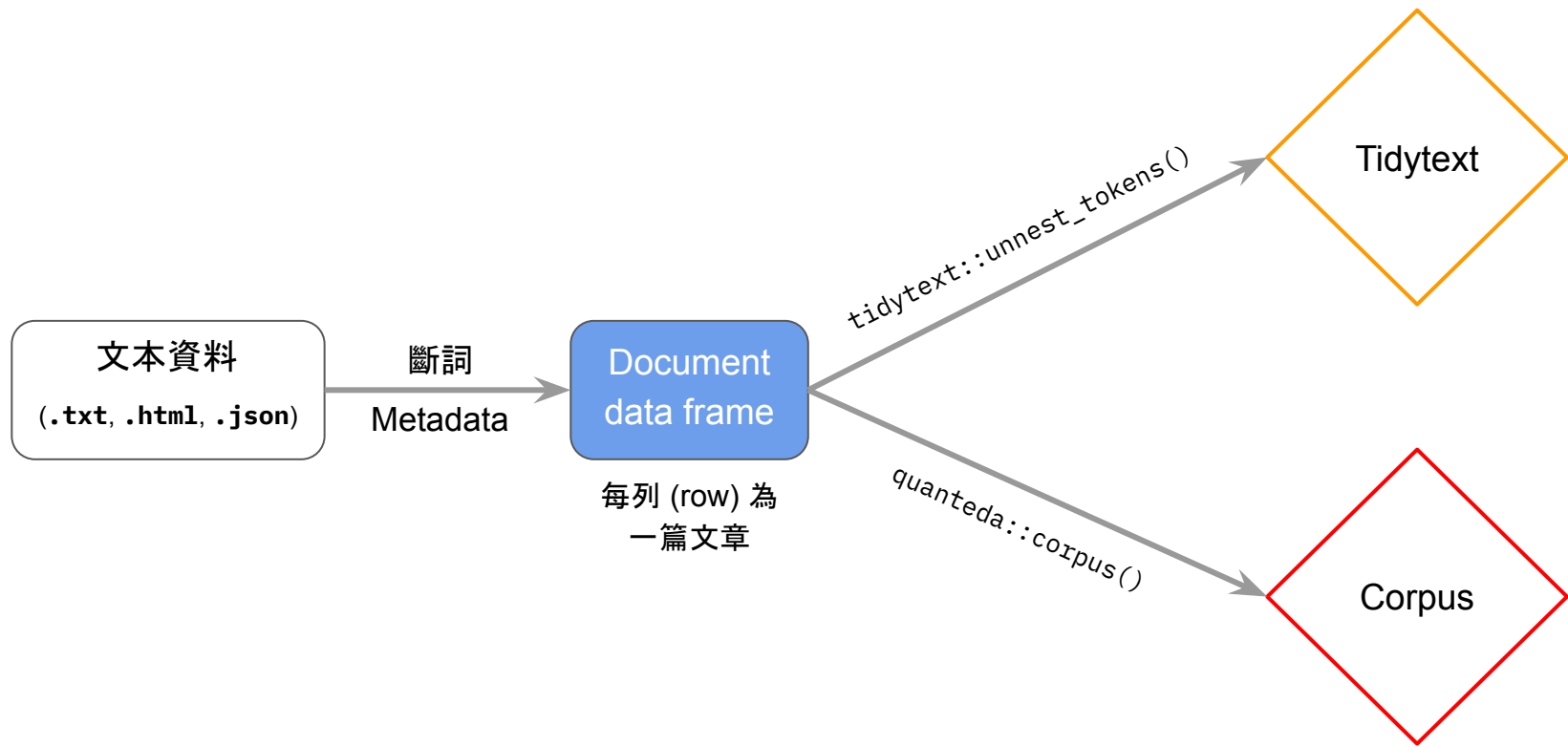


Lab 08

中文文本資料處理

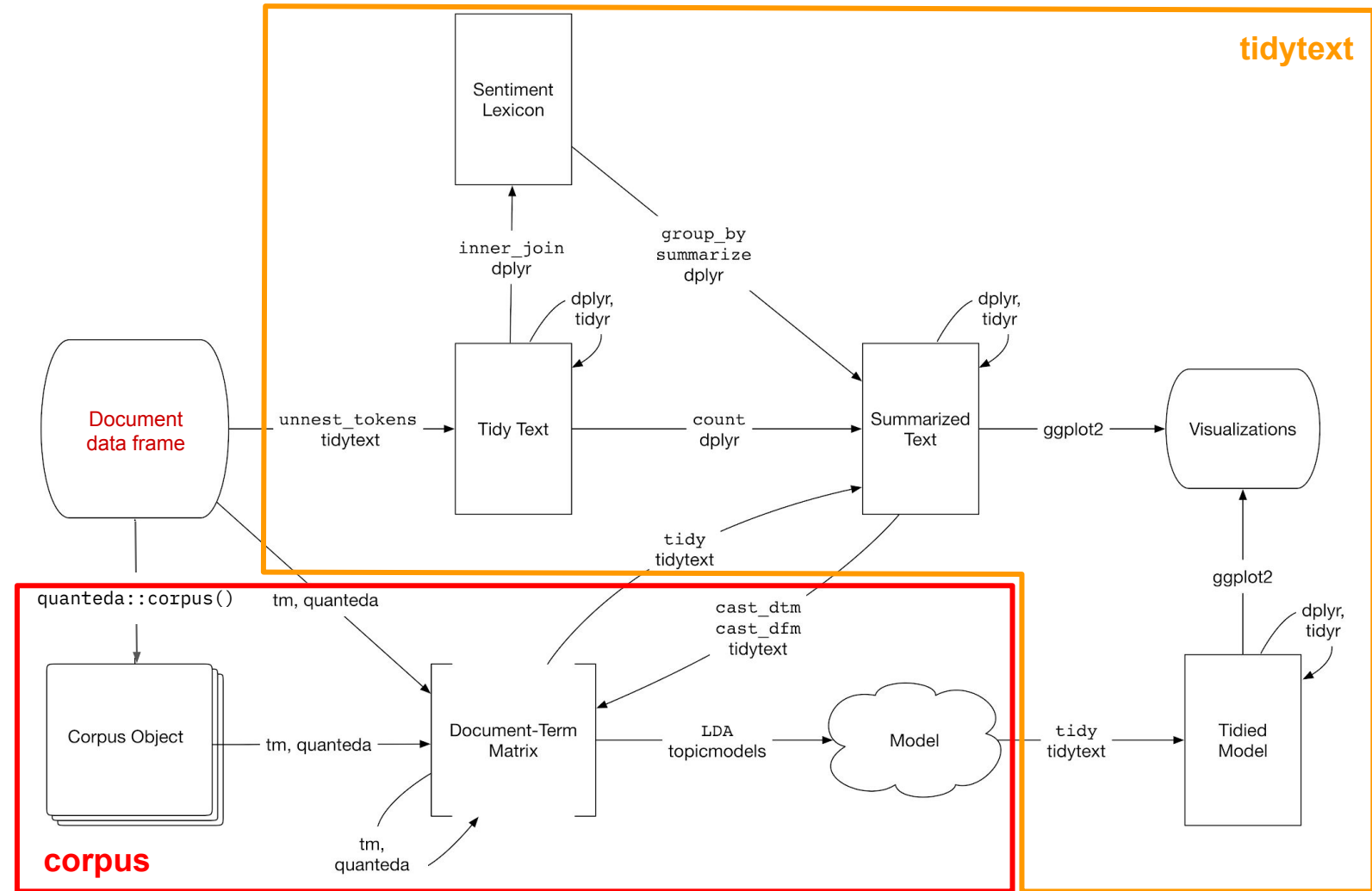
斷詞、tidytext、quanteda

R Text Mining Frameworks



Document data frame

doc_id	author	title	content
1	楊喚	夏夜	蝴蝶和蜜蜂們帶著花朵的蜜糖回來了羊隊和牛群告別了田野回家了火紅的太陽也滾著...
2	宋晶宜	雅量	朋友買了一件衣料綠色的底子帶白色方格當她拿給我們看時一位對圍棋十分感興趣的同學說啊好像棋盤似的我看倒有點像稿紙...
3	胡適	母親的教誨	每天天剛亮時我母親便把我喊醒叫我披衣坐起我從不知道她醒來坐了多久了她看我清醒了便對我說昨天我做錯了什麼事說錯了什麼話要我認錯要我用功讀書...
...



jiebaR

```
library(jiebaR)
seg <- worker(user = "path/to/user_dict")
segment("失業的熊讚陪柯文哲看銀翼殺手", seg)
```

Document data frame

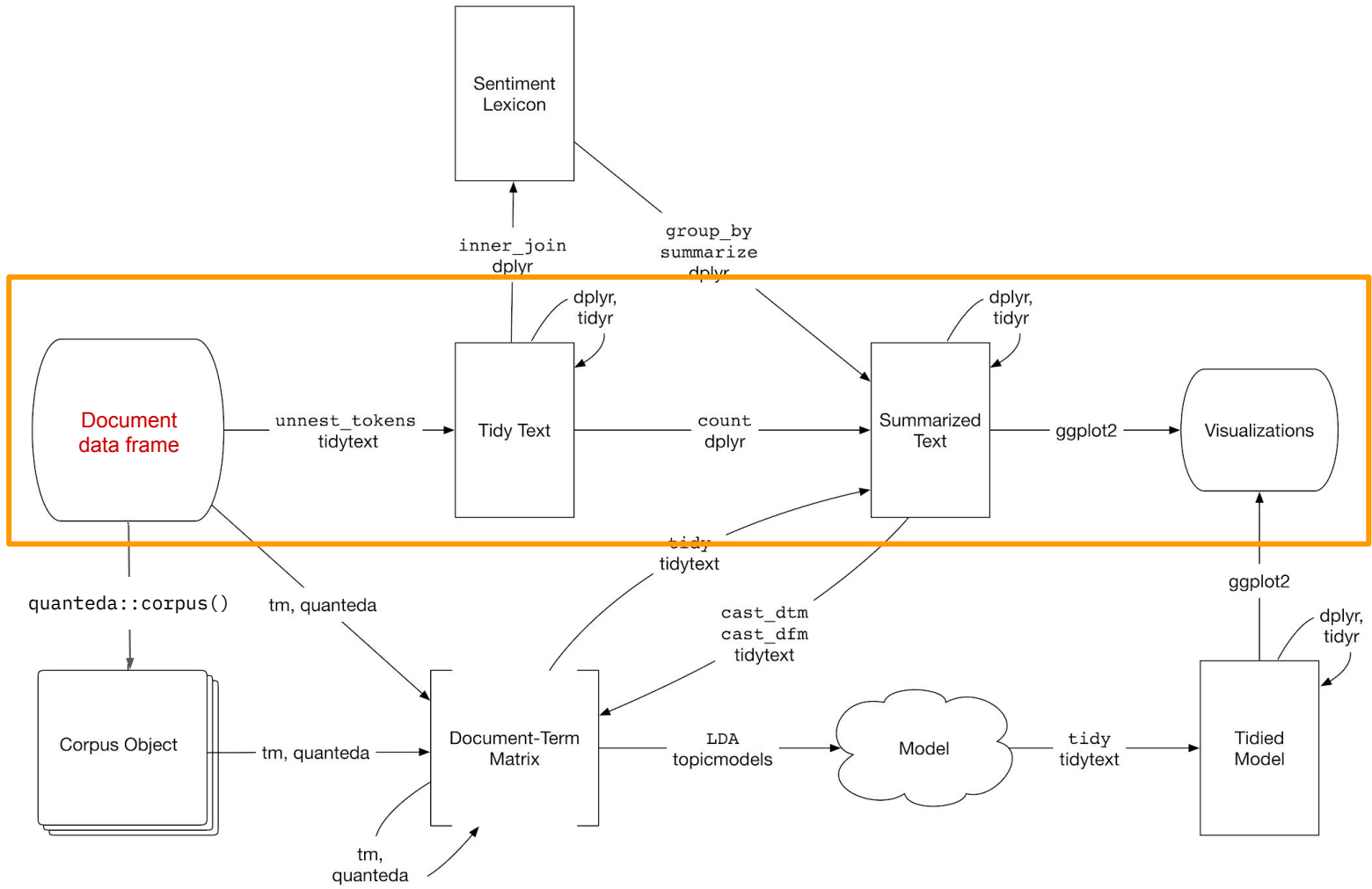
doc_id	author	title	content
1	楊喚	夏夜	蝴蝶和蜜蜂們帶著花朵的蜜糖回來了羊隊和牛群告別了 田野回家了火紅的太陽也滾著...
2	宋晶宜	雅量	朋友買了一件衣料綠色的底子帶白色方格當她拿給我們 看時一位對圍棋十分感興趣的同學說啊好像棋盤似的 我看倒有點像稿紙...
3	胡適	母親的教誨	每天天剛亮時我母親便把我喊醒叫我披衣坐起我從不 知道她醒來坐了多久了她看我清醒了便對我說昨天我 做錯了什麼事說錯了什麼話要我認錯要我用功讀書...

語料

```
docs_df <- readRDS("samesex_marriage.rds")
docs_df
```

id	topic	title	content
pro_292.txt	pro	國防大學生染愛滋被「逼」退學 醫:國軍應公開道歉	東森新聞記者嚴云岑台北報導 2016-8-19 愛滋病感染者權益促進會舉辦傳染病病人權益案例座談會...
anti_16.txt	anti	性別平等課本一點都不平等	不平等的公民課本 公民課綱中寫明要以平等尊重原則來撰寫課本內容 經查後發現撰寫...
pro_175.txt	pro	國際跨性別現身日 性別友善刻不容緩	近年來跨性別學生希望依照自身性別認同使用空間的聲音與訴求已浮現在台灣校園中 但在社會未正視跨性別者處境...
...

Tidytex framework



例子：詞頻表

unnest_tokens() + **dplyr** + **ggplot2**

quanteda framework

例子: Key Word in Context

`corpus()`

`tokenize_regex()`

`tokens()`

`kwic()`