

Building Research Software

- Damien Irving
- Kate Hertweck
- Luke Johnston
- Sara Mahallati
- Joel Ostblom
- Charlotte Wickham
- Greg Wilson

1. Tell us about yourself.

What are your qualifications for writing this book?

We are all research software engineers, practicing data scientists, and/or professional educators in those fields.

Do you have any unique characteristics or experiences that will make you stand out as the author?

One of us co-founded Software Carpentry, a non-profit organization that has taught programming and data analysis skills to over 50,000 people worldwide since 2010.

2. Tell us about the book.

What is the technology or idea that you're writing about?

This book introduces the core skills researchers need to develop robust software that others can use, to share their work with others, and to be productive as part of a research software team. We assume readers have some basic programming knowledge, and build on that to cover:

- using the Unix shell to manage work and make it repeatable
- effective use of version control
- building reusable software tools
- automating workflows
- configuring software
- testing and error handling
- creating productive, inclusive teams
- documenting, packaging, and releasing software for use by the wider research community

All of this material has been used and refined in workshops, some of it by multiple instructors over many years.

Why is it important now?

Software is now as important to research as telescopes and test tubes, but most researchers are never taught anything more than basic programming skills. As a result, they takes days or weeks to do what could be done in minutes or hours, have no idea how reliable their software is, and cannot share it with others. Filling in these gaps will help researchers get more done in less time and with less pain.

Tell us roughly how it works or what makes it different from its alternatives.

- Each set of tools and practices is covered in one chapter (except for the Unix shell and version control, which are each given two chapters to make exposition clearer).
- A running example demonstrates how to transform a simple data analysis problem into a reproducible, reusable project.

What type of book are you planning to write?

- The book is a tutorial suitable both for self-study and as a textbook in a one-semester course.
- It could fit into either the “In Action” or “In Practice” series.
- It includes exercises and solutions, all of which use real-world tools and examples.

3. Please give us 5 or 6 representative tasks in the domain of your book.

1. Create a Unix shell script to re-run a series of data analysis steps on a new input file with a single command.
2. Use branches, pull requests, and code review to add a feature to the Python programs invoked by that script.
3. Create a Makefile that automatically regenerates results and plots whenever a new data file is added to the project.
4. Add error handling and logging to the scripts so that they can be put into production.
5. Use style-checking tools to ensure the Python code is well written, and then release it as an installable package.

6. Adopt a license, a code of conduct, and efficient decision-making processes to ensure the project runs smoothly.
7. Document, package, and release software for use by the wider research community.

4. The minimally-qualified reader (MQR).

At the beginning of the book, we introduce three personas who characterize our audience:

Amira Khan completed a master’s in library science five years ago and has since worked for a small aid organization. She did some statistics during her degree, and has learned some R and Python by doing data science courses online, but has no formal training in programming. Amira would like to tidy up the scripts, data sets, and reports she has created in order to share them with her colleagues. These lessons will show her how to do this and what “done” looks like.

Jun Hsu completed a 4-month data science bootcamp last year after doing a PhD in Geology and now works for a company that does forensic audits. He uses a variety of machine learning and visualization packages, and would now like to turn some of his own work into an open source project. This book will show him how such a project should be organized and how to encourage people to contribute to it.

Sami Virtanen became a competent programmer during a bachelor’s degree in applied math and was then hired by the university’s research computing center. The kinds of applications they are being asked to support have shifted from fluid dynamics to data analysis; this guide will teach them how to build and run data pipelines so that they can pass those skills on to their users.

5. Q&A

What are the three or four most commonly-asked questions about this technology?

We wrote this book in part because the people who need it don’t even know what questions to ask. “How do I write a general rule in a Makefile?” presupposes that you know that Make exists. Similarly, people don’t ask, “How should a small group run a meeting?” because it doesn’t occur to them that running meetings is an improvable skill.

What are others interested in this topic asking in forums?

All of our chapters are driven by questions we have answered (some of them many times) in workshops or in discussions with colleagues. Our answers for

technical subjects frequently draw on Stack Overflow, particularly for topics like continuous integration and packaging; our answers for community management and project leadership draw on extensive experience in open source.

6. Tell us about the competition and the ecosystem.

What other books are available on this topic?

The Turing Way (<https://the-turing-way.netlify.app/>), produced by the Turing Institute in the UK, is the only up-to-date resource we know of with similar breadth. Other books cover only one topic in more depth than our audience needs (e.g., Ray & Ray's *Unix and Linux: Visual QuickStart Guide*), are aimed at commercial software developers (any number of books on testing, coding style, and dev ops), or go into data science rather than telling readers how to build software (e.g., Zhang's upcoming *A Tour of Data Science*).

How does the proposed book compare to them?

1. Our book has broader and deeper coverage, and includes exercises with solutions.
2. Our focus on programming best practices (e.g. principles like writing modular, reusable, testable code). Others teach basic Python and R syntax and how to complete isolated tasks with particular libraries, but today, people can pick up that information from a Google search. Our content is the bigger picture you can't glean from online code snippets or Stack Overflow answers.
3. We go through the process of building a data science workflow/package from scratch. Most other texts simply present information (like a reference book): we tackle a real data science problem from beginning to end, weaving in relevant information as we go.
4. Our book is a ready-to-go university semester course.

What resources would you currently recommend to someone wanting to learn this subject?

- The Software Carpentry and Data Carpentry lessons (for technical topics).
- Fogel's *Open Source Software* (for project management).

What are the most important web sites and companies associated with this topic?

- The Carpentries (<http://carpentries.org>)
- The Turing Way (<https://the-turing-way.netlify.app>)

Where do others interested in this topic gather online?

- The Carpentries have a very active online community (<http://carpentries.org>).
- The UK's Society of Research Software Engineering (<https://society-rse.org/>).

7. Book size and illustrations

- Approximate number of published pages: 470.
- Approximate number of diagrams and other graphics: 40.
- Approximate number of code listings: 150.

8. Contact information

Formal name

- Dr. Damien B. Irving
- Dr. Kate L. Hertweck
- Dr. Luke Johnston
- Sara Mahallati
- Joel Ostblom
- Dr. Charlotte Wickham
- Dr. Gregory V. Wilson

Name as it would appear on the cover

- Damien Irving
- Kate Hertweck
- Luke Johnston
- Sara Mahallati
- Joel Ostblom
- Charlotte Wickham
- Greg Wilson

Social media

- Damien Irving: <https://damienirving.github.io/>
- Kate Hertweck: <http://katehertweck.com/>
- Luke Johnston: <http://lukewjohnston.com/>
- Sara Mahallati: <https://scholar.google.ca/citations?user=8zCoa-UAAAAJ>
- Joel Ostblom: <https://joelostblom.com/>
- Charlotte Wickham: <https://www.cwick.co.nz/>
- Greg Wilson: <http://third-bit.com/>

9. Schedule

First draft complete.

10. Table of Contents

- 1 Introduction
 - 1.1 What's the big picture?
 - 1.2 Who are these lessons for?
 - 1.3 Syllabus
 - 1.4 Acknowledgments
- 2 The Basics of the Unix Shell
 - 2.1 Exploring Files and Directories
 - 2.2 Moving Around
 - 2.3 Creating New Files and Directories
 - 2.4 Moving Files and Directories
 - 2.5 Copy Files and Directories
 - 2.6 Deleting Files and Directories
 - 2.7 Wildcards
 - 2.8 Reading the Manual
 - 2.9 Combining Commands
 - 2.10 How Pipes Work
 - 2.11 Repeating Commands on Many Files
 - 2.12 Variable Names
 - 2.13 Redoing Things
 - 2.14 Creating New Filenames Automatically
 - 2.15 Summary
 - 2.16 Exercises
 - 2.17 Key Points
- 3 Going Further with the Unix Shell
 - 3.1 Creating New Commands
 - 3.2 Making Scripts More Versatile
 - 3.3 Turning Interactive Work into a Script
 - 3.4 Finding Things in Files

- 3.5 Finding Files
- 3.6 Configuring the Shell
- 3.7 Summary
- 3.8 Exercises
- 3.9 Key Points
- 4 Command Line Programs in Python
 - 4.1 Programs and Modules
 - 4.2 Handling Command-Line Options
 - 4.3 Documentation
 - 4.4 Counting Words
 - 4.5 Pipelining
 - 4.6 Positional and Optional Arguments
 - 4.7 Collating Results
 - 4.8 Writing Our Own Modules
 - 4.9 Plotting
 - 4.10 Summary
 - 4.11 Exercises
 - 4.12 Key Points
- 5 Git at the Command Line
 - 5.1 Setting Up
 - 5.2 Creating a New Repository
 - 5.3 Adding Existing Work
 - 5.4 Describing Commits
 - 5.5 Saving and Tracking Changes
 - 5.6 Synchronizing with Other Repositories
 - 5.7 Exploring History
 - 5.8 Restoring Old Versions of Files
 - 5.9 Ignoring Files
 - 5.10 Summary
 - 5.11 Exercises
 - 5.12 Key Points
- 6 Advanced Git
 - 6.1 What's a Branch?
 - 6.2 Creating a Branch
 - 6.3 What Curve Should We Fit?
 - 6.4 Verifying Zipf's Law
 - 6.5 Merging
 - 6.6 Handling Conflicts
 - 6.7 A Branch-Based Workflow
 - 6.8 Using Other People's Work
 - 6.9 Pull Requests
 - 6.10 Handling Conflicts in Pull Requests
 - 6.11 Summary
 - 6.12 Exercises
 - 6.13 Key Points
- 7 Automating Analyses

- 7.1 Updating a Single File
- 7.2 Managing Multiple Files
- 7.3 Updating Files When Programs Change
- 7.4 Reducing Repetition in a Makefile
- 7.5 Automatic Variables
- 7.6 Generic Rules
- 7.7 Defining Sets of Files
- 7.8 Documenting a Makefile?
- 7.9 Automating Entire Analyses
- 7.10 Summary
- 7.11 Exercises
- 7.12 Key Points
- 8 Program Configuration
 - 8.1 Configuration File Formats
 - 8.2 Matplotlib Configuration
 - 8.3 The Global Configuration File
 - 8.4 The User Configuration File
 - 8.5 Adding Command-Line Options
 - 8.6 A Job Control File
 - 8.7 Summary
 - 8.8 Exercises
 - 8.9 Key Points
- 9 Error Handling
 - 9.1 Exceptions
 - 9.2 Kinds of Errors
 - 9.3 Writing Useful Error Messages
 - 9.4 Reporting Errors
 - 9.5 Summary
 - 9.6 Exercises
 - 9.7 Key Points
- 10 Working in Teams
 - 10.1 Include Everyone
 - 10.2 Establish a Code of Conduct
 - 10.3 Include a License
 - 10.4 Planning
 - 10.5 Bug Reports
 - 10.6 Labeling Issues
 - 10.7 Prioritizing
 - 10.8 Meetings
 - 10.9 Making Decisions
 - 10.10 Handling Conflict
 - 10.11 Summary
 - 10.12 Exercises
 - 10.13 Key Points
- 11 Code Style, Review, and Refactoring
 - 11.1 Python Style

- 11.2 Order
- 11.3 Checking Style
- 11.4 Refactoring
- 11.5 Code Reviews
- 11.6 Python Features
- 11.7 Summary
- 11.8 Exercises
- 11.9 Key Points
- 12 Project Structure
 - 12.1 What is a project?
 - 12.2 What files should every project contain?
 - 12.3 How should I structure the rest of my project?
 - 12.4 How should I manage a mix of compiled programs and scripts?
 - 12.5 How should I document the software in a project?
 - 12.6 What should I document?
 - 12.7 How can I create a useful FAQ?
 - 12.8 How does documentation for data differ from documentation for code?
 - 12.9 How should I document the data in a project?
 - 12.10 How can I get started documenting my data?
 - 12.11 What data documentation should I create first?
 - 12.12 How should I manage data that can't be stored in version control?
 - 12.13 Summary
 - 12.14 Exercises
 - 12.15 Key Points
- 13 Testing
 - 13.1 Assertions
 - 13.2 Unit Testing
 - 13.3 Testing Frameworks
 - 13.4 Testing Floating-Point Values
 - 13.5 Testing Error Handling
 - 13.6 Integration Testing
 - 13.7 Regression Testing
 - 13.8 Test Coverage
 - 13.9 Continuous Integration
 - 13.10 When to Write Tests
 - 13.11 Summary
 - 13.12 Exercises
 - 13.13 Key Points
- 14 Python Packaging
 - 14.1 Creating a Python Package
 - 14.2 Virtual Environments
 - 14.3 Installing a Development Package
 - 14.4 What Installation Does
 - 14.5 Distributing Packages
 - 14.6 Documenting Packages
 - 14.7 Hosting Documentation Online

- 14.8 Summary
- 14.9 Exercises
- 14.10 Key Points
- 15 Finale
- References
- A License
- B Code of Conduct
- C Contributing
- D Glossary
- E Setting Up
- F Key Points
- G Solutions
- H YAML
- I Working Remotely