



Projects & Workflows: Introduction, Concepts, & Setting Up Your Data Toolbox

Harvard DataFest
2017

Bob Freeman, PhD, Dir Research TechOps, HBS

Daina Bouquin, Head Librarian, Harvard-Smithsonian Center for Astrophysics, FAS

Derek Miller, PhD, Asst Professor of English, FAS

Caroline Shamu, Asst Professor of BCMP, HMS

17 January 2017



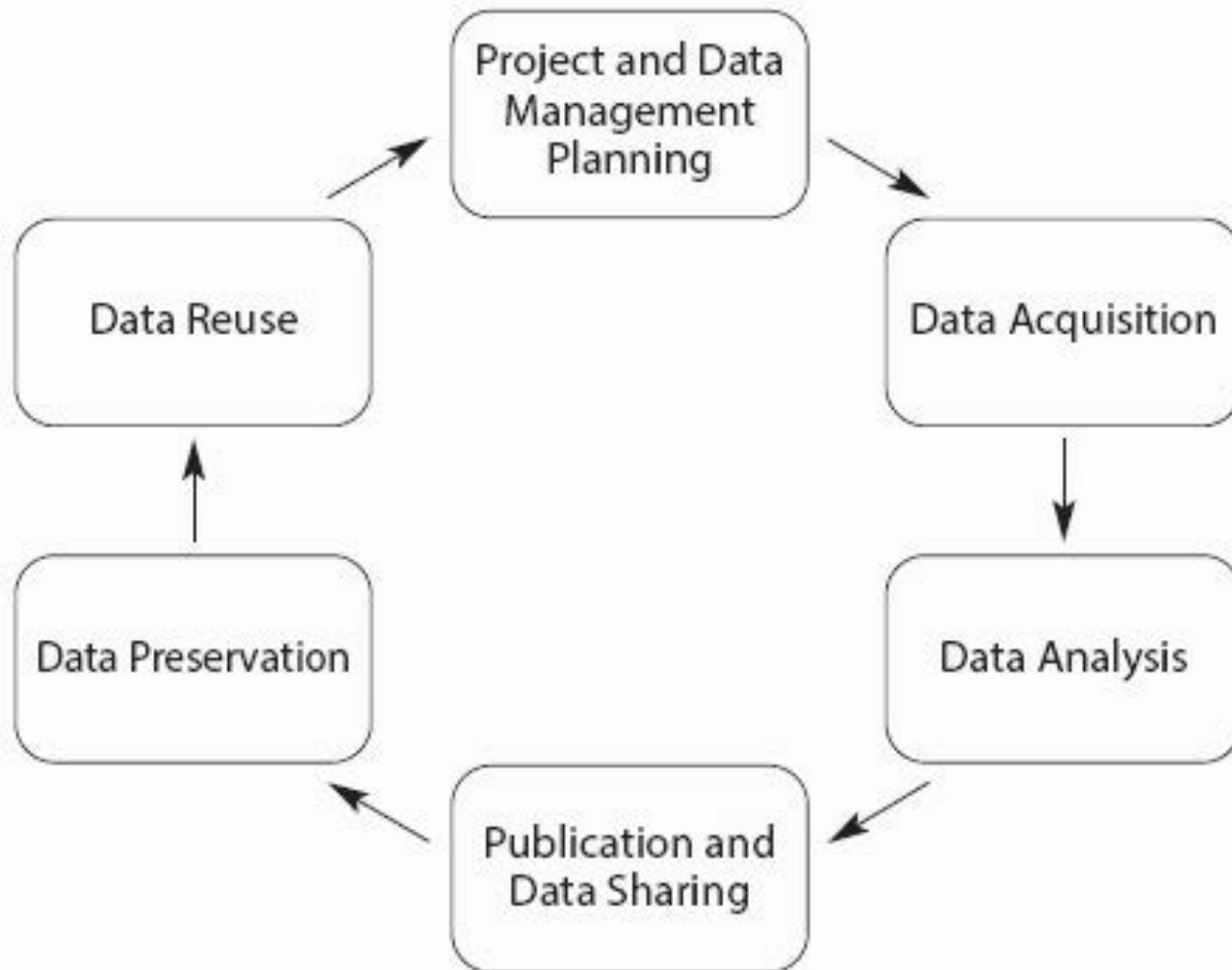
Overview

- Opening ideas for the afternoon session
- Why Plan?
- Why use frameworks and standards?
- What is in your Data Toolbox?

Planning & Organization

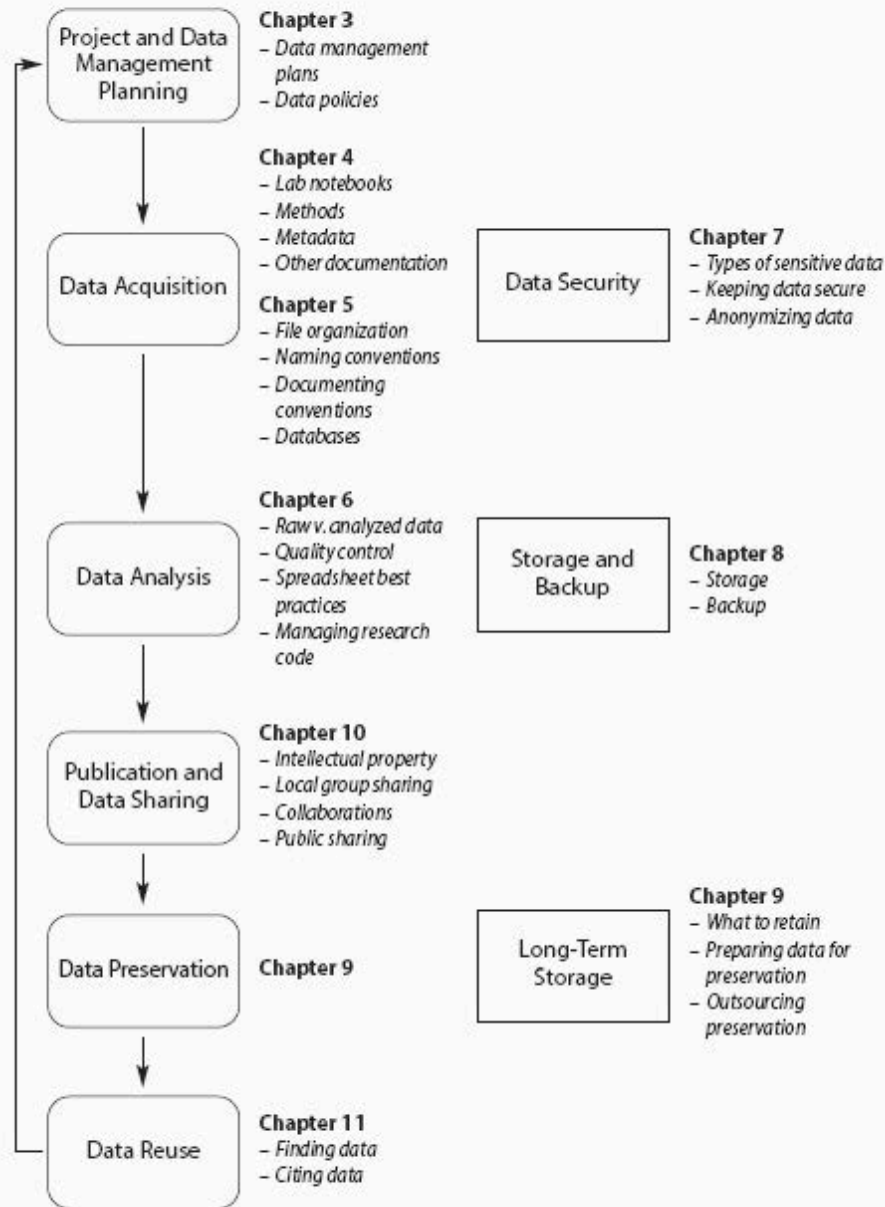


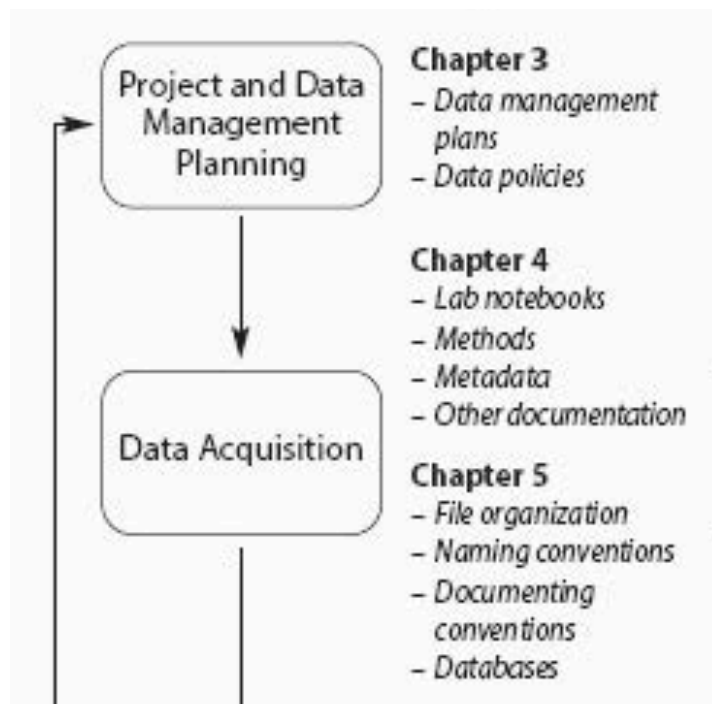
Modern Data LifeCycle



Data Lifecycle

Storage







2:15pm

● Tools for Reproducible Data Science (Hands-on) ¹
Christopher Gandrud

● Versioning your Data and Scripts (Hands-on)
Bob Freeman • Amir Karger • Radhika Khetani

4:00pm

● Finding, Accessing, Documenting: Biomedical Data
David Osterbur

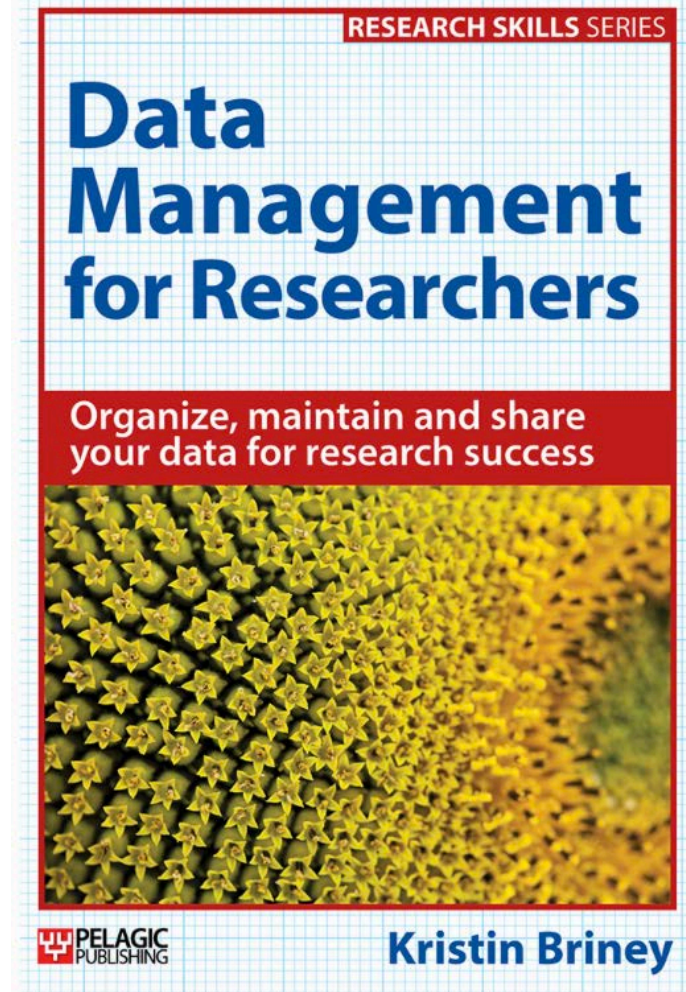
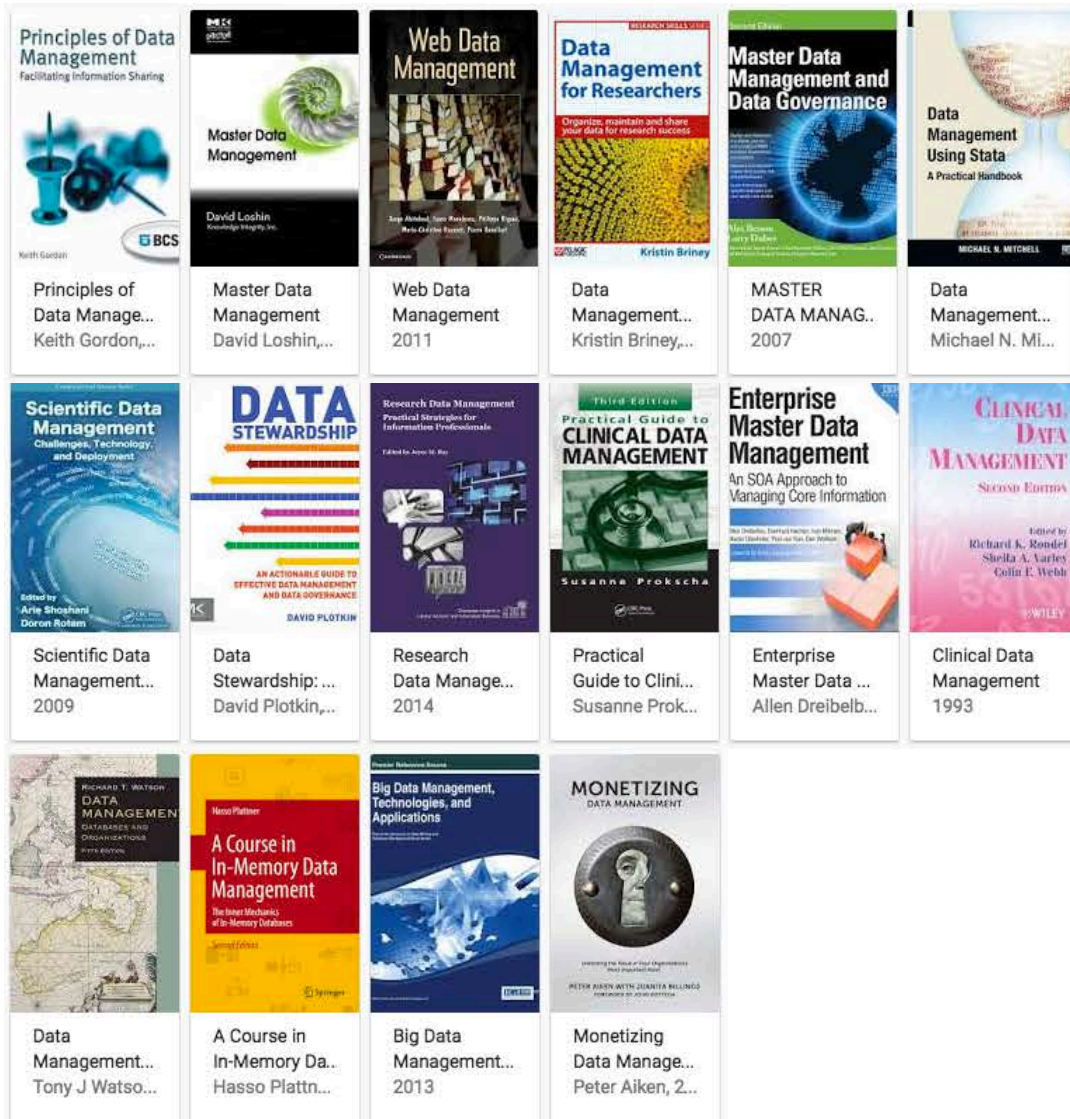
● Finding, Accessing, Documenting: Social Science Data
Alex Caracuzzo • Diane Sredl

● Understanding and Evaluating Survey Data Documentation: A User's Guide ¹
Chase Harrison

5:00pm

● Panel Discussion and Wrap-Up ¹
Hugh Truslow • Alex Caracuzzo • Bob Freeman • Christopher Gandrud • Chase Harrison • Amir Karger • Radhika Khetani • David Osterbur • Diane Sredl

Google Search for “Data Management”



Good Data Management...

Ask yourself...

- What types of data do I have? How much do I have?
- Do I use any third-party data?
- What data tools and technology are readily available to me?
- How much do I collaborate? Is this internal or external to my institution?
- How long must I keep my data?
- Will I share my data?
- Does my data have security concerns, such as personally identifiable information? What does my funder/institution/employer require?
- Is there anything particular in my research workflow that might affect my data management?
- What problems with my data do I often encounter

Good Data Management...

Your plan will cover...

- What data will you create?
- How will you document and organize your data?
- How will you store your data and, if necessary, keep it secure?
- How will you manage your data after the completion of the project?
- How will you make your data available for reuse, as necessary?

Good Data Management...

Your plan will cover...

- What data will you create?
- How will you document and organize your data?
- How will you store your data and, if necessary, keep it secure?
- How will you manage your data after the completion of the project?
- How will you make your data available for reuse, as necessary?

A New Hope[®]

Harvard Biomedical Data Management

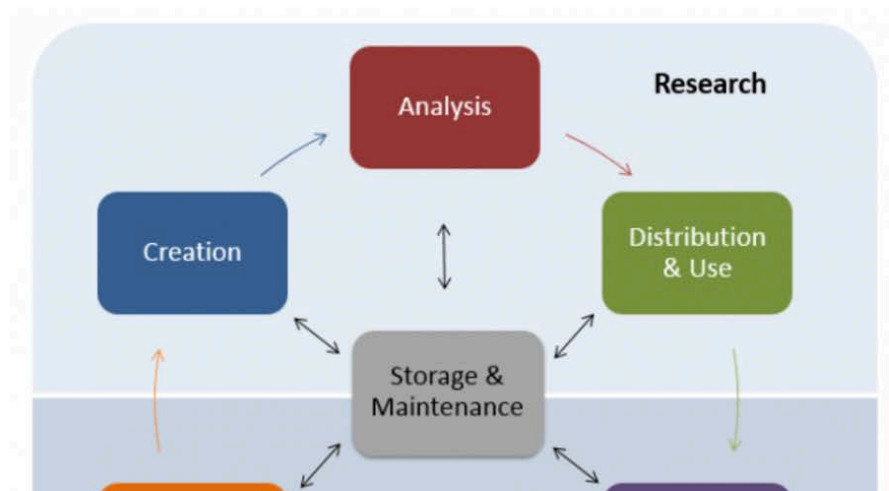
Best practices & support services for research data lifecycles

[About ▼](#)[Best Practices ▼](#)[Planning ▼](#)[Data Repositories ▼](#)[Storage ▼](#)[Policies ▼](#)[Harvard Open Access](#)

Data Management

Data Management is the process of providing the appropriate labeling, storage, and access for data at all stages of a research project. We recognize that best practices for each of these aspects of data management can and often do change over time, and are different for different stages in the data lifecycle.

Early and attentive management at each step of the data lifecycle will ensure the discoverability and longevity of your research.



Ask us your biomedical research data management questions!



The Francis A.
Countway Library of Medicine

[Receive Data Management Updates](#)

Organization

OPEN ACCESS Freely available online

PLOS COMPUTATIONAL BIOLOGY

Education

A Quick Guide to Organizing Computational Biology Projects

William Stafford Noble^{1,2*}

¹ Department of Genome Sciences, School of Medicine, University of Washington, Seattle, Washington, United States of America, ² Department of Engineering, University of Washington, Seattle, Washington, United States of America

Introduction

Most bioinformatics coursework focuses on algorithms, with perhaps some components devoted to learning programming skills and learning how to use existing bioinformatics software. Unfortunately, for students who are preparing for a research career, this type of curriculum fails to address many of the day-to-day organizational challenges associated with performing computational experiments. In practice, the principles behind organizing and documenting computational experiments are often learned on the fly, and this learning is strongly influenced by personal predilections as well as by chance interactions with collaborators or colleagues.

The purpose of this article is to describe one good strategy for carrying out computational experiments. I will not describe profound issues such as how to formulate hypotheses, design experiments, or draw conclusions. Rather, I will focus on relatively mundane issues such as organizing files and directories and documenting progress. These issues are important because poor organizational choices can lead to significantly slower research progress. I do not claim that the strategies I outline here are optimal. These are simply the principles and practices that I have developed over 12 years of bioinformatics research, augmented with various suggestions from other researchers with whom I have discussed these issues.

Principles

The core guiding principle is simple: Someone unfamiliar with your project should be able to look at your computer files and understand in detail what you did and why. This "someone" could be any of a variety of people: someone who read your published article and wants to try to reproduce your work, a collaborator who wants to understand the details of your experiments, a future student working in your lab who wants to extend your work after you have moved on to a new job, your research advisor, who may be interested in

OPEN ACCESS Freely available online

PLOS BIOLOGY

Community Page

Best Practices for Scientific Computing

Greg Wilson^{1*}, D. A. Aruliah², C. Titus Brown³, Neil P. Chue Hong⁴, Matt Davis⁵, Richard T. Guy^{6*}, Steven H. D. Haddock⁷, Kathryn D. Huff⁸, Ian M. Mitchell⁹, Mark D. Plumbley¹⁰, Ben Waugh¹¹, Ethan P. White¹²,
¹ Mozilla Foundation, Toronto, Ontario, Canada, ² Michigan, United States of America, ³ University of California, Berkeley, United States of America, ⁴ University of London, London, United Kingdom, ⁵ University of Wisconsin, Madison, Wisconsin, United States of America, ⁶ University of California, Berkeley, United States of America, ⁷ University of California, Berkeley, United States of America, ⁸ University of California, Berkeley, United States of America, ⁹ University of California, Berkeley, United States of America, ¹⁰ University of California, Berkeley, United States of America, ¹¹ University of California, Berkeley, United States of America, ¹² University of California, Berkeley, United States of America

¹ Mozilla Foundation, Toronto, Ontario, Canada, ² Michigan, United States of America, ³ University of California, Berkeley, United States of America, ⁴ University of London, London, United Kingdom, ⁵ University of Wisconsin, Madison, Wisconsin, United States of America, ⁶ University of California, Berkeley, United States of America, ⁷ University of California, Berkeley, United States of America, ⁸ University of California, Berkeley, United States of America, ⁹ University of California, Berkeley, United States of America, ¹⁰ University of California, Berkeley, United States of America, ¹¹ University of California, Berkeley, United States of America, ¹² University of California, Berkeley, United States of America

Introduction

Scientists spend an enormous amount of time using software. However, they do not do this efficiently. As a result, they waste time on practices that would be more maintainable if they used better practices for scientific computing. This article describes foundations in research that improve scientists' productivity.

Software is as important as telescopes and test tubes for computational problem solving. Scientists, more and more, are spending large amounts of time on developing new projects, combining data, and other computational tasks.

Scientists typically do not do this because doing so requires a lot of time. As a result, recent studies show that 30% or more of their time is spent on lack of exposure to basic writing, maintainable code, trackers, code reviews, and other computational tasks.

We believe that software apparatus [3] and should as any physical apparatus.

PLOS BIOLOGY

COMMUNITY PAGE

Computing Workflows for Biologists: A Roadmap

Ashley Shade^{1,2*}, Tracy K. Teal^{2,3}

Code and Data for the Social Sciences: A Practitioner's Guide

Matthew Gentzkow

Jesse M. Shapiro¹

Chicago Booth and NBER

March 10, 2014



An aerial photograph of a historic university campus, likely Harvard University. The image shows a dense cluster of red-brick buildings with white window frames and classical architectural features. A prominent central clock tower with a green dome and a bell tower is visible. The campus is surrounded by lush green trees. The text "Planning & Metadata Frameworks" is overlaid in white, sans-serif font across the center of the image.

Planning & Metadata Frameworks

Metadata Frameworks & Standards



Daina Bouquin

Head Librarian, Harvard-Smithsonian Center for Astrophysics



Derek Miller

Assistant Professor of English, Harvard University



Caroline Shamu

Assistant Professor, Harvard Medical School

An aerial photograph of a historic university campus, likely Harvard University. The image shows a dense cluster of red-brick buildings with white window frames and classical architectural features. A prominent central clock tower with a green dome and a bell tower section stands out. The campus is surrounded by lush green trees, and the overall scene is captured from a high angle, providing a comprehensive view of the architectural layout.

Data Toolbox



Data Toolbox

- Will post more information on the GitHub repository
- Invite others to send me their recommendations
- Important to know your (system) limitations...
 - Transitioning to cloud and cluster (HPC/HTC) systems doesn't have to be difficult!