# SevenBridges

# **Persistent Reproducible Reporting**

Nan Xiao, Seven Bridges

2017/05/20 @ China R Conference Beijing

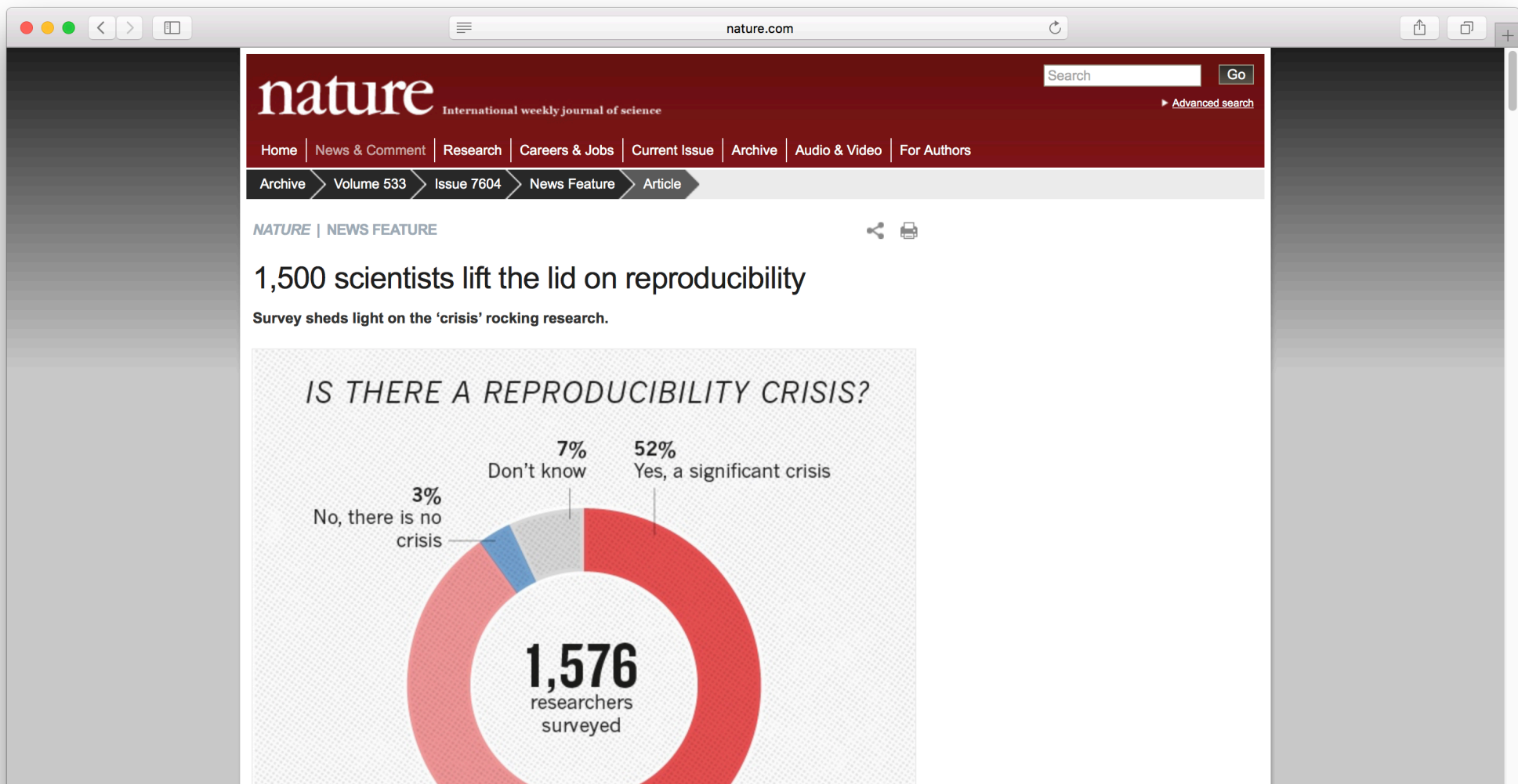# DOCUMENT-LEVEL REPRODUCIBILITY
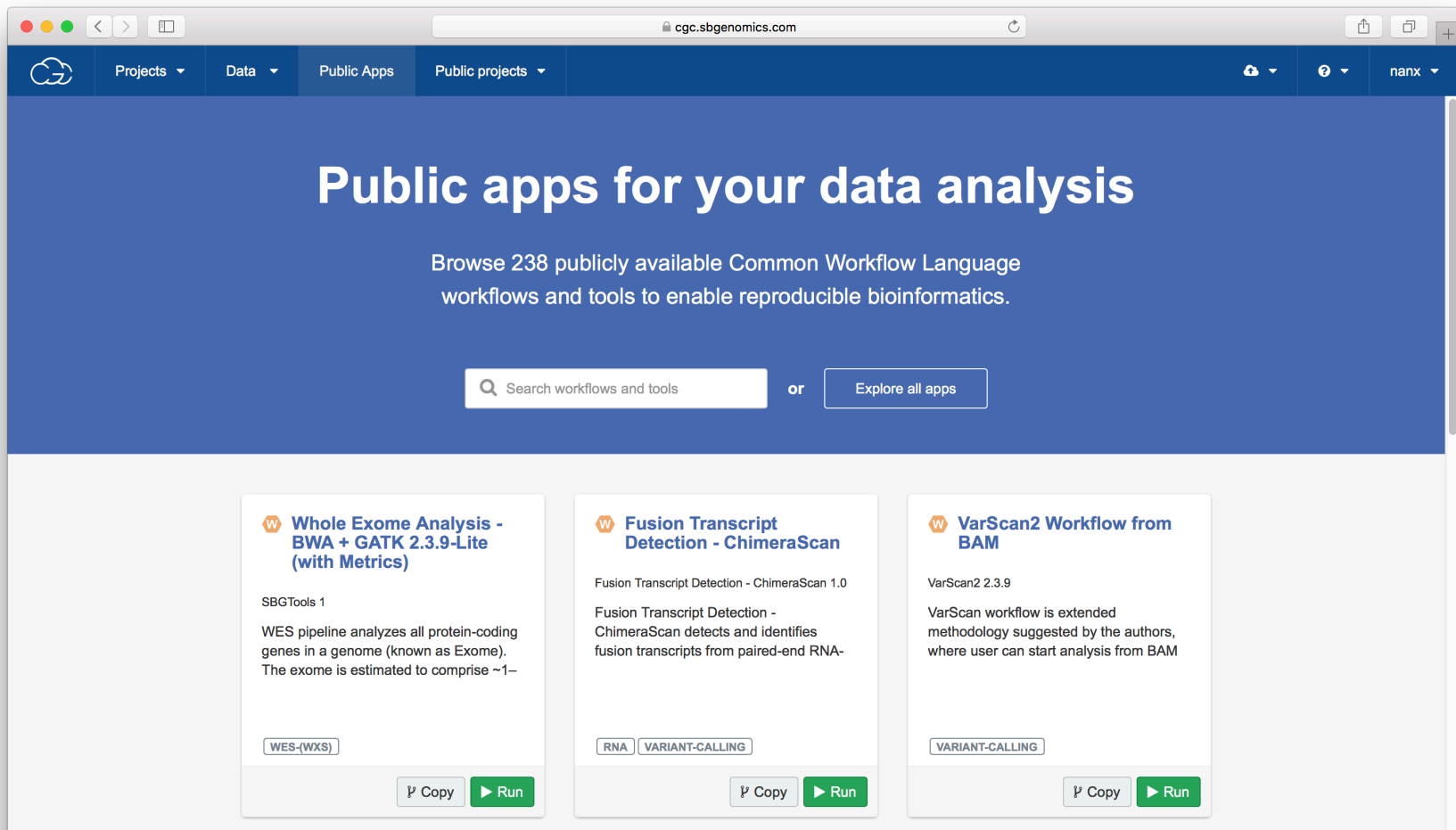
# R MARKDOWN + KNITR TO THE RESCUE

# REPRODUCIBILITY

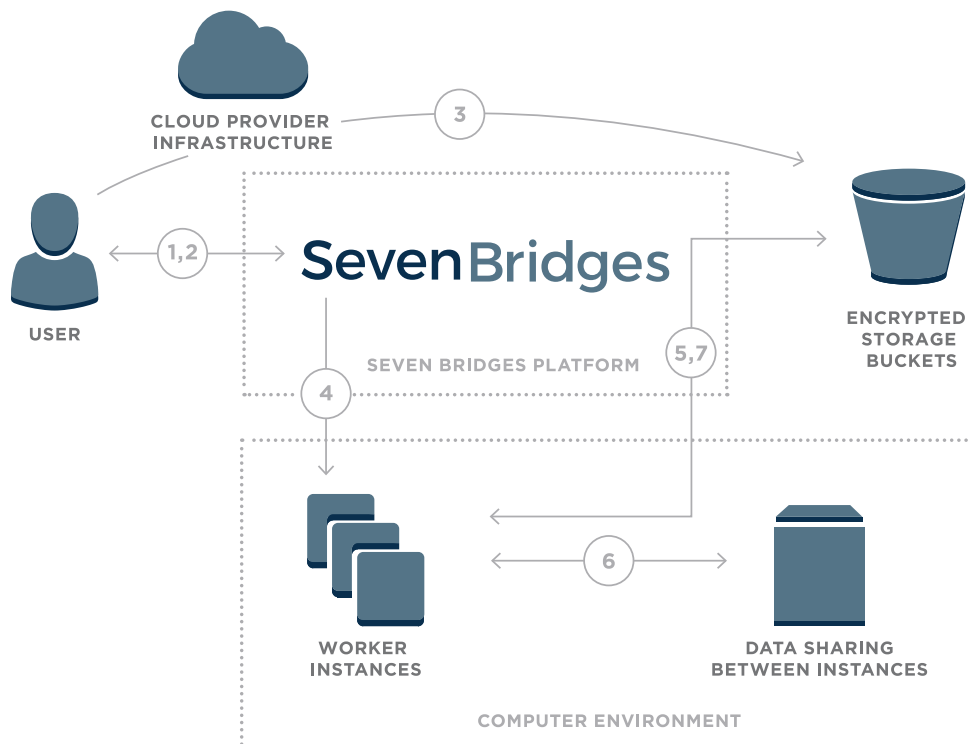... has always been a concern in both academia & industry.

# CANCER GENOMICS CLOUD (CGC)

- [www.cancergenomicscloud.org](http://www.cancergenomicscloud.org)
- Hundreds of automated analysis workflows for petabyte-scale data from The Cancer Genome Atlas

# PRODUCT & ENGINEERING INNOVATIONS IN CGC



CLOUD PROVIDER
INFRASTRUCTURE

3

USER

1,2

SevenBridges

SEVEN BRIDGES PLATFORM

4

5,7

ENCRYPTED
STORAGE
BUCKETS

WORKER
INSTANCES

6

DATA SHARING
BETWEEN INSTANCES

COMPUTER ENVIRONMENT
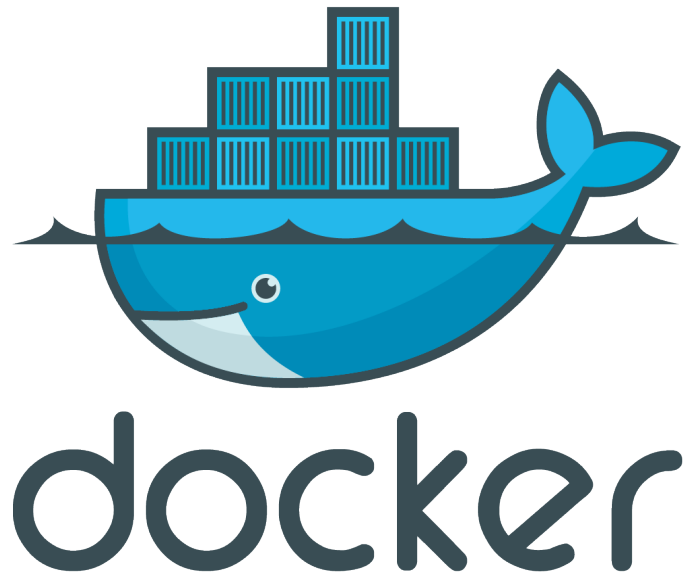
COMMON
WORKFLOW
LANGUAGE

Rabix

# CHALLENGE: OS-LEVEL REPRODUCIBILITY

How to ensure your reports are reproducible across <u>time</u> and <u>environments</u>, when the data, analysis tools, operating systems are all <u>evolving</u>?

# DOCKER

- Docker allows applications and their dependencies to be packaged into discrete runtime environments, called <u>containers</u>. Applications packaged in this way can be run from many diverse infrastructures.
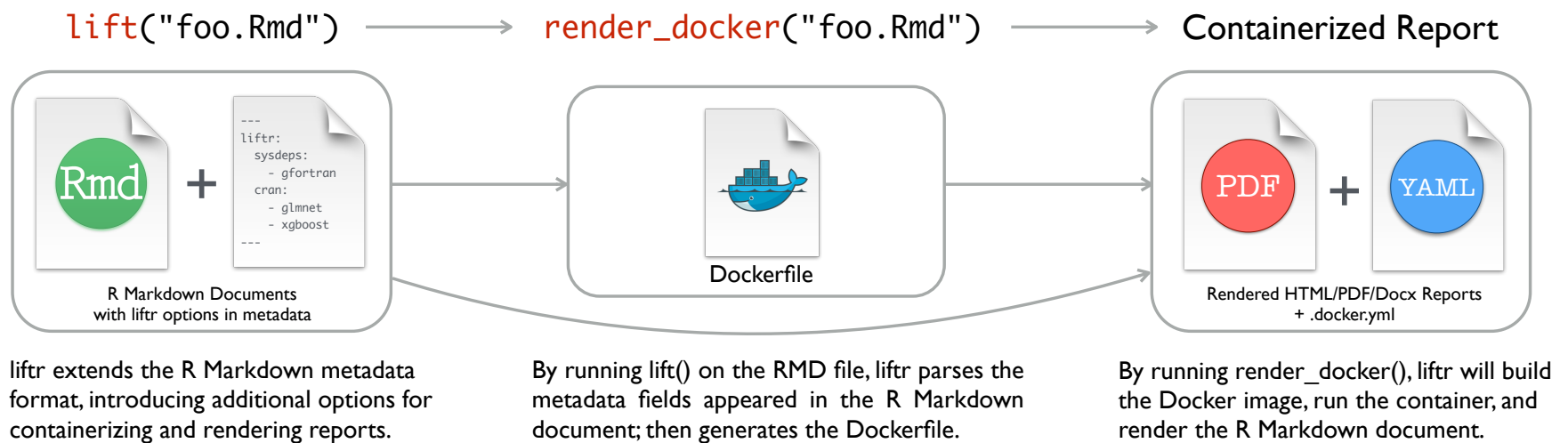
# liftr

OS-level <u>reproducibility</u> & <u>persistency</u> for reports.



docker + Rmd + knitr = liftr

# DOCKERIZE DOCUMENTS AS EASY AS 1-2-3

`lift("foo.Rmd")` → `render_docker("foo.Rmd")` → Containerized Report

```
---
liftr:
  sysdeps:
    - gfortran
  cran:
    - glmnet
    - xgboost
---
```

**R Markdown Documents**
**with liftr options in metadata**

**Dockerfile**

**Rendered HTML/PDF/Docx Reports**
**+ .docker.yml**

liftr extends the R Markdown metadata format, introducing additional options for containerizing and rendering reports.

By running lift() on the RMD file, liftr parses the metadata fields appeared in the R Markdown document; then generates the Dockerfile.

By running render_docker(), liftr will build the Docker image, run the container, and render the R Markdown document.

# DOCKERIZE DOCUMENTS AS EASY AS 1-2-3

```r
library("liftr")
input = "demo.Rmd"

lift(input)                    # Generate Dockerfile
render_docker(input)           # Render report with Docker

purge_image(input)             # Clean up Docker image
push_image(input)              # Push image to registry (devel)
```
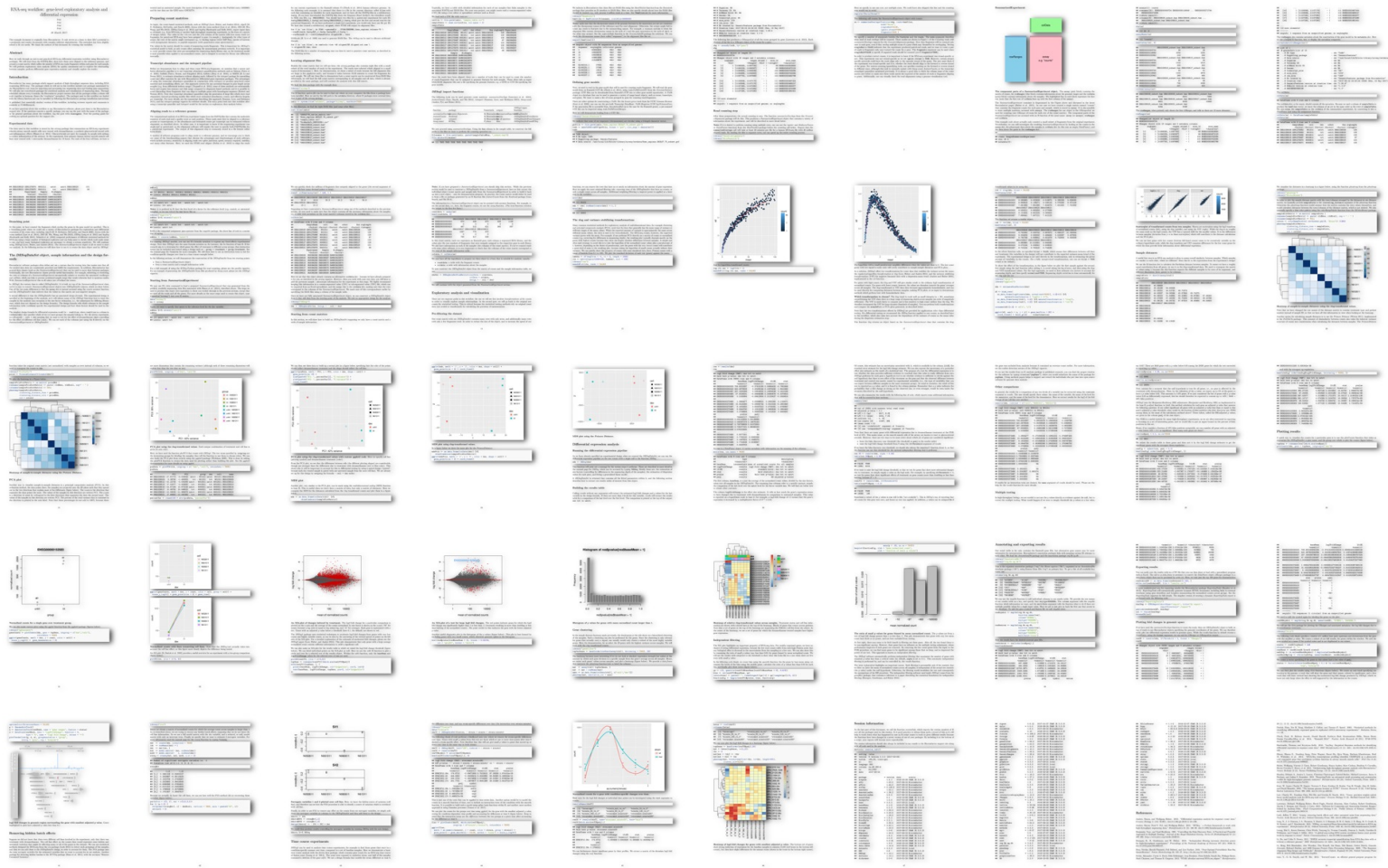
# DEMO: RNA-SEQ DATA ANALYSIS

Example workflow from Bioconductor.org

- RNA-Seq: biotechnology for measuring the expression of genes. It can help identify potential key genes in cancer.

- TBs of RNA-Seq data are generated. Computational tools and workflows are developed to analyze such data.

- We need to ensure such reports are reproducible through time, when datasets, analysis tools are both evolving.

- Code available from: bit.ly/liftrdemo

# STEP 1

Add liftr metadata to the
R Markdown document:

- Base image
- System dependencies
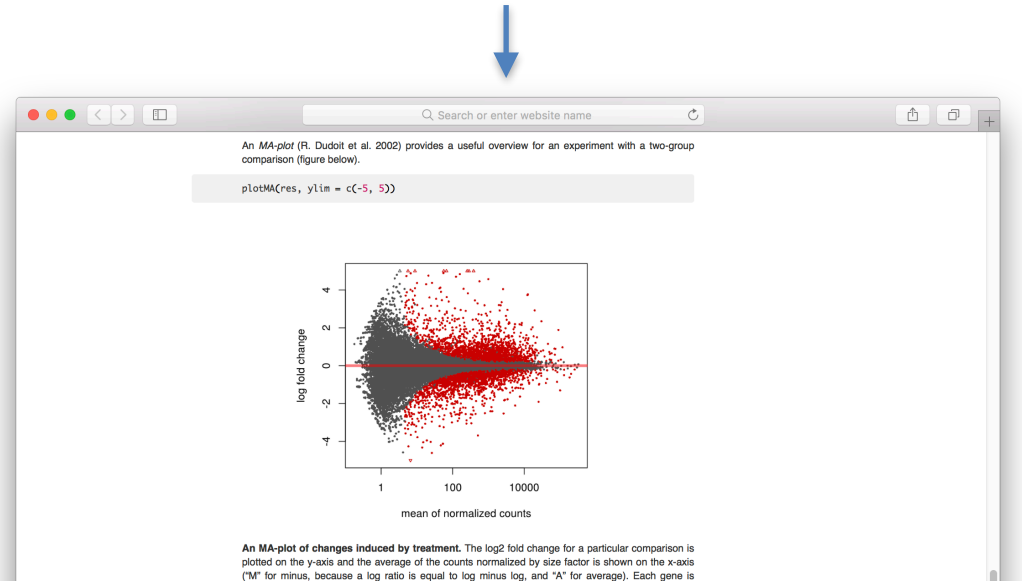- Package dependencies
- …

**STEP 2**

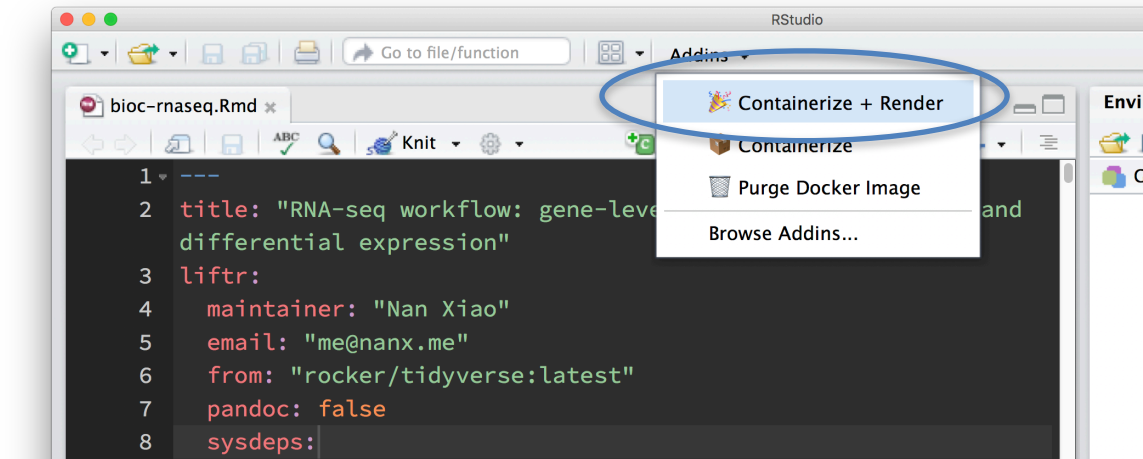# Generate Dockerfile with liftr::lift

# STEP 3

- **liftr::render_docker** will <u>build</u> the Docker image, <u>run</u> the container, and <u>render</u> into PDF/HTML/DOCX.

- Re-compile: <u>cached</u> Docker image layers are used to improve speed.

- Remove the used image, or push to Docker registry.

# FUTURE WORKS

We aim to expand the  R Markdown tool chain by exploring the next frontier: system-level reproducibility, and democratize reproducible report creation/sharing.

To achieve this, we need:

- Standard renderers + independent YAML configuration file
- Better IDE support (RStudio Addins)
- Better on-boarding experience: automatic dependency parsing
- Cloud-based rendering and containerization services for dynamic documents

# Q & A

Visit liftr.me for more info

Contact: me@nanx.me