

# Tidy Verbs for Fast Data Manipulation: : CHEAT SHEET



## Basics

Combining the merits of syntax elegance from **dplyr** and computing performance from **data.table**, **tidyfst** intends to provide users with state-of-the-art data manipulation tools with least pain. The **data.table** syntax (**dt[i, j, by]**) could be applied in tidyfst, while the tidy data principal is implemented everywhere.



Each **variable** is saved in its own **column**

Each **observation** is saved in its own **row**

**x %>% f(y)** becomes **f(x, y)**



## Data I/O

**CSV(considered to be general on various platforms):**

**fread("file.csv")** – read data from a flat file such as .csv or .tsv into R.

**fwrite(dt, "file.csv")** – write data to a flat file from R.

**FST(considered to be fast and memory efficient):**

**import\_fst("file.fst")** – read data from a flat file of .fst into R.

**export\_fst(dt, "file.fst")** – write data to a .fst file from R.

## Dealing with NAs

**drop\_na\_dt** – drop entries with NA in the column(s)

**replace\_na\_dt** – replace NA with other value

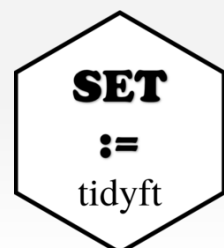
**delete\_na\_cols** – delete column(s) when NA volume or proportion is larger than a threshold

**delete\_na\_rows** – delete row(s) when NA volume or proportion is larger than a threshold

**fill\_na\_dt** – fill NA with previous or next observations

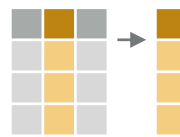
## tidyft

In **data.table**, there is an important feature: modification by reference. This could be really useful in high performance computation. This feature is implemented in tidyft, the mirror package of tidyfst.



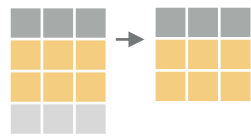
## Subset data

### BY COLUMN



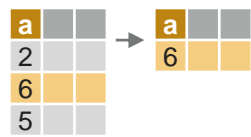
**select\_dt(dt, 2)** – get column(s) as data.table  
**pull\_dt(dt, 2)** – get column as a vector  
**select\_mix** – select columns flexibly

### BY ROW



Subset row based on:

- **slice\_dt** – Position
- **slice\_sample** – Randomly
- **slice\_max\_dt/slice\_min\_dt** – values
- **slice\_top\_dt/slice\_tail\_dt** – Position from top or bottom
- **distinct\_dt** – unique values
- **filter\_dt** – condition



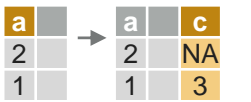
### LOGICAL OPERATORS TO USE IN filter\_dt

<	<=	is.na()	%in%		%like%
>	>=	!is.na()	!	&	%between%

## Update data

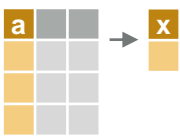


**mutate\_dt** – add or mutate column(s)  
**mutate\_vars** – update multiple columns



**mutate\_when** – update entries that meet certain condition

## Aggregation

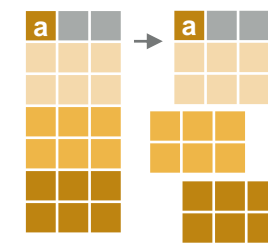


**summarise\_dt** – aggregate one column to one value  
**count\_dt/add\_count\_dt** – get unique counts  
**summarise\_vars** – aggregate multiple columns  
**summarise\_when** – conditional aggregation

## Reorder data

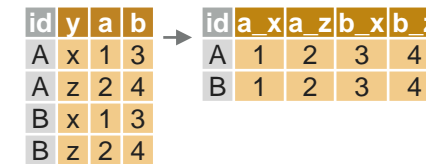
**arrange\_dt** – reorder data by row according to column value  
**relocate\_dt** – reorder data by column(s)

## Group computation

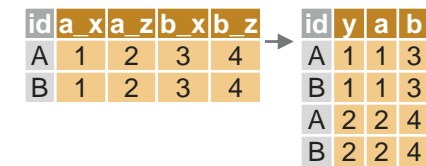


While tidyfst provides **group\_dt/group\_by\_dt/group\_by\_exe** to implement computation in groups. If you are handling big data, it is recommended to use the **by** parameter in the function when available, e.g. `summarise_dt(iris, avg = mean(Sepal.Length), by = Species)`

## Reshape data



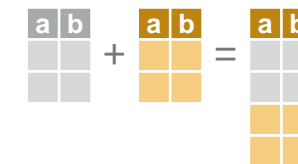
**wider\_dt** – Transform data from long to wide



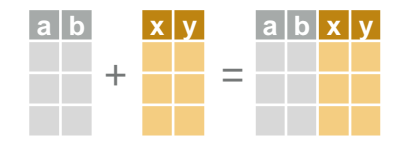
**longer\_dt** – Transform data from wide to long

## Merge data

### BIND

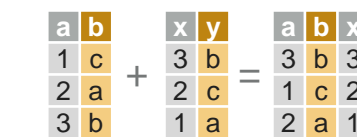


**rbind**



**cbind**

### JOIN



**left\_join\_dt / right\_join\_dt**  
**inner\_join\_dt / full\_join\_dt**  
**anti\_join\_dt / semi\_join\_dt**

### SET OPERATIONS



**intersect\_dt**  
**union\_dt**  
**setdiff\_dt**  
**setequal\_dt**