# Word Clouds?



## Text visualization tool

- **Words scaled** according to the a quantity (often a count or a proportion)
- Nice (**meaningless**) word **position**/order.

- For documents, vocabulary often pruned (stop words/rarely used words)

# Word Clouds?



## Text visualization tool

- **Words scaled** according to the a quantity (often a count or a proportion)
- Nice (**meaningless**) word **position**/order.

- For documents, vocabulary often pruned (stop words/rarely used words)

# Word Clouds in R





## `wordcloud`

- **base graphics**
- fast but crude placement algorithm in R

## `wordcloud2`

- `html` **graphics**
- placement algorithm in `js`
- lots of bell and whistles (arbitrary rotation, mask...)

# ggwordcloud



## A new ggplot2 geometry

- **ggplot2 ecosystem**
- a **new geometry** similar to `geom_text` and `geom_text_repel`
- fast placement algorithm in C
- more functionalities than `wordcloud` and `wordcloud2`...

# geom_text_wordcloud

```
library(ggwordcloud)
data("love_words_small")
set.seed(42)
ggplot(love_words_small,
       aes(label = word,
           size = speakers)) +
  geom_text_wordcloud() +
  scale_size_area(max_size = 24)
  theme_minimal()
```



### A dedicated text geometry

- geometry with the `geom_text` **syntax**:
  - `label` for the word
  - `size` for the count
- automatic placement without overlapping around a default $(0, 0)$ position

# geom_text_wordcloud

```
library(ggwordcloud)
data("love_words_small")
set.seed(42)
ggplot(love_words_small,
       aes(label = word,
           size = speakers)) +
  geom_text_wordcloud() +
  scale_size_area(max_size = 24)
  theme_minimal()
```



## A dedicated text geometry

- geometry with the `geom_text` **syntax**:
    - `label` for the word
    - `size` for the count
- automatic placement without overlapping around a default $(0, 0)$ position

# Text Scaling





**geom_text_wordcloud**

- **Font size proportional to** (the square root of) the count/proportion.
- Ink area depends on the number of letters...

**geom_text_wordcloud_area**

- **Ink area proportional to** the count/proportion.
- Perception not biased by number of letters.

## Rotation

- **Arbitrary rotation** with `angle` aesthetic.



## Facet

- Compatible with `ggplot2` **faceting system**.

# Shapes and Mask



square      triangle-forward      star

## Shapes

- Cloud may have **different base shapes**.



## Mask

- Words stay within a **prescribed mask**.

- Other functionalities: color, position, fonts…

# Algorithm



## Algorithm

- Compute text area using `textGrob` and deduce font size
- Draw text using again `textGrob` and compute a mask made of small bounding boxes
- Use a fast spiraling placement algorithm (in `C++`) to place the words without any boxes overlap

- **Bottleneck:** text size computation!

# Algorithm



## Algorithm

- Compute text area using `textGrob` and deduce font size
- Draw text using again `textGrob` and compute a mask made of small bounding boxes
- Use a fast spiraling placement algorithm (in `C++`) to place the words without any boxes overlap

- **Bottleneck:** text size computation!

# Word Zones



## Interest of a Grammar

- Can produce graphs which where not planned!

# Thank you to the R ecosystem



## Code and algorithm

- `ggplot2`: environment
- `ggrepel`: basis of the code
- `wordcloud/wordcloud2`: source of inspiration

## Environment

- `rstudio`...
- `usethis`: package skeleton
- `testthat`: unit testing
- `devtools/rhub`: package dev/testing before `CRAN`
- `pkgdown`: documentation

# ggwordcloud



## A word cloud geometry for `ggplot2`

- A lovely, easy to use and full of functionalities package.
- Available on `CRAN`.
- Source and bug tracker at
  `https://github.com/lepennec/ggwordcloud`
- **Website:** `https://lepennec.github.io/ggwordcloud/`