

# Reproducible Research Pipelines Using R and RStudio

Best Practices for Analysis and Dissemination

Melinda K. Higgins, PhD.

March 21, 2018: 12:30pm – 3:30pm, EST



# Workshop Outline

- Module 01: Reproducible Research
- Module 02: Workshop materials, RStudio Interface, Getting Started with R
- Module 03: Understanding R, Working with objects, R Scripts, R Packages
- — BREAK 10 minutes —
- Module 04: Creating Documents with R Markdown
- Module 05: Create Document = R Script + R Markdown
- Module 06: Customizing R Markdown (templates, parameters and automation)

# Module 01

## Reproducible Research

# Outline

- Timeline Reproducible Research & Transparency
- People
- Books
- Literate Programming > Dynamic Documentation > [R]Markdown
- The Big Picture

# Timeline Reproducible Research & Transparency<sup>1</sup>

YEAR	Event
1992	Jon Claerbout coined the term "reproducible research" in his book "EARTH SOUNDINGS ANALYSIS: Processing versus Inversion (PVI)" <sup>2</sup>
1996	CONSORT statement introduced standards for reporting clinical trials <sup>3</sup>
2004	International Committee of Medical Journal Editors (ICMJE) stated they would not publish a clinical trial that had not been registered. <sup>4</sup>
2005	Ioannidis, J. P. A. Why most published research findings are false. PLoS Med. 2, e124 (2005) <sup>5</sup>

1. Timeline partially based on PLOS Blog December 2016 <http://blogs.plos.org/absolutely-maybe/2016/12/05/reproducibility-crisis-timeline-milestones-in-tackling-research-reliability/>

2. <http://sepwww.stanford.edu/sep/jon/reproducible.html>

3. Begg C, Cho M, Eastwood S, Horton R, Moher D, Olkin I, Pitkin R, Rennie D, Schulz KF, Simel D, Stroup DF (1996). Improving the quality of reporting of randomized controlled trials. The CONSORT statement. JAMA 276:637-639.

4. [http://www.icmje.org/news-and-editorials/update\\_2005.html](http://www.icmje.org/news-and-editorials/update_2005.html)

5. <https://doi.org/10.1371/journal.pmed.0020124>

# Timeline Reproducible Research & Transparency

YEAR	Event
2007	FDA Amendments Act (FDAAA) required more types of clinical trials to be registered (final rules took effect January 2017) <sup>6</sup>
2009	Journal of Biostatistics institutes policy to work with authors to publish articles that meet a standard of reproducibility. <sup>7</sup>
2011	Alsheikh-Ali, et.al. (2011), report the low percentage of researchers satisfying the policies regarding the availability and sharing of their data. <sup>8</sup>

6. <https://clinicaltrials.gov/ct2/manage-recs/fdaaa>

7. <https://academic.oup.com/biostatistics/article/10/3/405/293660/Reproducible-research-and-Biostatistics> & [https://academic.oup.com/biostatistics/pages/General\\_Instructions](https://academic.oup.com/biostatistics/pages/General_Instructions)

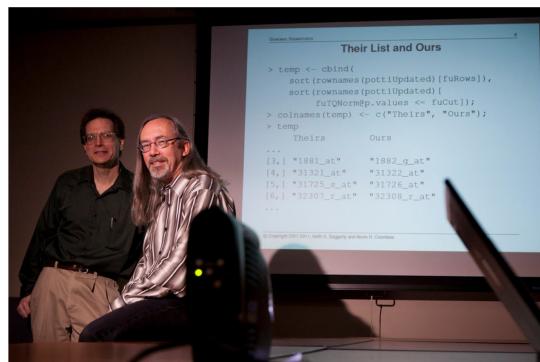
8. Alsheikh-Ali, A. A., Qureshi, W., Al-Mallah, M. H. & Ioannidis, J. P. Public availability of published research data in high-impact journals. PloS ONE 6, e24357, 2011; <https://doi.org/10.1371/journal.pone.0024357>

# Cancer Testing Falls Apart

The screenshot shows a news article from The New York Times. At the top, there's a navigation bar with links to various websites like nytimes.com, melindahiggins2000, Google, Online MBSR/Mindfulness, Druid Hills High School, and Campus Parent. Below the navigation, the main headline reads "How Bright Promise in Cancer Testing Fell Apart". Underneath the headline, it says "By GINA KOLATA JULY 7, 2011". To the left of the text, there are two small images: one of a man and one of a woman. To the right, there are three smaller images with captions: "A. Approves First Gene-Altering Leukemia Treatment. Costing \$5,000", "ECONOMIC SCENE Home Health Care: Shouldn't It Be Work Worth Doing?", and "ED.A. Cracks Down on Unscrupulous Stem Cell Clinics". There's also a sidebar with a photo of a man and the caption "TRUE & NUMBER New Fathers Are Than Ever".

## How Bright Promise in Cancer Testing Fell Apart

By GINA KOLATA JULY 7, 2011



Keith Baggerly, left, and Kevin Coombes, statisticians at M. D. Anderson Cancer Center, found flaws in research on tumors. Michael Stravato for The New York Times

The Duke saga began when a prestigious journal, *Nature Medicine*, [published a paper](#) on Nov. 6, 2006, by Dr. Anil Potti, a cancer researcher at Duke University Medical Center; Joseph R. Nevins, a senior scientist there; and their colleagues. They wrote about genomic tests they developed that looked at the molecular traits of a cancerous tumor and figured out which [chemotherapy](#) would work best.

First, though, he asked two statisticians at M. D. Anderson, Keith Baggerly and Kevin Coombes, to check the work. Several other doctors approached them with the same request.

Dr. Baggerly and Dr. Coombes found errors almost immediately. Some seemed careless — moving a row or a column over by one in a giant spreadsheet — while others seemed inexplicable. The Duke team shrugged them off as “clerical errors.”

And the Duke researchers continued to publish papers on their genomic signatures in prestigious journals. Meanwhile, they started three trials using the work to decide which drugs to give patients.

Dr. Baggerly and Dr. Coombes tried to sound an alarm. They got the attention of the National Cancer Institute, whose own investigators wanted to use the Duke system in a clinical trial but were dissuaded by the criticisms. Finally, they [published their analysis](#) in *The Annals of Applied Statistics*, a journal that medical scientists rarely read.

<http://www.nytimes.com/2011/07/08/health/research/08genes.html>

2010 Video Presentation by Keith A. Baggerly

[http://videolectures.net/cancerbioinformatics2010\\_baggerly\\_irrh/](http://videolectures.net/cancerbioinformatics2010_baggerly_irrh/)

# The Excel-Error Heard Around the World

<https://newrepublic.com/article/112951/rogoff-reinhart-and-world-excel-error-research>  
Ph melindahiggins2000 Google Online MBSR/Mindfu Druid Hills High Scho Campus Parent Portal A Series of R Worksh U14 Boys

NEW REPUBLIC

## The Weird and Very Real World of Excel-Error Research

The Rogoff-Reinhart blunder is a prominent example of a very common problem

BY ROBERT LONG | April 18, 2013

They're calling it the "Excel Error Heard Round the World": Kenneth Rogoff and Carmen Reinhart's widely cited paper about the relationship between public debt and economic growth was revealed Monday to have grossly misstated economic growth for high-debt countries, all because of a forehead-smackingly simple error in an Excel spreadsheet. ("It is sobering that such an error slipped into one of our papers despite our best efforts to be consistently careful," the paper's authors said on Wednesday.)

<https://newrepublic.com/article/112951/rogoff-reinhart-and-world-excel-error-research>

# Timeline Reproducible Research & Transparency

YEAR	Event
2012	Begley and Ellis reviewed 53 "landmark" studies and only 6 (11%) had the scientific findings confirmed. <sup>9</sup>
2013	Center for Open Science launches & by 2014 the Open Science Framework has 7000 users with more than 45,000+ and over 15 institutions by 2017 <sup>10</sup>
2014	NIH publishes their guidelines for addressing reproducibility <sup>11</sup>
2015	The Open Science Collaboration reports that they were only able to replicate between 1/3 to 1/2 of the results from 100 studies <sup>12</sup>

9. <http://www.nature.com/nature/journal/v483/n7391/full/483531a.html>

10. <https://cos.io/about/brief-history-cos-2013-2017/> & <https://osf.io/>

11. <https://www.nih.gov/research-training/rigor-reproducibility>

12. Science, 28 Aug 2015: Vol. 349, Issue 6251, aac4716; DOI: 10.1126/science.aac4716; <http://science.sciencemag.org/content/349/6251/aac4716>

# Wide-Spread Gene Name Errors



COMMENT | OPEN ACCESS

## Gene name errors are widespread in the scientific literature

Mark Ziemann, Yotam Eren and Assam El-Osta

Genome Biology 2016 17:177 | <https://doi.org/10.1186/s13059-016-1044-7> | © The Author(s). 2016

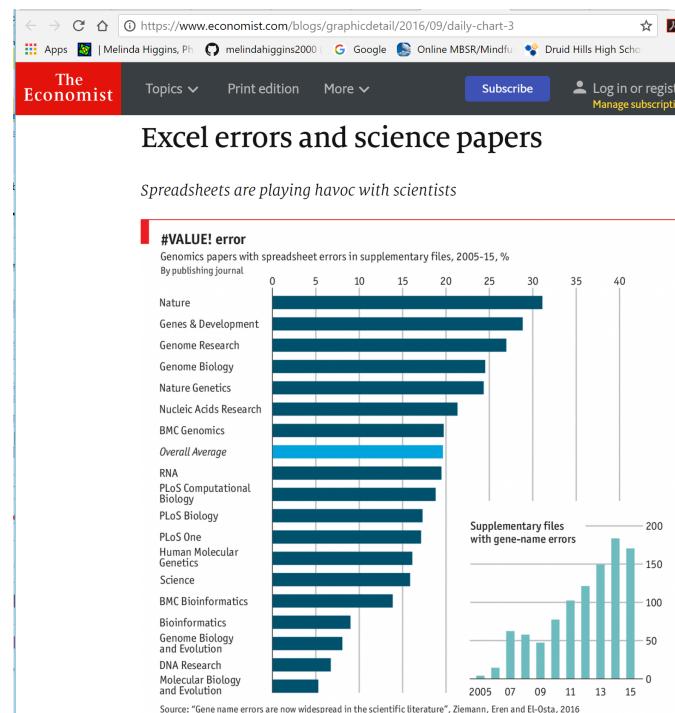
Published: 23 August 2016

### Abstract

The spreadsheet software Microsoft Excel, when used with default settings, is known to convert gene names to dates and floating-point numbers. A programmatic scan of leading genomics journals reveals that approximately one-fifth of papers with supplementary Excel gene lists contain erroneous gene name conversions.

<https://genomebiology.biomedcentral.com/articles/10.1186/s13059-016-1044-7>

# Wide-Spread Gene Name Errors



<https://www.economist.com/blogs/graphicdetail/2016/09/daily-chart-3>

# People

- Victoria Stodden <https://ischool.illinois.edu/people/faculty/vcs>
  - presentation on History of the Reproducibility Movement  
<https://web.stanford.edu/~vcs/talks/ICERM-Dec102012STODDEN.pdf>
  - co-author "Implementing Reproducible Research" book  
<https://www.crcpress.com/Implementing-Reproducible-Research/Stodden-Leisch-Peng/p/book/9781466561595>
- Roger Peng <http://www.biostat.jhsph.edu/~rpeng/index.html>
  - Associate Editor for Reproducible Research - Biostatistics Journal  
[https://academic.oup.com/biostatistics/pages/Editorial\\_Board](https://academic.oup.com/biostatistics/pages/Editorial_Board)
  - co-author "Implementing Reproducible Research" book  
<https://www.crcpress.com/Implementing-Reproducible-Research/Stodden-Leisch-Peng/p/book/9781466561595>

# People

- John P.A. Ioannidis [https://profiles.stanford.edu/john-ioannidis?  
tab=publications](https://profiles.stanford.edu/john-ioannidis?tab=publications)
  - Professor of Medicine and of Health Research and Policy at Stanford University School of Medicine and a Professor of Statistics at Stanford University School of Humanities and Sciences
- Christopher Gandrud [https://www.iq.harvard.edu/people/christopher-  
gandrud](https://www.iq.harvard.edu/people/christopher-gandrud)
  - research fellow at IQSS (Institute for Quantitative Social Science)
  - Book Author "Reproducible Research with R and RStudio"  
[https://www.crcpress.com/Reproducible-Research-with-R-and-R-  
Studio/Gandrud/p/book/9781466572843](https://www.crcpress.com/Reproducible-Research-with-R-and-R-Studio/Gandrud/p/book/9781466572843)

# People

- Yihui Xie
  - software engineer for RStudio <https://www.rstudio.com/about/>
  - author of "Dynamic Documents with R and knitr" <https://www.crcpress.com/Dynamic-Documents-with-R-and-knitr/Xie/p/book/9781482203530>
  - author of "Bookdown: Authoring Books and Technical Documents with R Markdown" book <https://www.crcpress.com/bookdown-Authoring-Books-and-Technical-Documents-with-R-Markdown/Xie/p/book/9781138700109> and bookdown R package <https://cran.r-project.org/web/packages/bookdown/index.html>
  - author of blogdown R package <https://cran.r-project.org/web/packages/blogdown/index.html>

# People

- Friedrich Leisch
  - Professor of Applied Statistics at the University of Natural Resources and Life Sciences, Vienna
  - developer of Sweave for creating dynamic reports  
<https://leisch.userweb.mwn.de/Sweave/>
  - co-author "Implementing Reproducible Research" book  
<https://www.crcpress.com/Implementing-Reproducible-Research/Stodden-Leisch-Peng/p/book/9781466561595>

# Books on Reproducibility and Tools of the Trade

## Image



## Book

Implementing Reproducible Research by Victoria Stodden, Friedrich Leisch, Roger D. Peng <https://www.crcpress.com/Implementing-Reproducible-Research/Stodden-Leisch-Peng/p/book/9781466561595>



Dynamic Documents with R and knitr (Chapman & Hall/CRC The R Series) 1st Edition by Yihui Xie <https://www.crcpress.com/Dynamic-Documents-with-R-and-knitr/Xie/p/book/9781482203530>



bookdown: Authoring Books and Technical Documents with R Markdown by Yihui Xie <https://www.crcpress.com/bookdown-Authoring-Books-and-Technical-Documents-with-R-Markdown/Xie/p/book/9781138700109> & read online <https://bookdown.org/yihui/bookdown/>

---

# more books

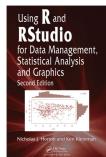
Image



Jennifer (Jenny)  
Bryan  
jennylbc

Book

Happy Git and GitHub for the useR by Jenny Bryan; read online  
<http://happygitwithr.com/>



Using R and RStudio for Data Management, Statistical Analysis, and Graphics, Second Edition by Nicholas J. Horton & Ken Kleinman  
<https://www.crcpress.com/Using-R-and-RStudio-for-Data-Management-Statistical-Analysis-and-Graphics/Horton-Kleinman/p/book/9781482237368>; also see [Project MOSAIC](#), <http://mosaic-web.org/>



ModernDive: An Introduction to Statistical and Data Sciences via R by Chester Ismay and Albert Y. Kim; read online <https://ismayc.github.io/moderndiver-book/> & Getting used to R, RStudio, and R Markdown by Chester Ismay <https://ismayc.github.io/rbasics-book/>

.

... and lots more ... see <https://bookdown.org/>

# Literate Programming > Dynamic Documentation > [R]Markdown

YEAR      Event

1992	"Literate Programming" is introduced by Donald Knuth as "that (which) combines a programming language with a documentation language, thereby making programs more robust, more portable, more easily maintained, and arguably more fun to write than programs that are written only in a high-level language. <b>The main idea is to treat a program as a piece of literature, addressed to human beings rather than to a computer.</b> " <a href="http://www-cs-faculty.stanford.edu/~knuth/lp.html">http://www-cs-faculty.stanford.edu/~knuth/lp.html</a>
2002	Friedrich Leisch introduces SWEAVE a program for "Dynamic generation of statistical reports using literate data analysis" <a href="https://leisch.userweb.mwn.de/Sweave/">https://leisch.userweb.mwn.de/Sweave/</a>

# Literate Programming > Dynamic Documentation > [R]Markdown

YEAR      Event

2004	John Gruber created the <b>Markdown</b> language in 2004 in collaboration with Aaron Swartz - their goal was to "write using an easy-to-read, easy-to-write plain text format, and optionally convert it to structurally valid XHTML (or HTML)" <a href="https://daringfireball.net/projects/markdown/">https://daringfireball.net/projects/markdown/</a>
2012	Yihui Xie releases <b>knitr</b> R package released - <b>knitr</b> was inspired by <b>SWEAVE</b>
2014	<b>rmarkdown</b> R package released - extends <b>Markdown</b> to work with R/RStudio environment

# The WEB System by Donald Knuth

The first published literate programming environment was **WEB**. Donald Knuth introduced it back in 1981 combining his TeX typesetting system with the Pascal programming language.

## Literate Programming

---

**Donald E. Knuth**

Computer Science Department, Stanford University, Stanford, CA 94305, USA

*"I chose the name **WEB** partly because it was one of the few three-letter words of English that hadn't already been applied to computers. But as time went on, I've become extremely pleased with the name, because I think that a complex piece of software is, indeed, best regarded as a web that has been delicately pieced together from simple materials. ... If we express a program as a web of ideas, we can emphasize its structural properties in a natural and satisfying way."<sup>13</sup>*

13. <http://www.literateprogramming.com/knuthweb.pdf>

# More Literate Programming Tools

Since WEB was introduced in 1981, many other programs implementing literate programming have emerged over time including:

- CWEB also created by Donald Knuth with Silvio Levy which was adapted for the C and C++ programming language instead of Pascal
- Axiom developed by IBM
- Noweb
- Literate
- Funnel WEB
- Molly
- Codnar
- Jupyter Notebook (formerly IPython Notebook)
- R Notebooks

# SWEAVE by Friedrich Leisch

## Sweave: Dynamic Generation of Statistical Reports Using Literate Data Analysis

What is Sweave?

*"Sweave is a tool that allows to embed the R code for complete data analyses in latex documents. The purpose is to create dynamic reports, which can be updated automatically if data or analysis change. Instead of inserting a prefabricated graph or table into the report, the master document contains the R code necessary to obtain it. When run through R, all data analysis output (tables, graphs, etc.) is created on the fly and inserted into a final latex document. The report can be automatically updated if data or analysis change, which allows for truly reproducible research."<sup>14</sup>*

14. Friedrich Leisch. Sweave: Dynamic generation of statistical reports using literate data analysis. In Wolfgang Härdle and Bernd Rönz, editors, Compstat 2002 - Proceedings in Computational Statistics, pages 575-580. Physica Verlag, Heidelberg, 2002. ISBN 3-7908-1517-9.

# The next evolution <- knitr



In 2012 Yihui Xie, created and released the **knitr** package for R to extend the capabilities of SWEAVE beyond LaTeX.

*"The knitr package was designed to be a transparent engine for dynamic report generation with R, solve some long-standing problems in Sweave, and combine features in other add-on packages into one package."*<sup>15</sup>

15. <https://yihui.name/knitr/>

# The next evolution <- ... + rmarkdown



- In 2014, RStudio released **rmarkdown** to extend the **markdown** language originally intended to write documents for the "web" (*i.e.* *HTML*).<sup>16</sup>
- **rmarkdown** leverages Pandoc ("universal document converter")<sup>17</sup> to convert between formats: from HTML (readable by web browsers) to DOC (such as from Microsoft Word or Google Docs) to ODT (Libre Office) to PDF (portable document format) to others like EPUB (e-books), HTML5 slide shows (slidy, ioslides), and TeX based documents and slides (Beamer).

16. <https://daringfireball.net/projects/markdown/syntax>

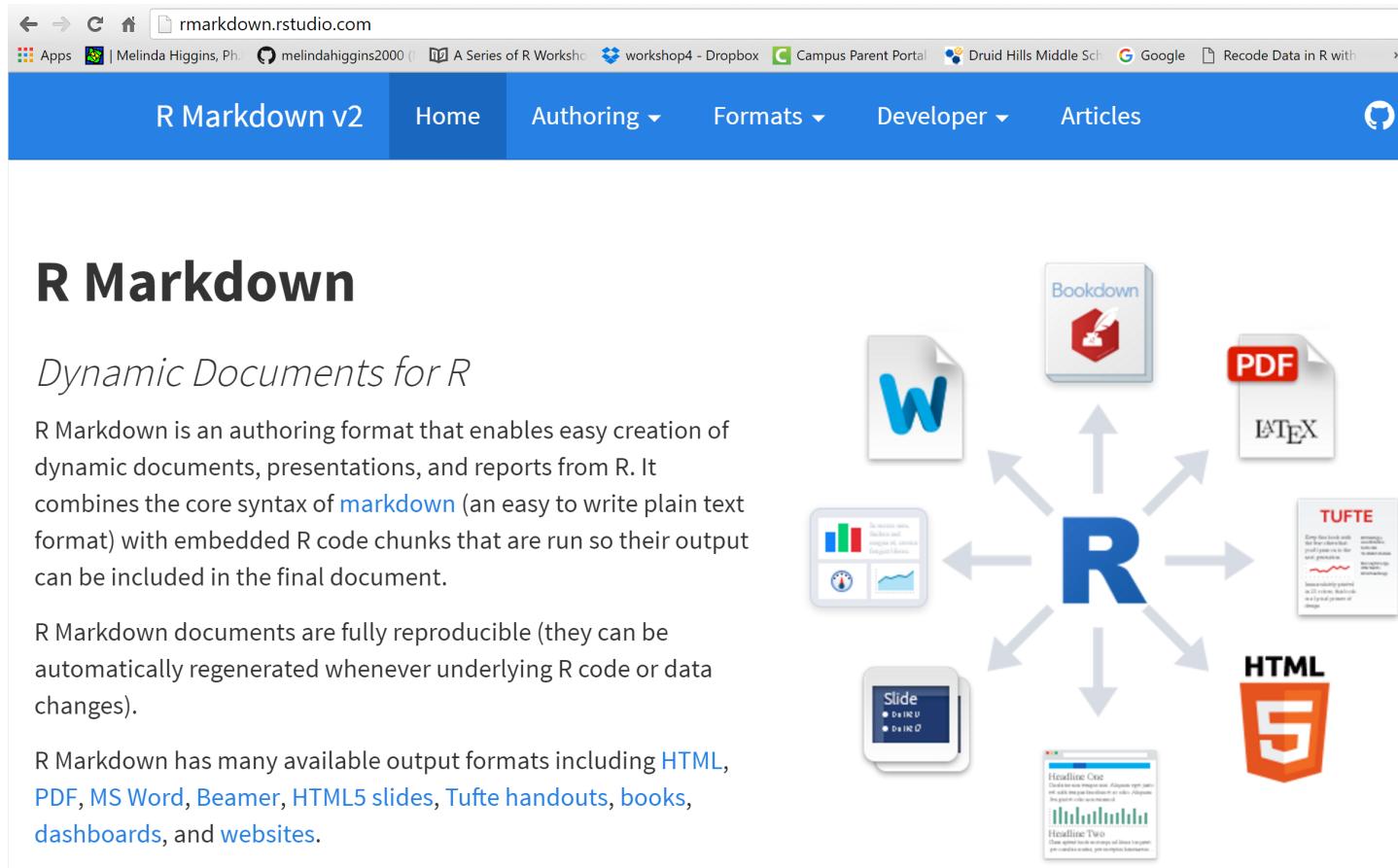
17. <http://pandoc.org/index.html>

# Pandoc <https://pandoc.org/>

...often called the *Swiss-Army knife* for converting files from one format to another. Pandoc can convert documents in markdown, reStructuredText, textile, HTML, DocBook, LaTeX, MediaWiki markup, TWiki markup, OPML, Emacs Org-Mode, Txt2Tags, Microsoft Word docx, LibreOffice ODT, EPUB, or Haddock markup to

- HTML formats: XHTML, HTML5, Slidy, reveal.js, Slideous, S5, DZSlides.
- Word processor formats: Microsoft Word docx, OpenOffice/LibreOffice ODT, OpenDocument XML
- Ebooks: EPUB version 2 or 3, FictionBook2
- Documentation formats: DocBook, TEI Simple, GNU TexInfo, Groff man pages, Haddock markup
- Page layout formats: InDesign ICML
- Outline formats: OPML
- TeX formats: LaTeX, ConTeXt, LaTeX Beamer slides
- PDF via LaTeX
- Lightweight markup formats: Markdown (including CommonMark), reStructuredText, AsciiDoc, MediaWiki markup, DokuWiki markup, Emacs Org-Mode, Textile
- Custom formats: written in lua.

# The RStudio "HUB"



The screenshot shows the R Markdown v2 website at [rmarkdown.rstudio.com](http://rmarkdown.rstudio.com). The top navigation bar includes links for Home, Authoring, Formats, Developer, and Articles. Below the navigation, a large section titled "R Markdown" is displayed, with a subtitle "Dynamic Documents for R". A text block explains that R Markdown is an authoring format for dynamic documents, presentations, and reports. It highlights its reproducibility and the ability to embed R code chunks. Another section lists available output formats: HTML, PDF, MS Word, Beamer, HTML5 slides, Tufte handouts, books, dashboards, and websites. To the right of the text, there is a diagram illustrating the central role of R in generating various document types.

**R Markdown**

*Dynamic Documents for R*

R Markdown is an authoring format that enables easy creation of dynamic documents, presentations, and reports from R. It combines the core syntax of [markdown](#) (an easy to write plain text format) with embedded R code chunks that are run so their output can be included in the final document.

R Markdown documents are fully reproducible (they can be automatically regenerated whenever underlying R code or data changes).

R Markdown has many available output formats including [HTML](#), [PDF](#), [MS Word](#), [Beamer](#), [HTML5 slides](#), [Tufte handouts](#), [books](#), [dashboards](#), and [websites](#).

**R**

Bookdown, PDF, L<sup>A</sup>T<sub>E</sub>X, TUFTE, HTML5, Slide

# Reproducible Principles - Process & Structure

- Organization
- Clear Documentation
- Standardized
- Centralized
- Efficiency

# 10 Simple Rules for Reproducible Computational Research<sup>18</sup>

1. For every result, keep track of how it was produced
2. Avoid Manual Data Manipulation Steps
3. Archive the Exact Version of All External Programs Used
4. Version Control All Custom Scripts
5. Record All Intermediate Results, When Possible in Standardized Formats
6. For Analyses that include randomness, note underlying random seeds
7. Always Store Raw Data Behind Plots
8. Generate Hierarchical Analysis Output, Allowing Layers of Increasing Detail to be Inspected
9. Connect Textual Statements to Underlying Results
10. Provide Public Access to Scripts, Runs, and Results

18. Sandve, G.K.; Nekrutenko, A.; Taylor, J.; Hovig, E. (2013) "Ten Simple Rules for Reproducible Computational Research" PLOS Computational Biology, 9(10).<https://doi.org/10.1371/journal.pcbi.1003285>

# Standard Practices

*Think about your own work...*

- What do you want to automate?
- What could you re-use?
  - code, files, formatting, graphics, logos, header, footer, boilerplate
- What should you share with your team?
- What do you find yourself doing over and over?
  - correcting or reformatting
- If you won the lottery today (and left your job), what do you need to tell your replacement so they can pick up where you left off and complete your current tasks?

# Journalism - 538.com

538.com <http://fivethirtyeight.com/> hosts stories and opinion pieces covering poll analyses, politics, economics, health, popular culture, and sports. The founder, Nate Silver, and the 538 team are best known for their political polling and forecasting during the United States Presidential and related elections since 2008. ESPN now owns 538.com (as of 2013) retaining Nate Silver as the Editor-in-Chief.

Most of their articles provide references and links to the original data sources plus details on how their figures, analyses and statistical models were developed. They also host the data, code and details behind their analyses on Github <https://github.com/fivethirtyeight/>.

We will work with some of these datasets in our exercises later in this course and work with the `fivethirtyeight` R package <https://cran.r-project.org/web/packages/fivethirtyeight/>.

# Telling Stories with Data

Andrew Flowers (economist, data scientist, journalist and former writer for [fivethirtyeight.com](http://fivethirtyeight.com)) presented "Finding and Telling Stories with R" at the 2017 RStudio Conference (Orlando, FL).

The webinar recording of his presentation is available online  
[https://www.rstudio.com/resources/videos/finding-and-telling-stories-with-r/.](https://www.rstudio.com/resources/videos/finding-and-telling-stories-with-r/)

In his presentation, he highlights the various aspects of "data journalism" and importance of workflow, data processing and transparency in analysis and communication - all key aspects of reproducibility. Andrew Flowers is also a contributor to the **fivethirtyeight** R package.

# Transparency - Journal of Biostatistics

"Our reproducible research policy is for papers in the journal to be kite-marked **D** if the data on which they are based are freely available, **C** if the authors' code is freely available, and **R** if both data and code are available, and our Associate Editor for Reproducibility is able to use these to reproduce the results in the paper. Data and code are published electronically on the journal's website as Supplementary Materials."

[https://academic.oup.com/biostatistics/pages/General\\_Instructions](https://academic.oup.com/biostatistics/pages/General_Instructions)

Example of an article marked **R**:

- Air pollution and health in Scotland: a multicity study; by Duncan Lee; Claire Ferguson ; and Richard Mitchell; *Biostatistics*, Volume 10, Issue 3, 1 July 2009, Pages 409–423, <https://doi.org/10.1093/biostatistics/kxp010>

# Speed - 2001 outbreak of *E.Coli 0104:H4*

In 2001 there was an outbreak of *E.Coli 0104:H4* that killed 50 people in Europe  
<http://dx.doi.org/10.5524/100001>.

Researchers at BGI (*formally the Beijing Genomics Institute*) worked in collaboration with the Medical Center in Hamburg-Eppendorf to rapidly sequence the genome of the pathogen. Given the severity of the outbreak, the team announced and released the genome via Twitter to the world-wide community of microbial genomicists.

A Github repository was established <https://github.com/ehec-outbreak-crowdsourced/BGI-data-analysis/wiki> to "crowdsource" analysis and research to find a treatment.

People started contributing their work in under 24 HOURS and within 5 DAYS!! a bacterial agent was proposed to kill the pathogen. *This case, highlights the importance of these methods and work practices not only for speed and efficiency but also in rapidly addressing problems and developing solutions that can save lives.*

# Documentation

- main component is text
- well written
- good organization and flow
- easily accessible
- understood by team members at all levels
- code + text + figures combined [e.g. literate programming]
- at end, formatting styles applied via "markup/markdown"

# Organization

- projects grow
- supporting documentation and files - numerous
- relationships change and can grow more complex
- need file organization and naming schemes
- file names should be:
  - readable by the computer, easy to search, easy to sort (especially by date and author if needed)
  - human readable with logical naming schemes and contain enough info so human knows what is in the file/what the file is for
  - and short enough to be reasonably manageable
- consider user-based access and security (partitioned by "need to know" *[users with editing and write permissions versus users with read-only access]*)

Research Compendium Example <https://github.com/ropensci/rrrpkgs>

# Automation

- at a minimum, a diagram or instructions for workflow should be documented on how the components are to be assembled for your final product
- write code/scripts to automate
  - data raw to processed output
  - creating and removing temporary files
  - creating tables, figures, other components
  - assembling the components into final documents, products
  - rendering documents into multiple/desired formats

# Dissemination - Why?

- store and share your data and code – even if it is only for your future reference
- sometimes expectation/requirement of funding agency, publisher
- increased visibility, you as source - default subject matter expert
- speed of collaboration - faster advancement of science, knowledge
- good will with community/public

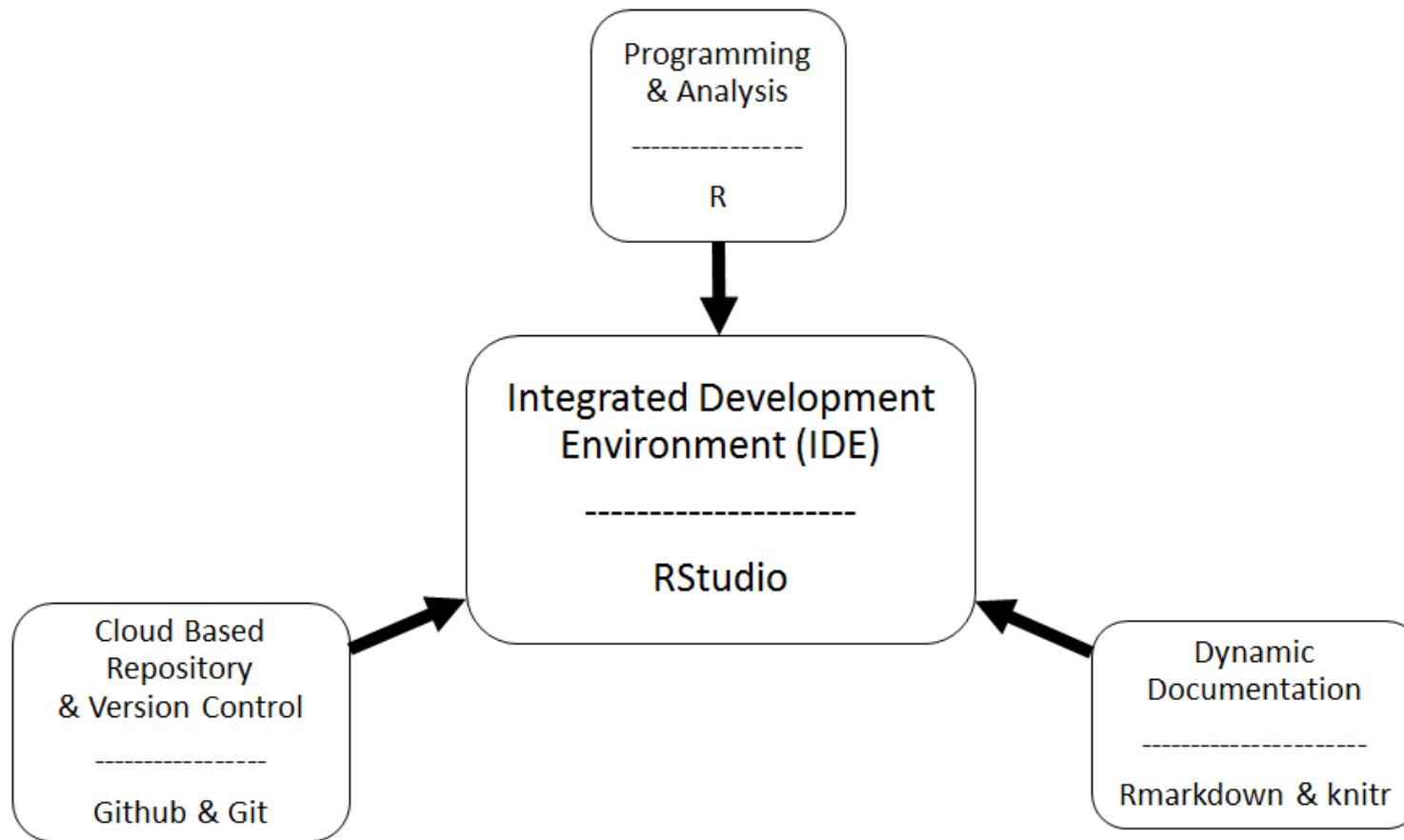
# Dissemination - How?

- Cloud-based "File Storage"
  - Dropbox <https://www.dropbox.com/>
  - Google drive <https://www.google.com/drive/>
  - Github (better with version control and tracking) <https://github.com/>
- Data repositories
  - GenBank <https://www.ncbi.nlm.nih.gov/genbank/>
  - PDB <https://www.rcsb.org/pdb/home/home.do>
- In addition to Github
  - Bitbucket <https://bitbucket.org/>
  - Dryad <http://datadryad.org/>
  - Figshare <https://figshare.com/>
  - Zenodo <https://zenodo.org/>

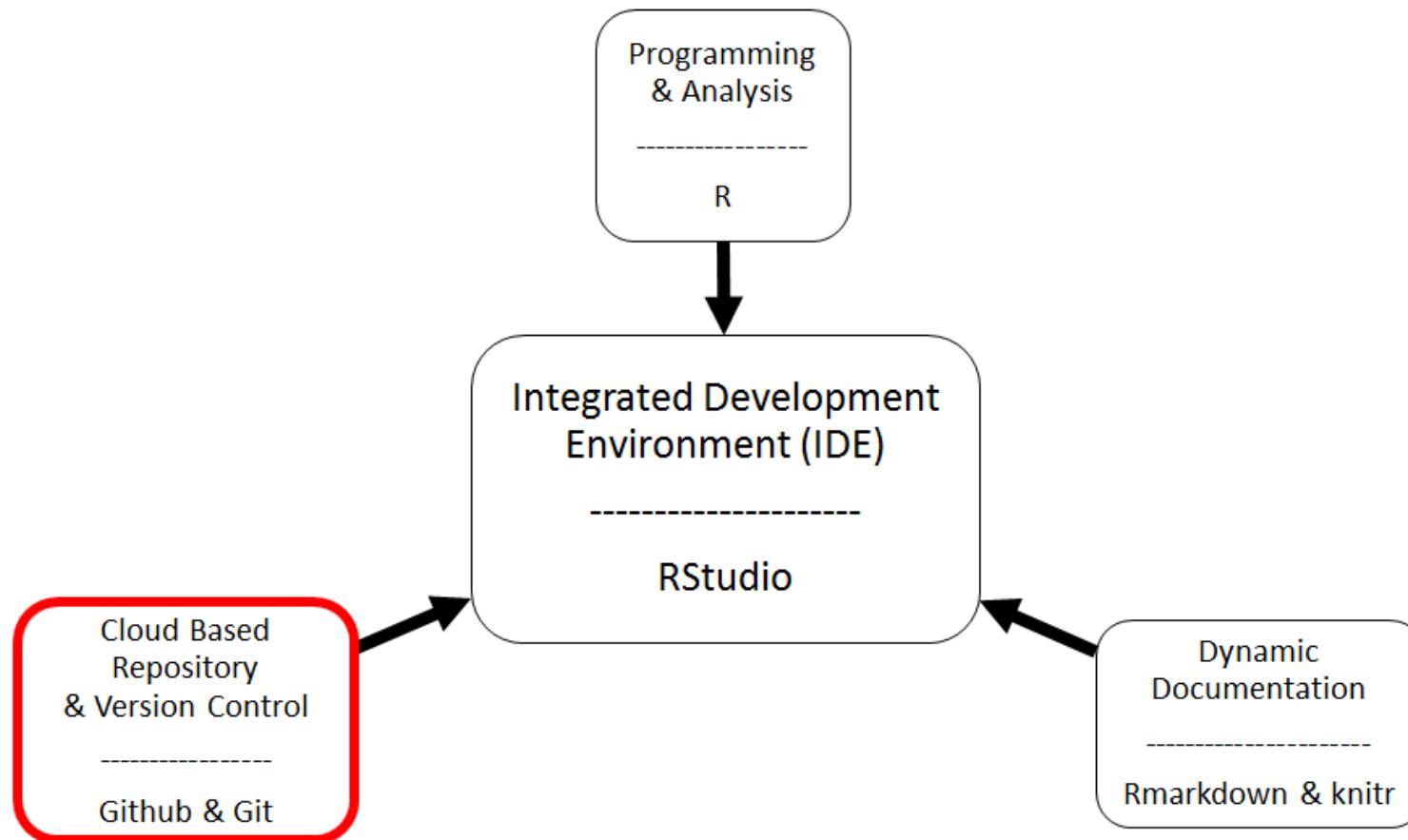
# Dissemination - Who? (e.g. *stakeholders*)

- Yourself
- Your organization - internal reports
- Journals - articles, manuscripts
- Books
- Blogs/Websites
- RSS feeds
- Rpubs <https://rpubs.com/>
- Gitbook <https://www.gitbook.com/>
- Bookdown <https://bookdown.org/yihui/bookdown/>

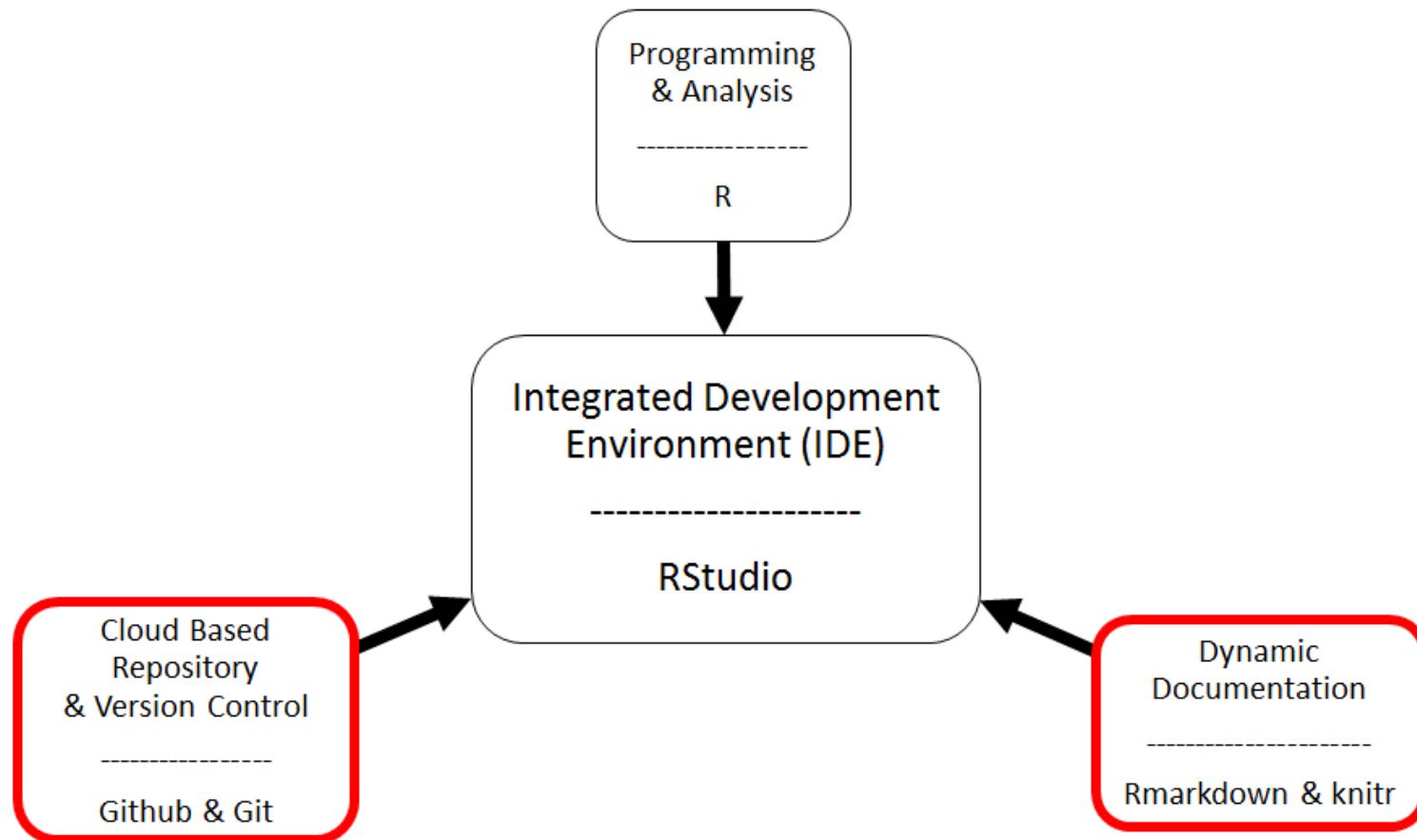
# The Big Picture



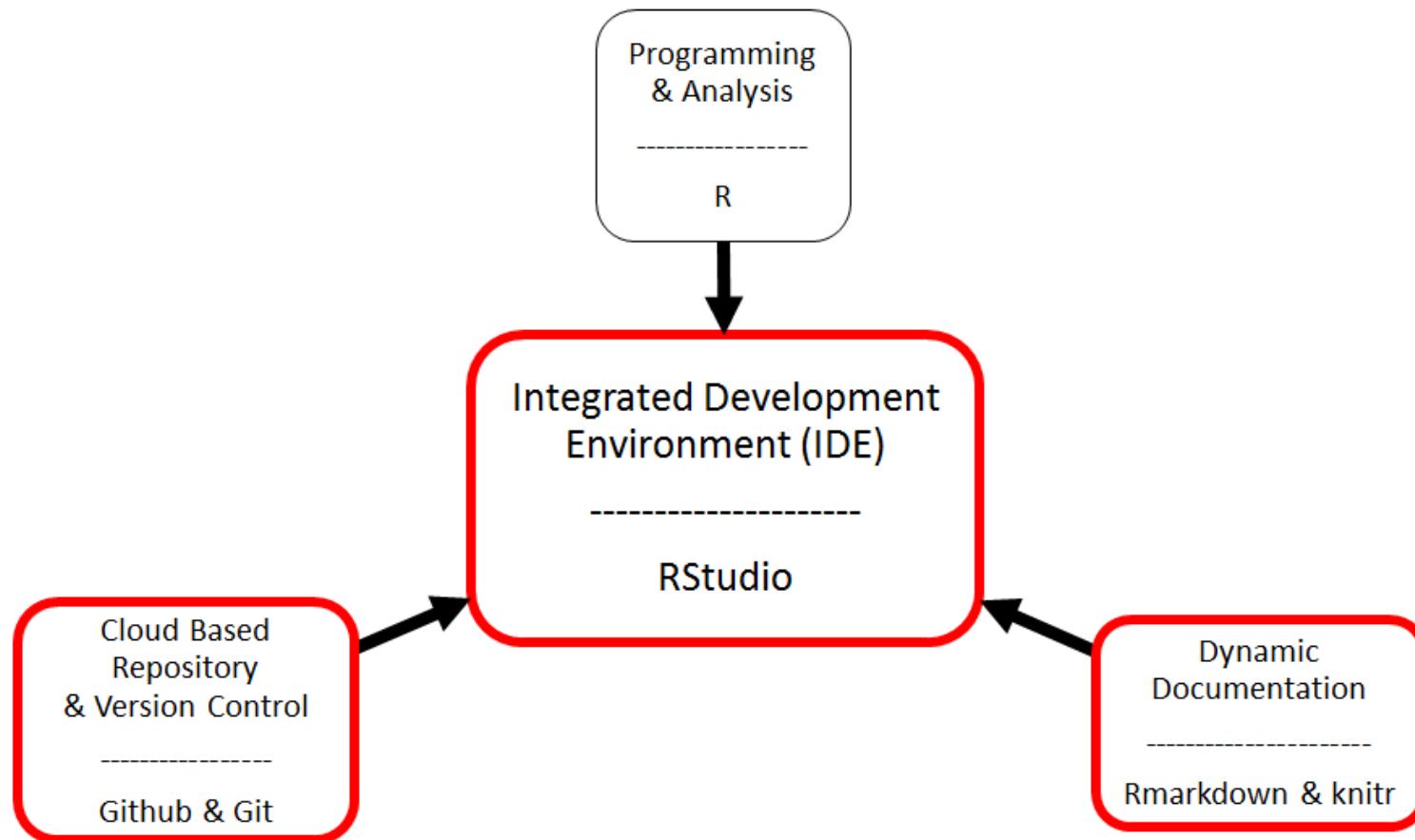
# The Big Picture



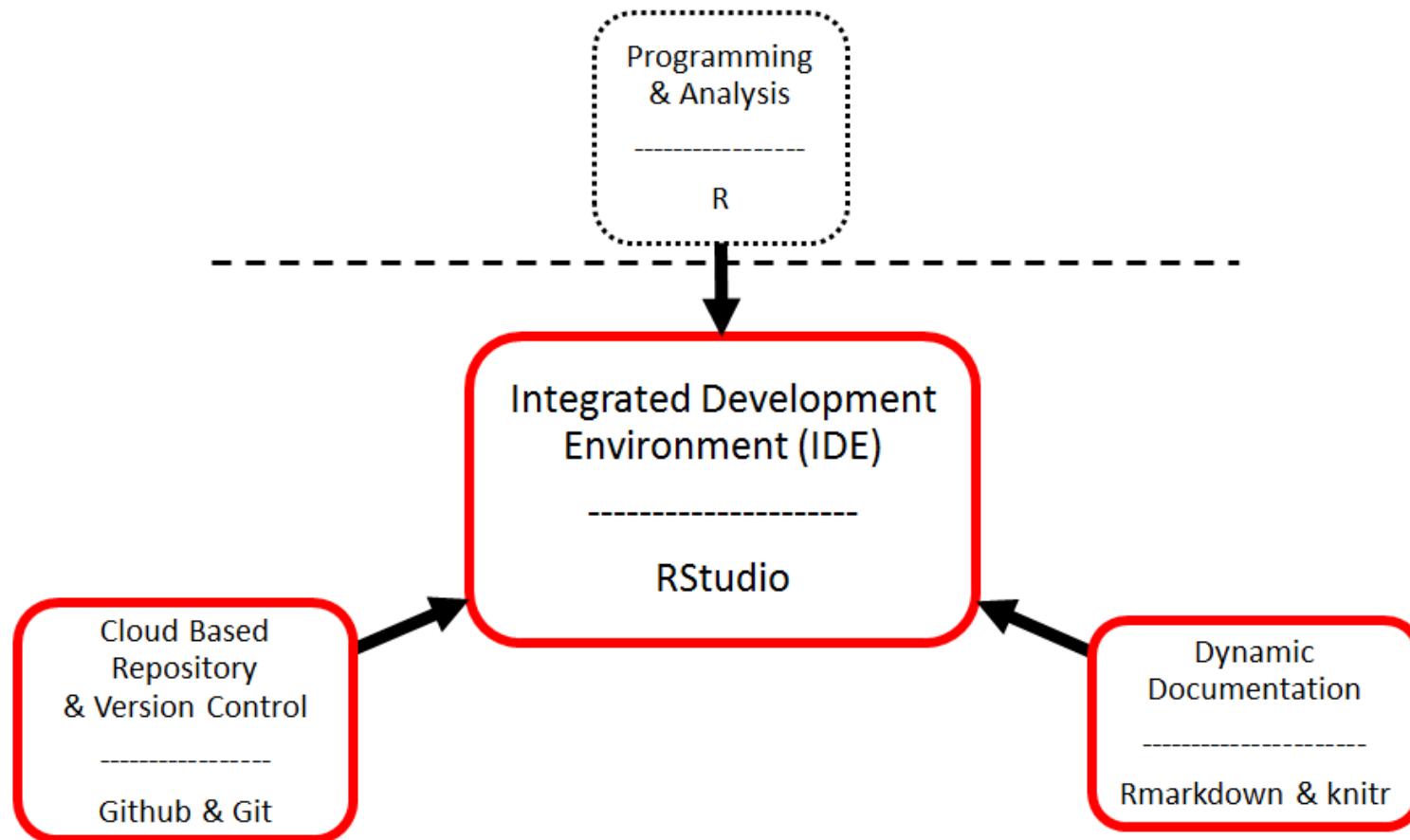
# The Big Picture



# The Big Picture



# The Big Picture



# Next in Lesson 02 ...

Literate Programming

&

Dynamic Documentation