strings together.*

---

**Description**

This is a method for the tidy `unite()` generic.

**Usage**

```
## S3 method for class 'dplyr_step'
unite(data, col, ..., sep = "_", remove = TRUE, na.rm = FALSE)
```

**Arguments**

|                     |  |
|---------------------|--|
| <code>data</code>   | A data frame.  |
| <code>col</code>    | The name of the new column, as a string or symbol.<br>This argument is passed by expression and supports <a href="#">quasiquotation</a> (you can unquote strings and symbols). The name is captured from the expression with <a href="#">rlang::ensym()</a> (note that this kind of interface where symbols do not represent actual objects is now discouraged in the tidyverse; we support it here for backward compatibility). |
| <code>...</code>    | <a href="#">&lt;tidy-select&gt;</a> Columns to unite   |
| <code>sep</code>    | Separator to use between values.   |
| <code>remove</code> | If TRUE, remove input columns from output data frame.  |
| <code>na.rm</code>  | If TRUE, missing values will be removed prior to uniting each value.   |

**Examples**

```
library(tidy)

df <- lazy_dt(expand_grid(x = c("a", NA), y = c("b", NA)))
df

df %>% unite("z", x:y, remove = FALSE)

# Separate is almost the complement of unite
df %>%
  unite("xy", x:y) %>%
  separate(xy, c("x", "y"))
# (but note `x` and `y` contain now "NA" not NA)
```

# Index

?dplyr\_by, [11](#), [18](#), [25](#), [30](#), [31](#)  
?join\_by, [17](#)

anti\_join(), [16](#)  
arrange(), [2](#), [11](#)  
arrange.dtplyr\_step, [2](#)  
as.data.frame(), [15](#)  
as.data.frame.dtplyr\_step  
    (collect.dtplyr\_step), [3](#)  
as.data.table.dtplyr\_step  
    (collect.dtplyr\_step), [3](#)  
as\_tibble.dtplyr\_step  
    (collect.dtplyr\_step), [3](#)

collect(), [15](#)  
collect.dtplyr\_step, [3](#)  
complete(), [5](#)  
complete.dtplyr\_step, [4](#)  
compute.dtplyr\_step  
    (collect.dtplyr\_step), [3](#)  
count(), [5](#)  
count.dtplyr\_step, [5](#)  
cross\_join(), [17](#)

data.table::as.data.table(), [15](#)  
data.table::CJ(), [8](#)  
data.table::dcast(), [22](#)  
data.table::fcoalesce(), [27](#)  
data.table::fintersect(), [14](#)  
data.table::fsetdiff(), [14](#)  
data.table::funion(), [14](#)  
data.table::melt(), [20](#)  
data.table::merge.data.table(), [16](#)  
data.table::nafill(), [9](#)  
data.table::setcolorder(), [16](#), [18](#)  
data.table::setkey(), [16](#)  
data.table::setnames(), [16](#), [26](#)  
data.table::tstrsplit(), [28](#)  
data.table::unique.data.table(), [6](#)  
desc(), [3](#)

distinct(), [6](#)  
distinct.dtplyr\_step, [6](#)  
drop\_na.dtplyr\_step, [7](#)

expand(), [5](#)  
expand.dtplyr\_step, [8](#)  
extract(), [21](#)

fill.dtplyr\_step, [9](#)  
filter.dtplyr\_step, [11](#)  
full\_join(), [16](#)

group\_by(), [11](#), [12](#), [18](#), [25](#), [30](#), [31](#)  
group\_by.dtplyr\_step, [12](#)  
group\_map(), [13](#)  
group\_map.dtplyr\_step  
    (group\_modify.dtplyr\_step), [13](#)  
group\_modify(), [13](#)  
group\_modify.dtplyr\_step, [13](#)  
grouped\_dt(lazy\_dt), [15](#)

head(), [14](#)  
head.dtplyr\_step, [14](#)

inner\_join(), [16](#)  
intersect(), [14](#)  
intersect.dtplyr\_step, [14](#)

join\_by(), [17](#)

lazy\_dt, [4](#), [13](#), [15](#)  
lazy\_dt(), [3](#), [5–8](#), [11–15](#), [17–20](#), [23–28](#),  
    [30–32](#)  
left\_join(), [16](#)  
left\_join.dtplyr\_step, [16](#)

mutate(), [18](#)  
mutate.dtplyr\_step, [18](#)

nest.dtplyr\_step, [19](#)

order(), [2](#)

`pivot_longer.dtplyr_step`, 20  
`pivot_wider.dtplyr_step`, 22  
`pull()`, 15

quasiquote, 33

`reframe()`, 24, 31, 32  
`reframe.dtplyr_step`, 24  
`relocate()`, 19, 25  
`relocate.dtplyr_step`, 25  
`rename()`, 26  
`rename.dtplyr_step`, 26  
`rename_with()`, 26  
`rename_with.dtplyr_step`  
    (`rename.dtplyr_step`), 26  
`replace_na.dtplyr_step`, 27  
`right_join()`, 16  
`rlang::as_function()`, 8  
`rlang::ensym()`, 33

`select()`, 16, 27  
`select.dtplyr_step`, 27  
`semi_join()`, 16  
`separate()`, 21  
`separate.dtplyr_step`, 28  
`setdiff()`, 14  
`setdiff.dtplyr_step`  
    (`intersect.dtplyr_step`), 14  
`slice()`, 29  
`slice.dtplyr_step`, 29  
`slice_head.dtplyr_step`  
    (`slice.dtplyr_step`), 29  
`slice_max.dtplyr_step`  
    (`slice.dtplyr_step`), 29  
`slice_min.dtplyr_step`  
    (`slice.dtplyr_step`), 29  
`slice_tail.dtplyr_step`  
    (`slice.dtplyr_step`), 29  
`summarise()`, 31  
`summarise.dtplyr_step`, 31

`tail()`, 14  
`tail.dtplyr_step` (`head.dtplyr_step`), 14  
`tbl()`, 12  
`tbl_dt` (`lazy_dt`), 15  
`tibble::as_tibble()`, 15  
`tidyr::nest()`, 19  
`tidyr::separate()`, 28  
`transmute()`, 32  
`transmute.dtplyr_step`, 32

`ungroup()`, 12  
`ungroup.dtplyr_step`  
    (`group_by.dtplyr_step`), 12  
`union()`, 14  
`union.dtplyr_step`  
    (`intersect.dtplyr_step`), 14  
`union_all()`, 14  
`union_all.dtplyr_step`  
    (`intersect.dtplyr_step`), 14  
`unite.dtplyr_step`, 33

`vctrs::vec_as_names()`, 8, 21, 23