

Demographic (donut) plots

Overview

This script creates donut plots that show survey participation rates across various demographics. In this script, I save various data tables to files, and then I read these data files in my notebook figure2.qmd to create the final figure, where all these donut plots are stitched together.

Import packages and utilities

```
project_root <- here::here() # requires that you be somewhere in the
# project directory (not above it)
# packages
suppressMessages(source(file.path(project_root, "scripts/packages.R")))
# functions and objects used across scripts
suppressMessages(source(file.path(project_root, "scripts/utils.R")))
```

Define functions

create_df_for_plotting

- Arguments:
 - data: The data frame with the survey data. This function will count rows to determine the count for each category, so don't tally up the group counts in advance. It will also remove rows where this column is an empty string, so no need to clean the df at all.
 - column: The name of the column with the data to be plotted.
- Details:

- This function creates a new data frame suitable for plotting as a donut plot. It essentially creates a single stacked bar with all our data, and then plots that on a polar coordinate system to make it a donut. Method from <https://r-graph-gallery.com/128-ring-or-donut-plot.html>. This function extracts the data of interest from the larger data frame, and puts them into a new data frame along with relevant breakpoints for and label locations.

- Outputs:

- `long_data`: a new data frame with columns values, Freq, fraction, ymax, ymin, labelPosition, and label.

```
create_df_for_plotting <- function(data, column) {  
  if (!column %in% names(data)) {  
    stop("Column not found in data frame")  
  }  
  # Extract specified column and remove rows where this column  
  # is an empty string, indicating the participant didn't answer  
  values <- data[[column]][nzchar(data[[column]])]  
  # Count occurrences of each unique value  
  values_table <- table(values)  
  
  # Convert to data frame and compute fractions  
  long_data <- as.data.frame(values_table) %>%  
    mutate(fraction = Freq / sum(Freq)) %>%  
    arrange(desc(fraction))  
  
  # Compute the cumulative percentages (top of each rectangle)  
  long_data$ymax <- cumsum(long_data$fraction)  
  long_data$ymin <- c(0, head(long_data$ymax, n = -1))  
  
  # Compute label position  
  long_data$labelPosition <- (long_data$ymax + long_data$ymin) / 2  
  
  # Create label column  
  long_data$label <- paste0(long_data$Freq)  
  
  # Make the text wrap so the legend is more compact  
  long_data <- long_data %>%  
    mutate(values = str_wrap(values, width = 20))  
  
  # from scripts/utils.R  
  long_data <- reorder_factor_by_column(  

```

```

    long_data,
    values,
    Freq
  )

  return(long_data)
}

```

donut_chart

- Arguments:
 - df: The data frame ready for plotting, such as the one produced by create_df_for_plotting.
- Details:
 - This function creates a donut plot.
- Outputs:
 - A ggplot object.

```

donut_chart <- function(
  df,
  cpalette = COLORS, # default: colors from utils.R
  legendpos = "right",
  title_bottom_margin = 15,
  legend_top_margin = 0,
  label_dist = 4.3, # controls how close the labels are to the ring
  white_labels = FALSE
) {

  if (white_labels) {
    label_x <- 3.5
    label_size <- 8
    label_color <- "white"
  } else {
    label_x <- label_dist
    label_size <- 9
    label_color <- "black"
  }

  ggplot(

```

```

df,
aes(
  ymax = ymax,
  ymin = ymin,
  xmax = 4,
  xmin = 3,
  fill = values
)
) +
geom_rect() +
# Add labels
geom_text(
  x = label_x,
  aes(y = labelPosition, label = label),
  size = label_size,
  color = label_color
) +
scale_fill_manual(values = cpalette) +
theme_void() +
coord_polar(theta = "y", clip = "off") +
xlim(c(2, label_x + 0.05)) +
theme(
  legend.text = element_text(size = 18),
  legend.title = element_blank(),
  legend.position = legendpos,
  legend.key.spacing.y = unit(0.5, "lines"),
  legend.margin = margin(t = unit(legend_top_margin, "lines")),
  plot.title = element_text(
    hjust = 0.5,
    size = 28,
    margin = margin(b = title_bottom_margin)
  ),
  plot.margin = margin(t = 20, r = 5, b = 10, l = 5),
  plot.background = element_rect(fill = "white", color = "white")
)
}

```

Load data

```
other_quant <- load_qualtrics_data("clean_data/other_quant.tsv")
status <- load_qualtrics_data("clean_data/contributor_status_Q3.tsv")

raw_data <- cbind(status, other_quant)
nrow(raw_data)
```

```
[1] 332
```

```
head(raw_data)
```

	Past	Future	campus	favorite_solution	field_of_study
1	True	True	UC Santa Barbara	Sustainability grants	Math and CS
2	True	True	UC Santa Barbara	Containerization	Life sciences
3	True	True	UC Santa Barbara	Computing environments	Humanities
4	True	True	UC Santa Barbara	Sustainability grants	Math and CS
5	True	True	UC Santa Barbara	Documentation help	Life sciences
6	False	True	UC Santa Barbara		Math and CS

	job_category	staff_categories
1	Faculty	
2	Post-Doc	
3	Other research staff	
4	Faculty	
5	Faculty	
6	Other research staff	

Wrangle data

First, filter out people who were neither past nor future contributors. We'll use this filtered data frame moving forward.

```
# Filter duds
data <- raw_data %>%
  filter(!(Past == "" | Future == "")) %>%
  filter(!(Past == "False" & Future == "False"))
```

Record experienced vs. aspiring, which we'll plot at the end of this notebook.

```
status_clean <- data %>%
  mutate(status = if_else(Past == "True", "Experienced", "Aspiring"))
```

Create a “population” variable. This will make it easier to change all the plot titles if we decide to redo these plots for just the experienced contributors.

```
population <- "Respondents"
```

Donut charts of participation by groups

This script creates 4 plots, called p1-p4.

Plot #1 is of campus.

```
campus_data <- create_df_for_plotting(data, "campus")
campus_data <- campus_data %>%
  filter(
    !startsWith(as.character(values), "I'm not affiliat")
  )

p1 <- donut_chart(campus_data) +
  labs(title = sprintf("Campus of %s", population))
```

Save data to file for fine-tuning of plot.

```
cd <- campus_data
names(cd)[1] <- "categories"
names(cd)[2] <- "n_responses"
write_df_to_file(
  cd,
  file.path("data_for_plots/fig_campus_donut.tsv")
)
```

Plot #2 is academics' field of study. Qualitative answers are not included in this plot.

```
field_data <- create_df_for_plotting(data, "field_of_study")
p2 <- donut_chart(field_data) +
  labs(title = sprintf("Academic %s' Fields of Study", population))
```

Save data to file for fine-tuning of plot.

```
fd <- field_data
names(fd)[1] <- "categories"
names(fd)[2] <- "n_responses"
write_df_to_file(
  fd,
  file.path("data_for_plots/fig_fields_donut.tsv")
)
```

Plot #3 is of all respondents' job categories.

```
job_data <- create_df_for_plotting(data, "job_category")

p3 <- donut_chart(job_data) +
  labs(title = sprintf("%s' Job Categories", population))
```

Save the count data to a file for fine-tuning of plot. Since this plot is more complicated, I am not pre-calculating the label positions; I am just saving the simple count data.

```
jd <- job_data %>% select(c(values, Freq))
names(jd) <- c("categories", "n_responses")
write_df_to_file(
  jd,
  file.path("data_for_plots/fig_job_donut.tsv")
)
```

Plot #4 is of staff categories (IT, Research Support, etc.). This one is more complicated. I'm not using my `create_df_for_plotting` function for this group because I want to combine the jobs that have only 1 or 2 responses into the existing "Other" category. So the code is similar but not the same.

```
staff_data <- data[["staff_categories"]][nzchar(data[["staff_categories"]])]
# Count occurrences of each unique value
staff_data <- as.data.frame(table(staff_data))
names(staff_data) <- c("job", "count")
staff_data$job <- as.character(staff_data$job)

staff_data_clean <- staff_data %>%
  mutate(job = if_else(count < 3, "Other", job)) %>% # relabel rare jobs as "Other"
  group_by(job) %>% # gather all "Other" rows together
```

```

    summarise(Freq = sum(count), .groups = "drop")

staff_long_data <- as.data.frame(staff_data_clean) %>%
  mutate(fraction = Freq / sum(Freq)) %>%
  arrange(desc(fraction))

# Compute the cumulative percentages (top of each rectangle)
staff_long_data$ymax <- cumsum(staff_long_data$fraction)
staff_long_data$ymin <- c(0, head(staff_long_data$ymax, n = -1))

# Compute label position
staff_long_data$labelPosition <- (staff_long_data$ymax + staff_long_data$ymin) /
  2

# Create label column
staff_long_data$label <- paste0(staff_long_data$Freq)

# Rename this one column to match the donut_chart function
names(staff_long_data)[names(staff_long_data) == "job"] <- "values"

# Wrap text
staff_long_data <- staff_long_data %>%
  mutate(values = str_wrap(values, width = 20))

# from scripts/utils.R
staff_long_data <- reorder_factor_by_column(
  staff_long_data,
  values,
  Freq
)

p4 <- donut_chart(staff_long_data) +
  labs(title = sprintf("Staff %s' Work Areas", population))

```

Save data to file for fine-tuning of plot.

```

sd <- staff_long_data
names(sd)[1] <- "categories"
names(sd)[2] <- "n_responses"
# I put newline characters in some of these categories so they would
# fit on the plot; remove them before saving to file
sd$categories <- sapply(sd$categories, function(x) gsub("\n", " ", x))

```



```
write_df_to_file(
  sd,
  file.path("data_for_plots/fig_staff_donut.tsv")
)
```

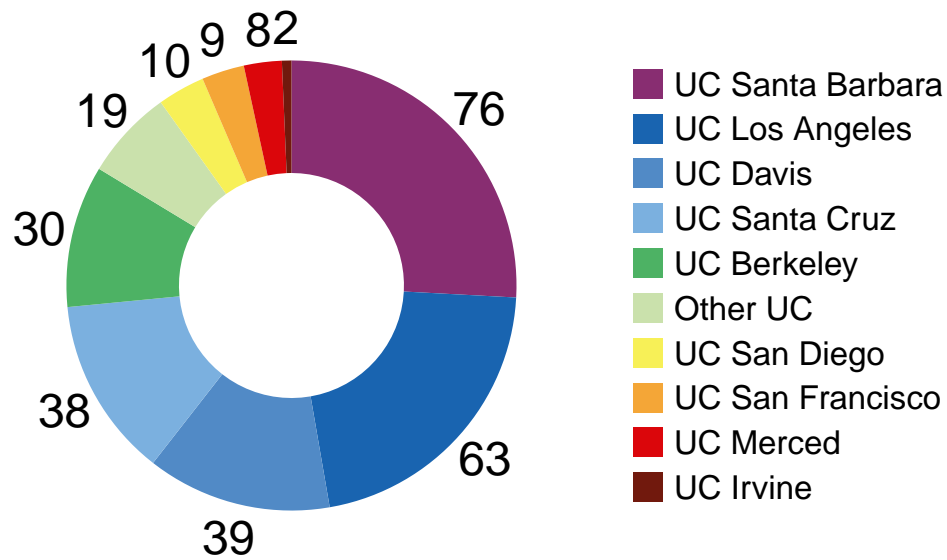
Refine plot aesthetics

I'm envisioning a figure where I'd like to have some variety in the color palettes and the legend position. I'll put three donuts onto one canvas in Powerpoint, and then make some additional tweaks: I'll add some lines to illustrate which slices of the donut I'm focusing on, and also I will need to fix this bug where the legend keys stretch to multiple lines. (<https://github.com/tidyverse/ggplot2/issues/3669>) (A fix has been merged, but it's not yet available in the latest release of ggplot2.)

Plot 1: campus

```
p1 <- donut_chart(
  campus_data,
  # Paul Tol's "rainbow" color scheme
  cpalette = c(
    "#882d71",
    "#1964b0",
    "#518ac6",
    "#7bb0df",
    "#4db264",
    "#cae1ac",
    "#f7f057",
    "#f4a637",
    "#db060b",
    "#71190d"
  )
) +
  labs(title = sprintf("Campuses of %s", population))
p1
```

Campuses of Respondents

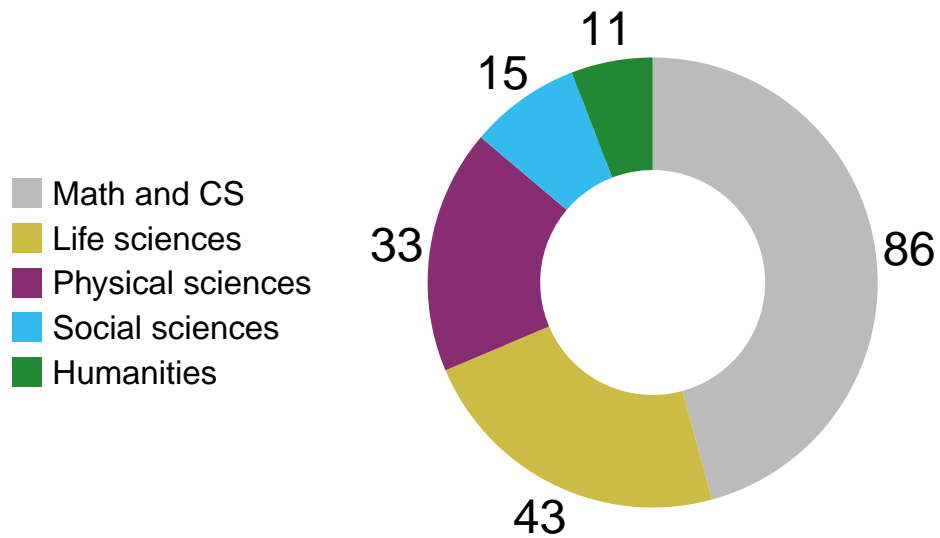


```
# Function from utils.R
save_plot(sprintf("donut_p1_%s.tiff", population), 10, 6, p = p1)
```

Plot 2: fields of study

```
p2 <- donut_chart(
  field_data,
  cpalette = c(
    '#BBBBBB',
    '#CCBB44',
    '#882d71',
    '#33BBEE',
    '#228833'
  ),
  legendpos="left"
) +
  labs(title = sprintf("Academic %s' Fields of Study", population))
p2
```

Academic Respondents' Fields of Study

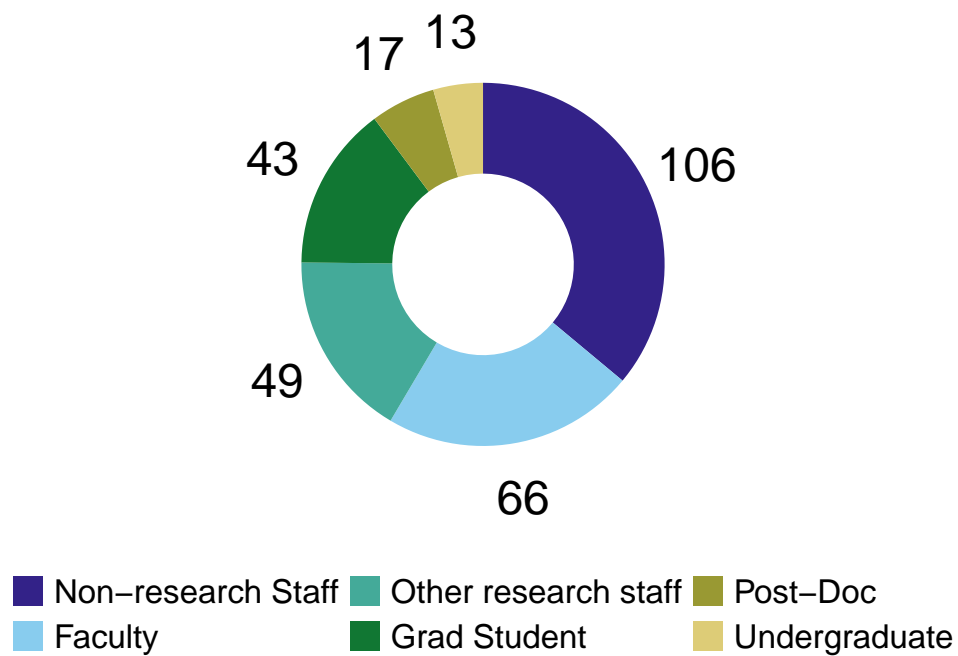


```
# Function from utils.R  
save_plot(sprintf("donut_p2_%s.tiff", population), 10, 6, p = p2)
```

Plot3: job categories

```
p3 <- donut_chart(  
  job_data,  
  legendpos = "bottom",  
  title_bottom_margin = 0,  
  legend_top_margin = 0,  
  label_dist = 4.6,  
  white_labels = FALSE  
) +  
  labs(title = sprintf("%s' Job Categories", population))  
)  
p3
```

Respondents' Job Categories



Save plot

```
# Function from utils.R
save_plot(
  sprintf("donut_p3_%s.tiff", tolower(population)),
  12,
  8,
  p = p3
)
```

Plot 4: non-research staff categories

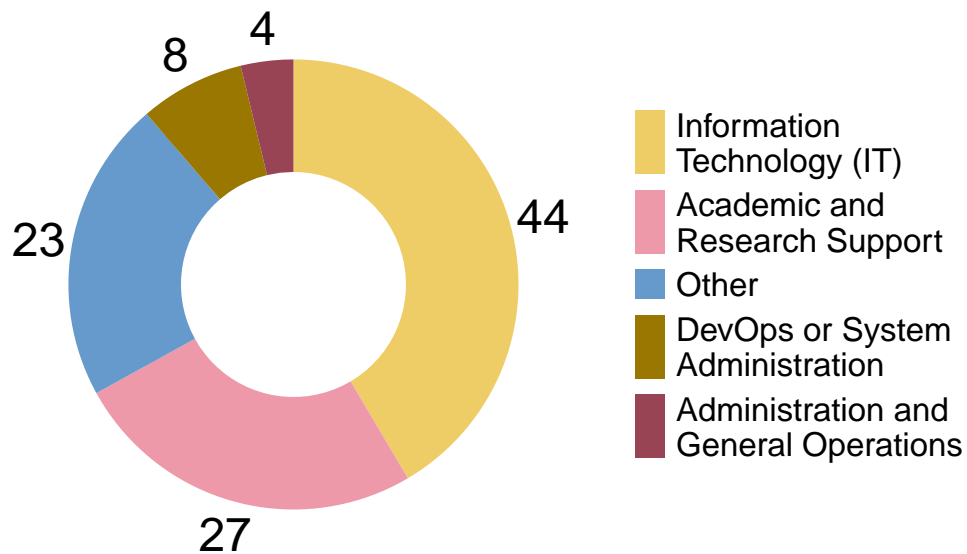
```
p4 <- donut_chart(
  staff_long_data,
  cpalette = c(
    '#EECC66',
    '#EE99AA',
    '#6699CC',
    '#997700',
  )
)
```

```

'#994455'
)
) +
  labs(title = sprintf("Staff %s' Work Areas", population))
p4

```

Staff Respondents' Work Areas



```

# Function from utils.R
save_plot(
  sprintf("donut_p4_%s.tiff", tolower(population)),
  10,
  6,
  p = p4
)

```

Experienced vs. aspiring contributors

Let's also make a simple donut plot with aspiring vs. experienced contributors.

```

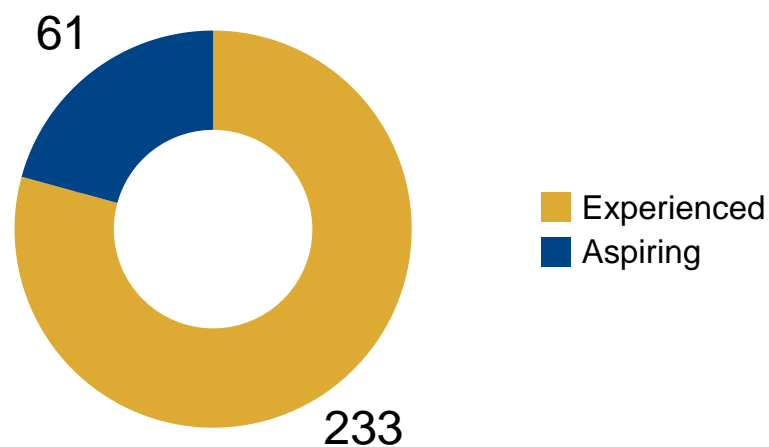
status_data <- create_df_for_plotting(status_clean, "status")

p5 <- donut_chart(
  status_data,
  label_dist = 4.5,
  cpalette = c(
    "#ddaa33",
    "#004488"
  ),
  title_bottom_margin = 0
) +
  labs(title = sprintf("%s' Statuses as\nOpen Source Contributors", population))

p5

```

Respondents' Statuses as Open Source Contributors



```

# Function from utils.R
save_plot(
  sprintf("donut_p5_%s.tiff", tolower(population)),
  10,

```

```
6,  
p = p5  
)
```

Save data to file for fine-tuning of plot.

```
statd <- status_data  
names(statd)[1] <- "categories"  
names(statd)[2] <- "n_responses"  
write_df_to_file(  
  statd,  
  file.path("data_for_plots/fig_status_donut.tsv")  
)
```

```
sessionInfo()
```

```
R version 4.4.2 (2024-10-31)  
Platform: aarch64-apple-darwin20  
Running under: macOS Sequoia 15.6.1
```

```
Matrix products: default
```

```
BLAS: /Library/Frameworks/R.framework/Versions/4.4-arm64/Resources/lib/libRblas.0.dylib
```

```
LAPACK: /Library/Frameworks/R.framework/Versions/4.4-arm64/Resources/lib/libRlapack.dylib;
```

```
locale:
```

```
[1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
```

```
time zone: America/Los_Angeles
```

```
tzcode source: internal
```

```
attached base packages:
```

```
[1] tools      grid      stats      graphics  grDevices  datasets  utils  
[8] methods    base
```

```
other attached packages:
```

```
[1] treemapify_2.5.6      tidyr_1.3.1           svglite_2.2.1  
[4] stringr_1.5.1         scales_1.4.0          readr_2.1.5  
[7] pwr_1.3-0             patchwork_1.3.2       ordinal_2023.12-4.1  
[10] lme4_1.1-37          Matrix_1.7-1          languageserver_0.3.16  
[13] here_1.0.1            gtools_3.9.5          ggforce_0.5.0  
[16] FSA_0.10.0            fpc_2.2-13            forcats_1.0.0
```

[19]	factoextra_1.0.7	ggplot2_3.5.2	emmeans_1.11.2
[22]	dplyr_1.1.4	corrplot_0.95	ComplexHeatmap_2.22.0
[25]	cluster_2.1.8.1	BiocManager_1.30.26	

loaded via a namespace (and not attached):

[1]	Rdpack_2.6.4	rlang_1.1.6	magrittr_2.0.3
[4]	clue_0.3-66	GetoptLong_1.0.5	matrixStats_1.5.0
[7]	compiler_4.4.2	flexmix_2.3-20	systemfonts_1.2.3
[10]	png_0.1-8	callr_3.7.6	vctrs_0.6.5
[13]	pkgconfig_2.0.3	shape_1.4.6.1	crayon_1.5.3
[16]	fastmap_1.2.0	labeling_0.4.3	rmarkdown_2.29
[19]	ggfittext_0.10.2	tzdb_0.5.0	ps_1.9.1
[22]	nloptr_2.2.1	purrr_1.1.0	xfun_0.53
[25]	modeltools_0.2-24	jsonlite_2.0.0	tweenr_2.0.3
[28]	parallel_4.4.2	prabclus_2.3-4	R6_2.6.1
[31]	stringi_1.8.7	RColorBrewer_1.1-3	boot_1.3-31
[34]	diptest_0.77-2	numDeriv_2016.8-1.1	estimability_1.5.1
[37]	Rcpp_1.1.0	iterators_1.0.14	knitr_1.50
[40]	IRanges_2.40.1	splines_4.4.2	nnet_7.3-19
[43]	tidyselect_1.2.1	yaml_2.3.10	doParallel_1.0.17
[46]	codetools_0.2-20	processx_3.8.6	lattice_0.22-6
[49]	tibble_3.3.0	withr_3.0.2	evaluate_1.0.4
[52]	polyclip_1.10-7	xml2_1.4.0	circlize_0.4.16
[55]	mclust_6.1.1	kernlab_0.9-33	pillar_1.11.0
[58]	renv_1.1.5	foreach_1.5.2	stats4_4.4.2
[61]	reformulas_0.4.1	generics_0.1.4	rprojroot_2.1.1
[64]	S4Vectors_0.44.0	hms_1.1.3	minqa_1.2.8
[67]	xtable_1.8-4	class_7.3-22	glue_1.8.0
[70]	robustbase_0.99-4-1	mvtnorm_1.3-3	rbibutils_2.3
[73]	colorspace_2.1-1	nlme_3.1-166	cli_3.6.5
[76]	textshaping_1.0.1	gtable_0.3.6	DEoptimR_1.1-4
[79]	digest_0.6.37	BiocGenerics_0.52.0	ucminf_1.2.2
[82]	ggrepel_0.9.6	rjson_0.2.23	farver_2.1.2
[85]	htmltools_0.5.8.1	lifecycle_1.0.4	GlobalOptions_0.1.2
[88]	MASS_7.3-61		