# Document as a Data

Tidyverse Korea

이광춘

# Table of Contents
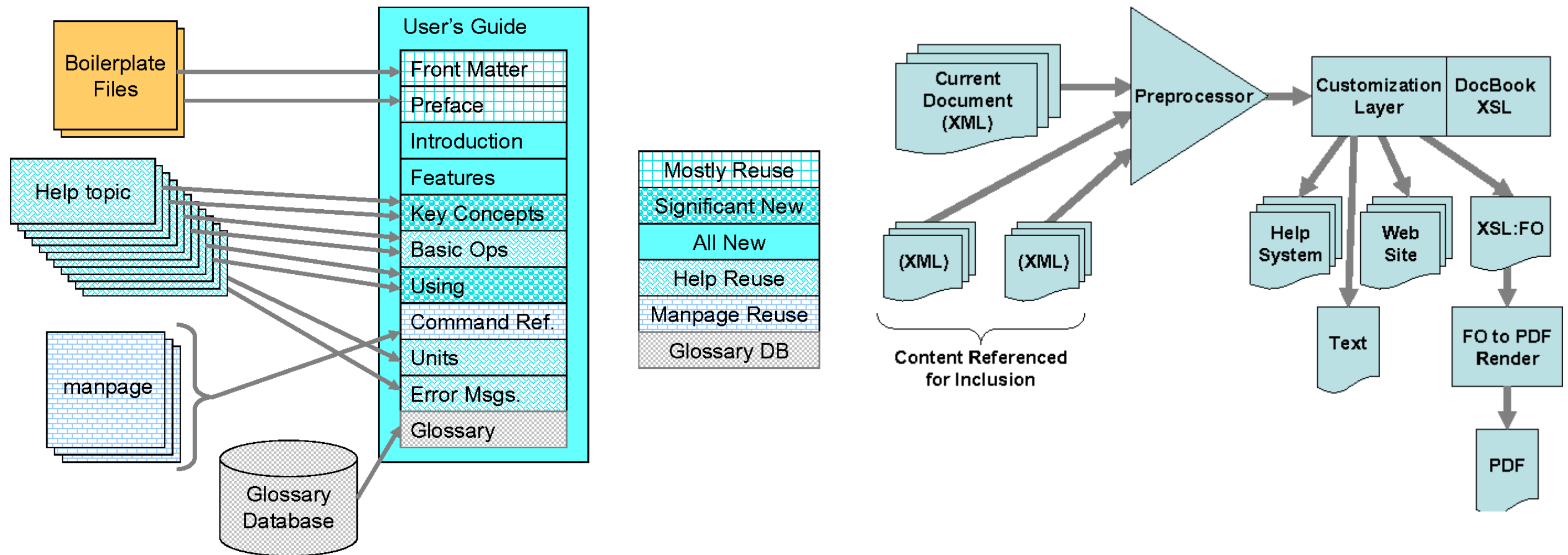
# Document Anatomy – Help Manual

## Anatomy of Hewlett-Packard User Manual Documents
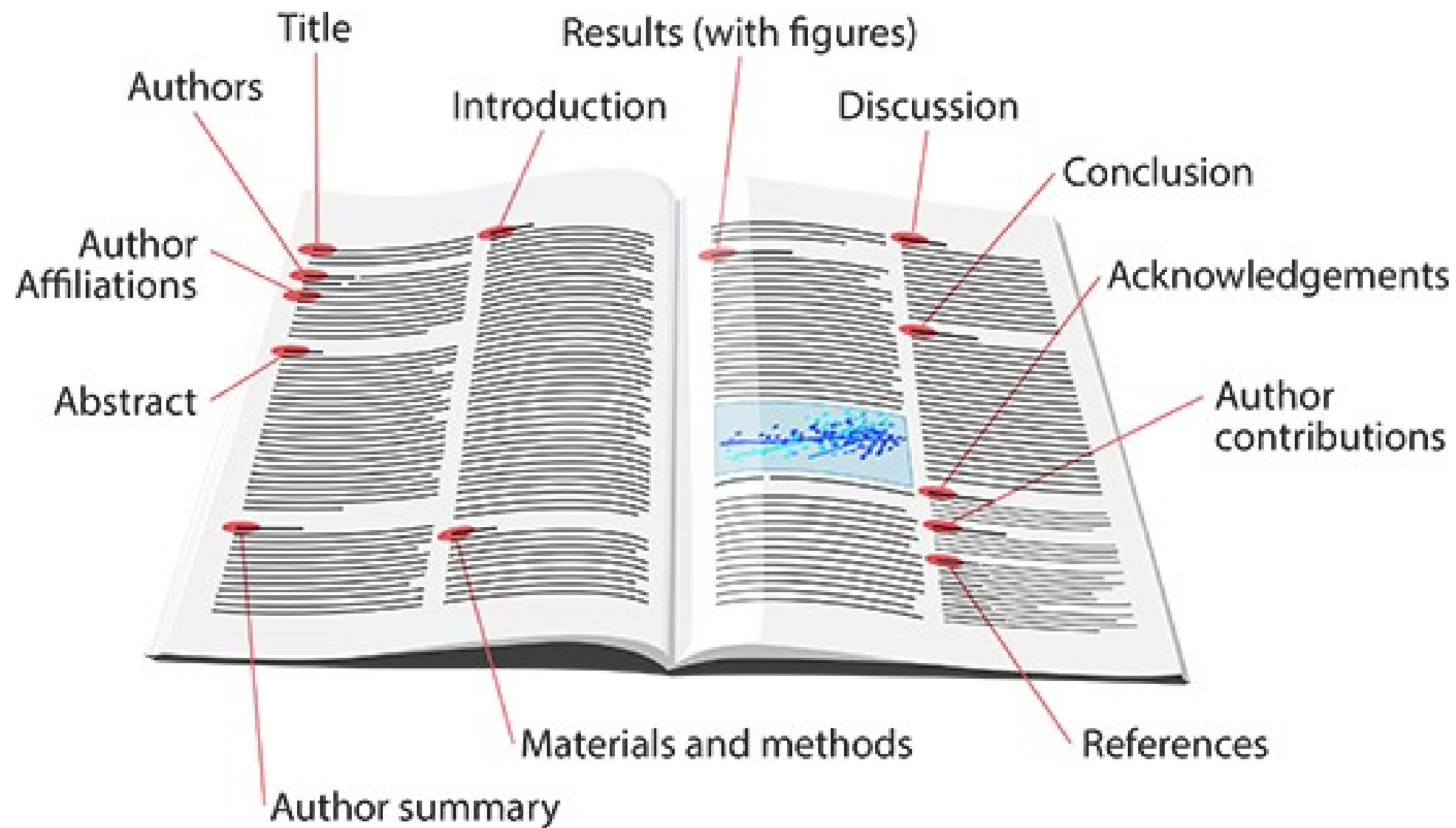
Kathy Haramundanis
 and Larry Rowland (Hewlett-Packard Company), "Experience paper: a content reuse documentation design experience", SIGDOC '07 Proceedings of the 25th annual ACM international conference on Design of communication Pages 229-233, El Paso, Texas, USA — October 22 - 24, 2007
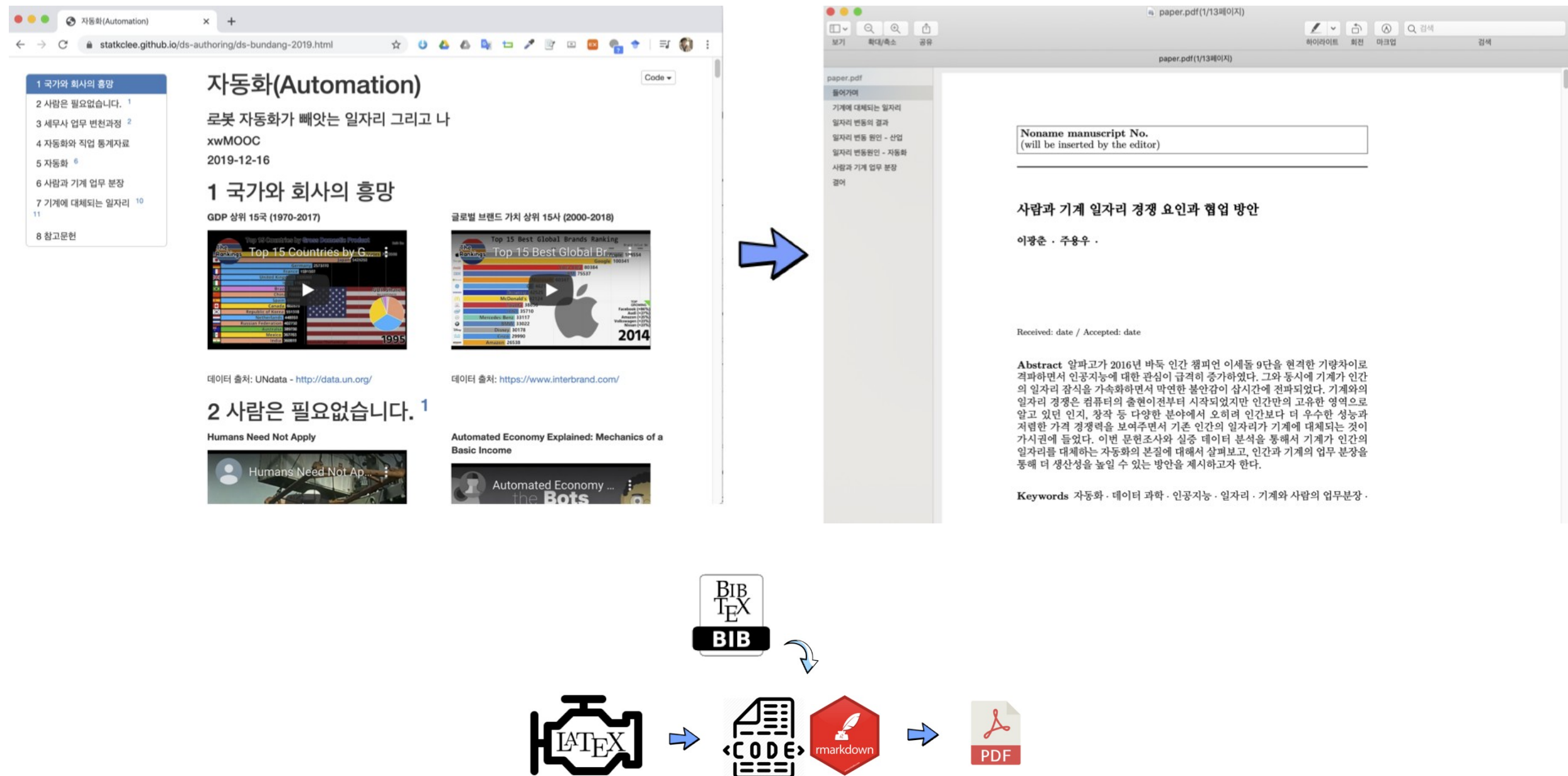
# Document Anatomy

## Anatomy of a research paper

# Document Anatomy - Papers

## Creating a research paper through Coding



출처: https://statkclee.github.io/comp_document/automation-kasdba.html
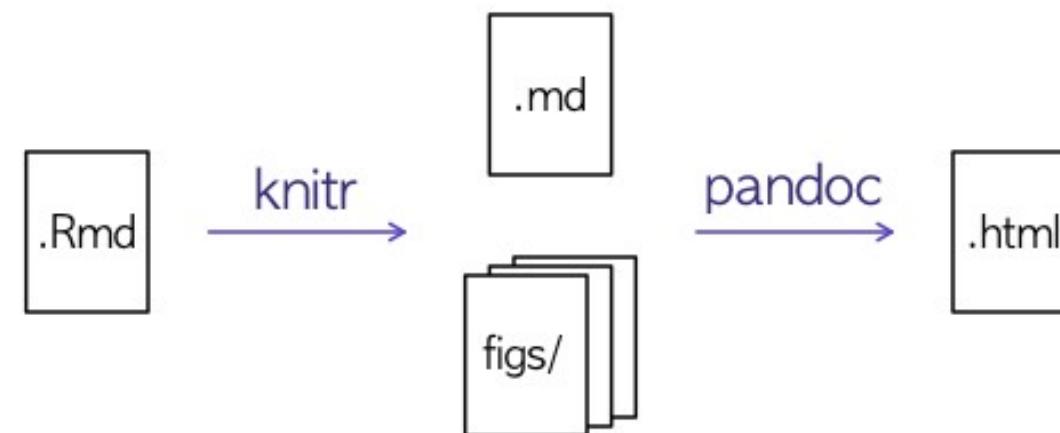
# Compendium

## Document as a Code

```
.
├── .Rprofile
├── .gitignore
├── README.md
├── analysis
│   └── archive
│   └── markdown
├── data
│   ├── documentation
│   ├── handmade
│   ├── html_reports
│   ├── processed
│   ├── public
│   └── source
├── etl
├── publish
├── scratch
├── viz
└── {{cookiecutter}}.Rproj
```

## Version Control



## Literate Programming

# Compendium – Reproducible Research

### 데이터 분석



### 분석 보고서



### 분석 팩키지



### 재현가능한 팩키지



출처: https://statkclee.github.io/comp_document/cd_compendium.html

# OCR

## OCR (Optical Character Recognition)
이미지 형태의 문서를 텍스트 형태의 문서로 변환해주는 기술



- 텍스트
- 이미지
- 표
- 그래프
- 문서구조
- 참고문헌
- …

# OCR 프로세스



**Preprocessing** → **Text Detection** → **Text Recognition** → **Information Extraction** → **Database dump**

| <원본 이미지> | <전처리 이미지> | <잘라낸 영역> | <글자 외곽선을 통해 추출한 이미지> |
|---|---|---|---|

**문서 종류**

**전화번호**

감사합니다.