

Tidy and Seamless in R

- 2020-10-14
- 유충현, Tidyverse Korea



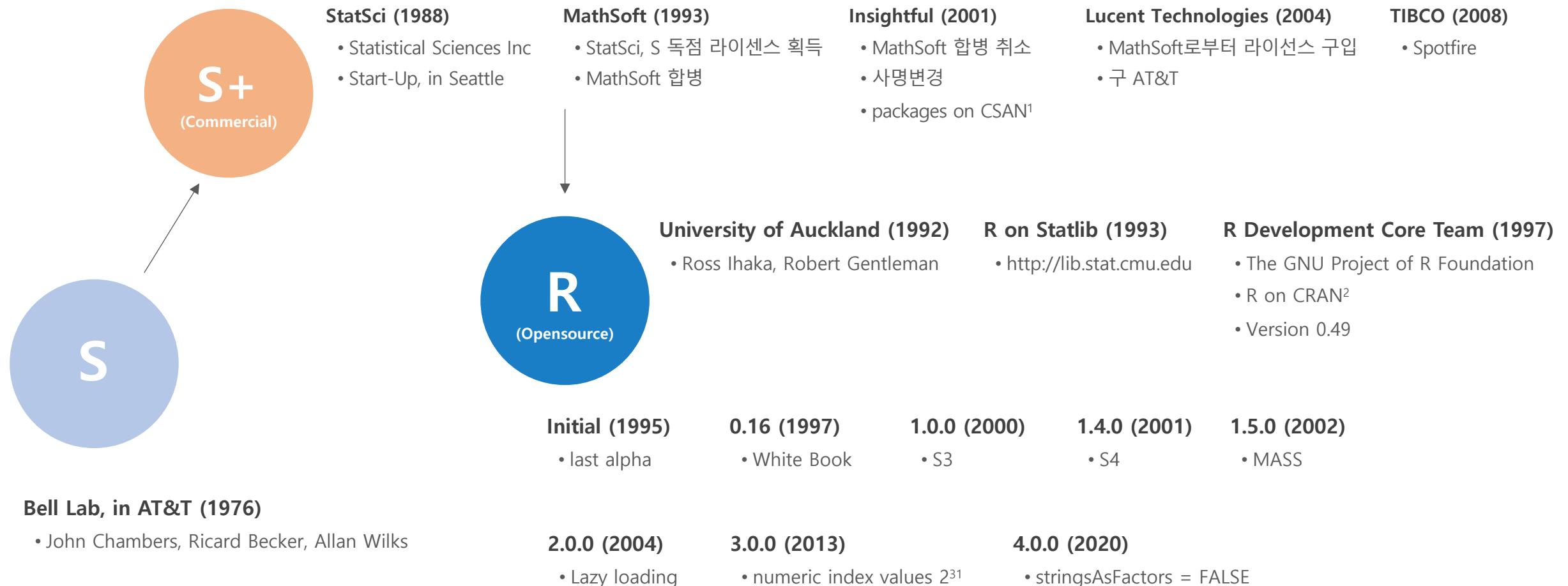
CONTENTS

1. R History
2. Paradigm Shift
3. Tidy & Seamless

Appendix

- Appendix 1. Seamless Example - Rmarkdown
- Appendix 2. Seamless Example – R Script
- Appendix 3. Seamless Example – Shell Script

1. R History – include S and S+



[https://en.wikipedia.org/wiki/S_\(programming_language\)](https://en.wikipedia.org/wiki/S_(programming_language))

<https://en.wikipedia.org/wiki/S-PLUS>

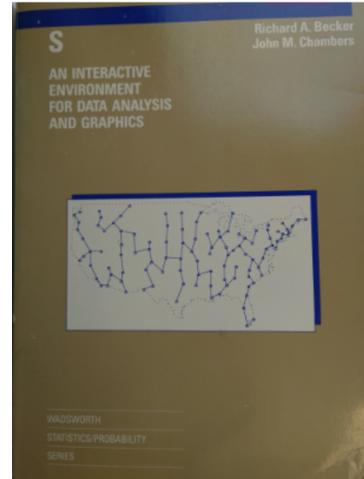
[https://en.wikipedia.org/wiki/R_\(programming_language\)](https://en.wikipedia.org/wiki/R_(programming_language))

1 : CSAN(The Comprehensive S Archive Network)

2 : CRAN (The Comprehensive R Archive Network)

S : An Interactive Environment for Data Analysis and Graphics

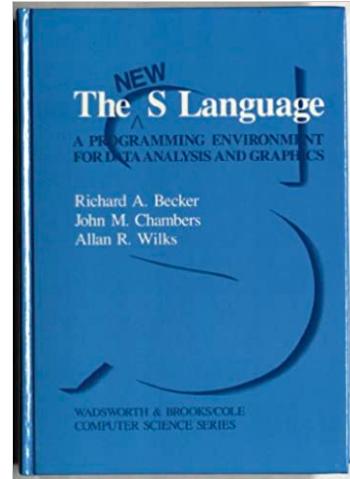
(Richard Becker, John Chambers)



1984 Brown Book

The New S Language

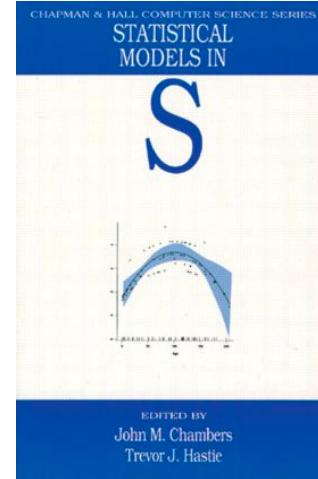
(Richard Becker, John Chambers, Allan Wilks)



1988 Blue Book

Statistical Models in S

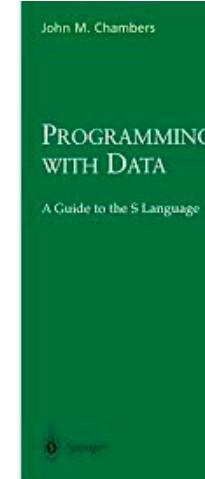
(John Chambers, Trevor Hastie)



1992 White Book

Programming with Data

(John Chambers)



1998 Green Book



R6: Encapsulated object-oriented programming for R



S2

start S3

completed S3

S4

R5

R6

- macros to functions
- 객체 컨셉
- Graphics device
 - X11, Postscript
- FORTRAN to C

- macros to functions
- data frame 객체 정의
- S3 Class (OOP)
 - Method dispatch
 - Generic function

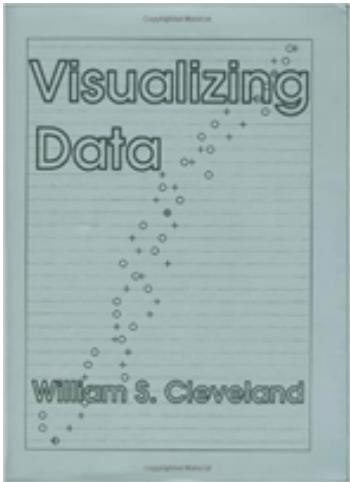
- formal S4 Class
 - methods 패키지
- Connections

- R5 Class
 - reference class
 - using S4

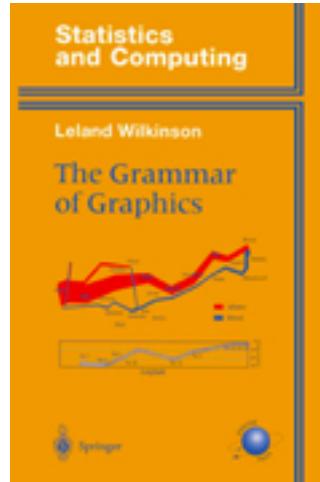
- R6 Class
 - public and private method
 - active binding

1. R History – Graphics

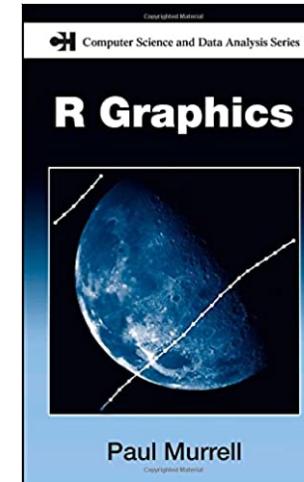
Visualizing Data
(William Cleveland)



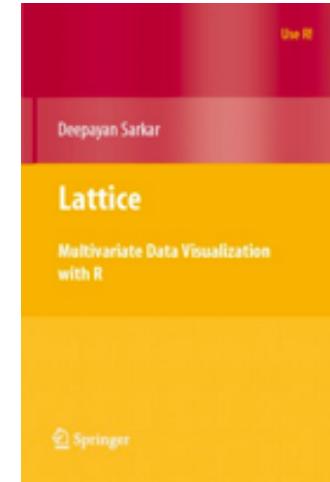
The Grammar of Graphics
(Leland Wilkinson)



R Graphics
(Paul Murrell)



Lattice
(Deepayan Sarkar)



ggplot2
(Hadley Wickham)



1993

1999

2005

2008

2009

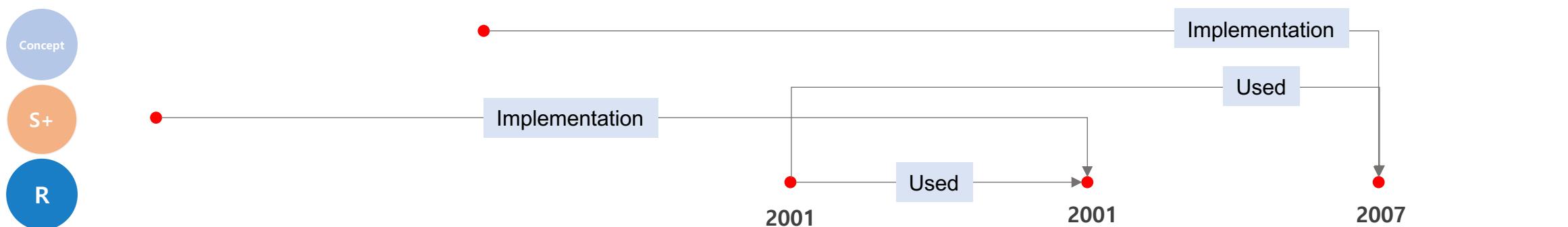
trellis package in S

Grammar Concept

grid package

lattice package

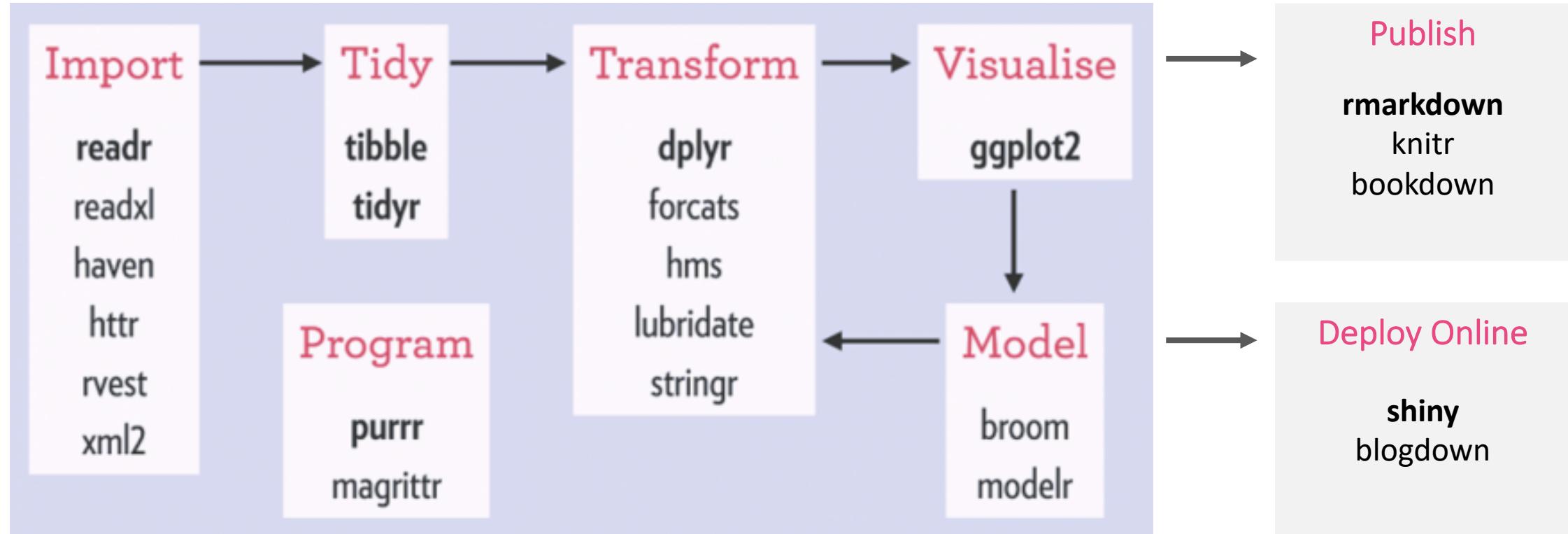
ggplot2 package





	<u>Before</u>	<u>After</u>
 Data Scale	Toy Data data.frame read.* In-Memory	Enterprise Data tibble read_* In-Database
 Package Developer	Statistician R Development Core Team Statistics	S/W Developer Hadley 사단 Machine Learning
 Distribution	Commercial S-Plus (CSAN)	Open Source R (CRAN, GitHub)
 Analytics Methods	R Graphics, lattice each packages	ggplot2 tidyverse, tidymodel

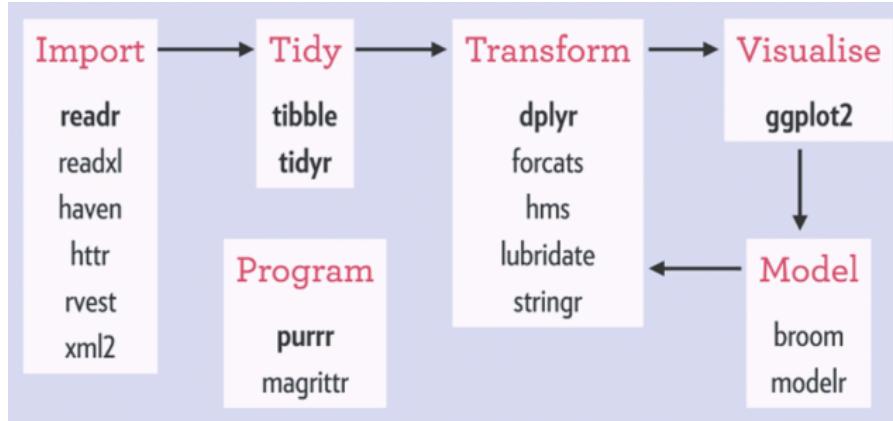
Seamless with tidyverse



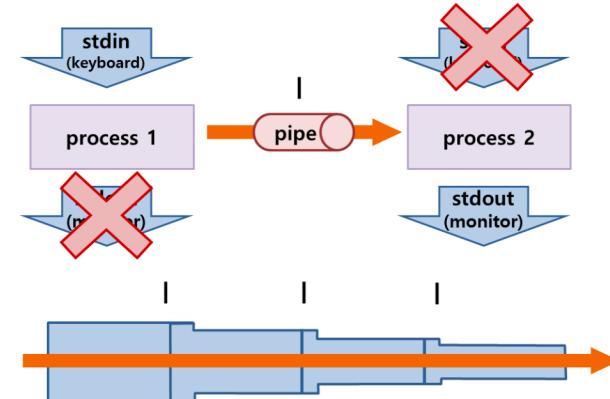
<https://www.kdnuggets.com/2017/10/tidyverse-powerful-r-toolbox.html>

Within Seamless

From import to model using tidyverse



Shell script using pipe



Pipe using magrittr package



Between Seamless



Rmarkdown을 이용한 Seamless 프로그래밍

- 시나리오 : 파일의 용량 증가에 따른 전통적인 `read.csv()`와 `readr` 패키지의 `read_csv()`, Python `pandas` `read_csv()`의 성능 비교
- 소스 : `simulation_import.Rmd`



1. R Chunk

- Seed data 생성



2. Bash Chunk

- 복수배 파일 생성



3. R Chunk

- 성능측정
- `read.csv()`, `read_csv()`



4. Python Chunk

- 성능측정
- `pandas's read_csv()`



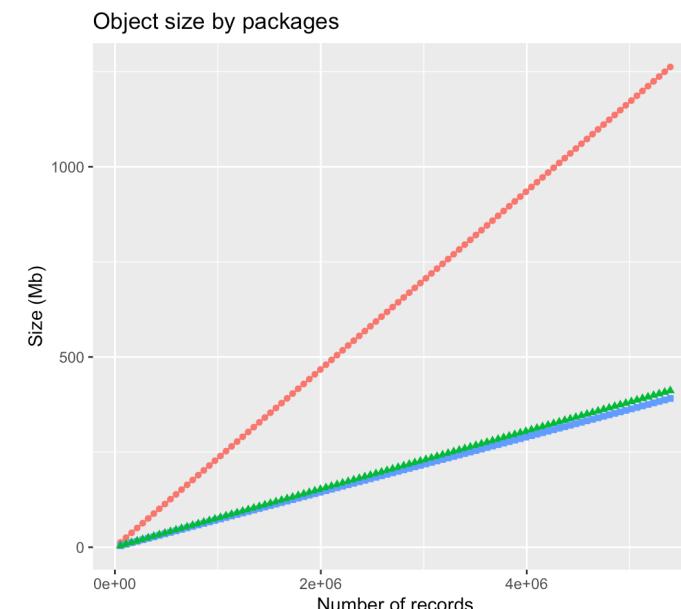
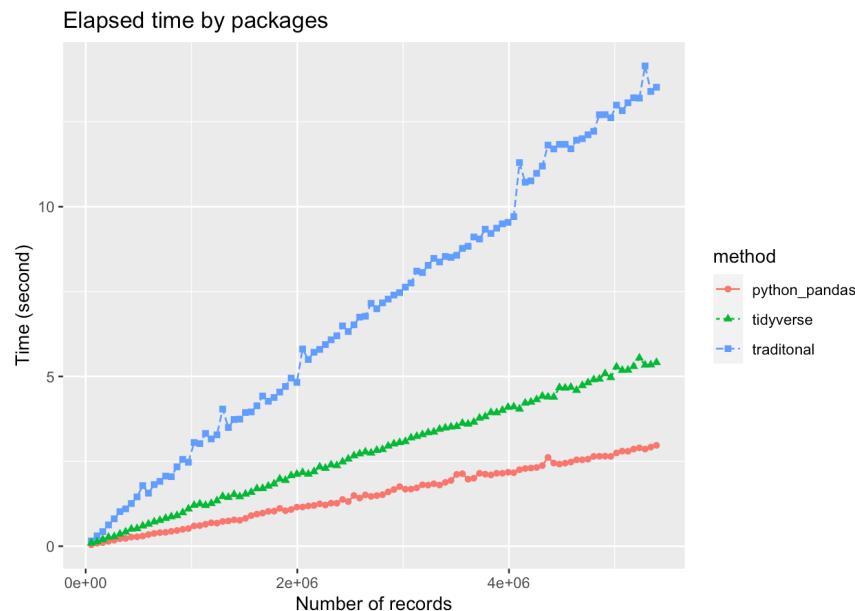
5. R Chunk

- 성능 비교
- 시각화



6.Bash Chunk

- 파일 삭제



Appendix 2. Seamless Examples

reticulate package로 Python 연동하기

- 시나리오 : R에서 생성한 대통령 연설문 Term Frequency 데이터를 Python에서 word cloud 그리기
- 소스 : call_python_from_R.R, speech.rda



1. R

- 데이터 읽어 들이기
- Term Frequency 벡터 조작

2. Python 호출

- R 벡터를 python의 dictionary로 변환
- word cloud 그리기



Shell script에서 R 연동하기

- 시나리오 : /var/log/system.log 파일에서 서비스별 호출 회수를 구하고, 이를 시각화
- 소스 : get_summary_system.sh



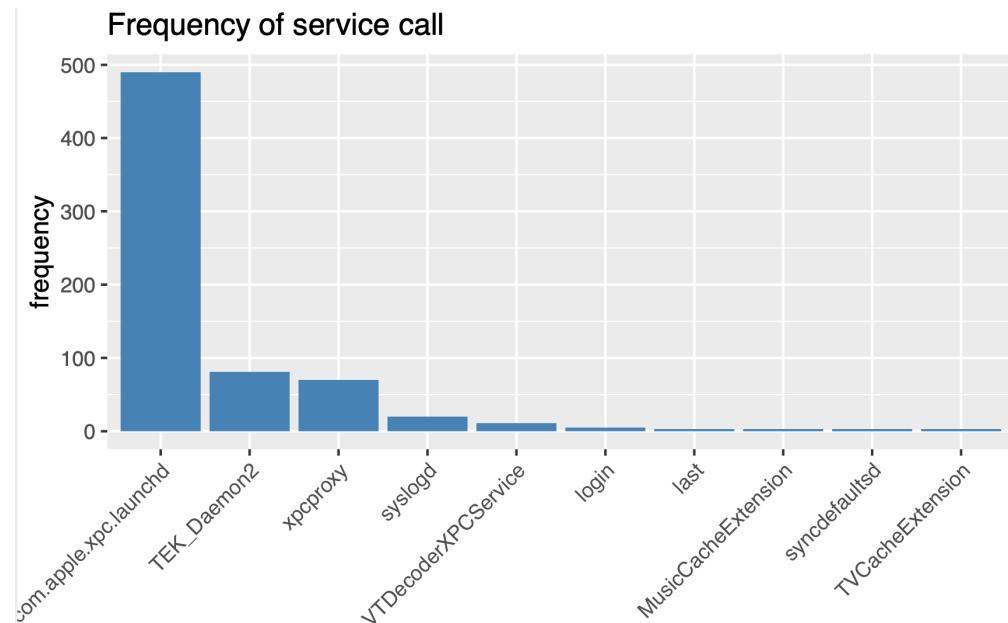
1. Shell Script

- pipe를 이용해서 서비스별 호출 회수 계산
- 결과를 파일에 저장



2. R 호출

- 파일을 읽어 data frame 생성
- ggplot2로 bar chart 그리기





E. O. D