Before we start...



- SABS VM issue ensure you're using new version!
 - If you haven't 'v2' in the VM name under VirtualBox, you'll need to upgrade
 - i.e. download, 'Import Appliance' in VirtualBox, delete old one
 - Help with issues after this intro
- You can find it here:

https://bit.ly/NewSABSVM

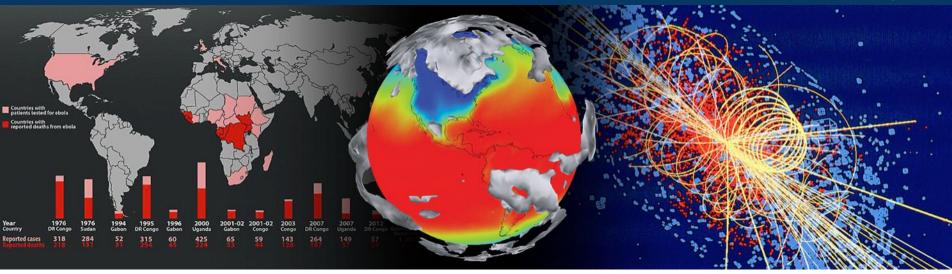
Introduction to Software Engineering Concepts

Steve Crouch
Software Sustainability Institute
s.crouch@software.ac.uk

11th October 2021

Modern research is impossible without software

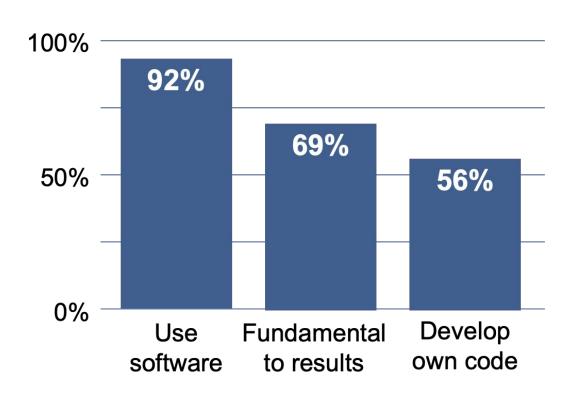




From thrown-together scripts, through an abundance of complex spreadsheets, to the millions of lines of code behind large-scale infrastructure, there are few areas where software does not play a fundamental part in research

Why should we care about software?





SSI survey of researchers, 2014[1]

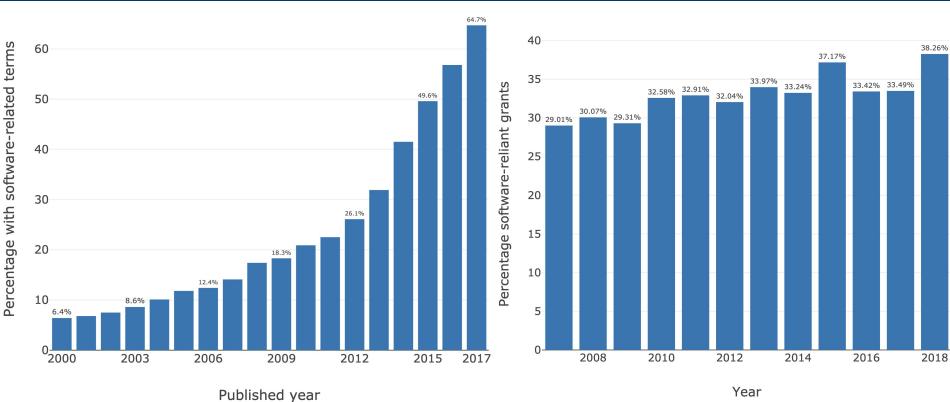
15 Russell Group Universities

Their software use and background

417 respondents

Why should we care about software?

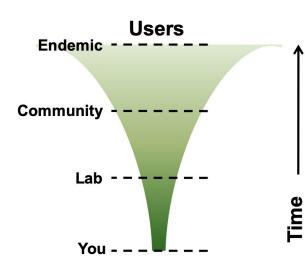




The software you write is important!



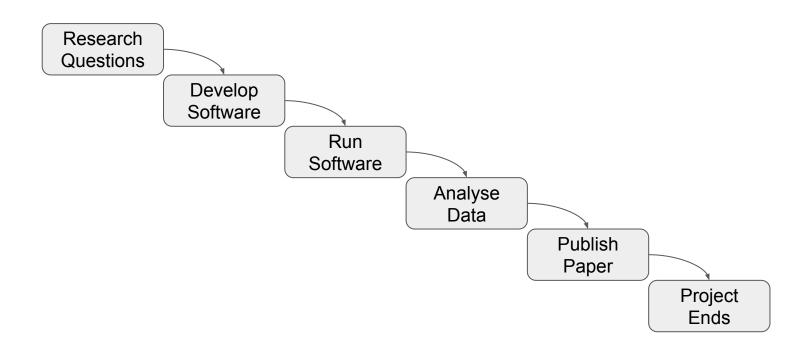
- Software inherently contains value
 - o Produces results, contains lessons learnt, effort
- Difficult to gauge to what extent it might be used in the future
 - O By who?
 - O Which parts?
 - Which projects?
 - Reproducibility from publications!



Can it/should it be reusable by others? ...including yourself?

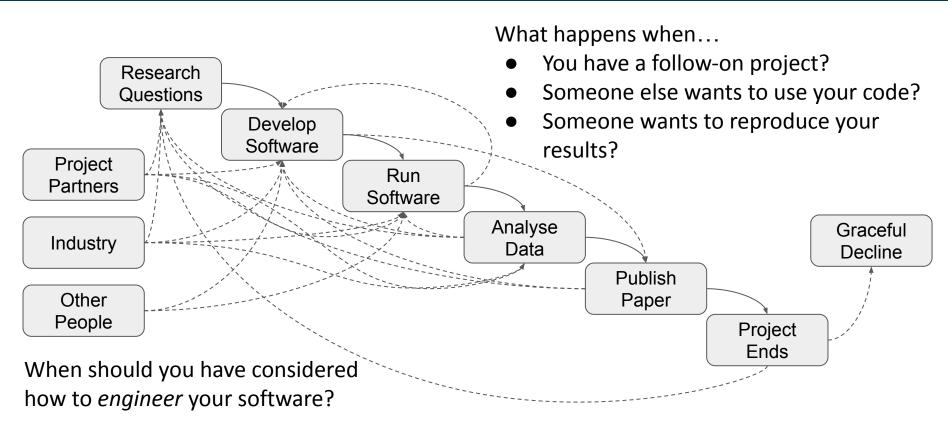
A typical research software lifecycle





In reality...





When?





"The best time to plant a tree is 20 years ago.

The second best time is now."



Programming vs Engineering



Programming / Coding

- Focus is on one aspect of software development
- Writes software for themselves
- Mostly an individual activity
- Writes software to fulfil research goals (ideally from a design)

Engineering

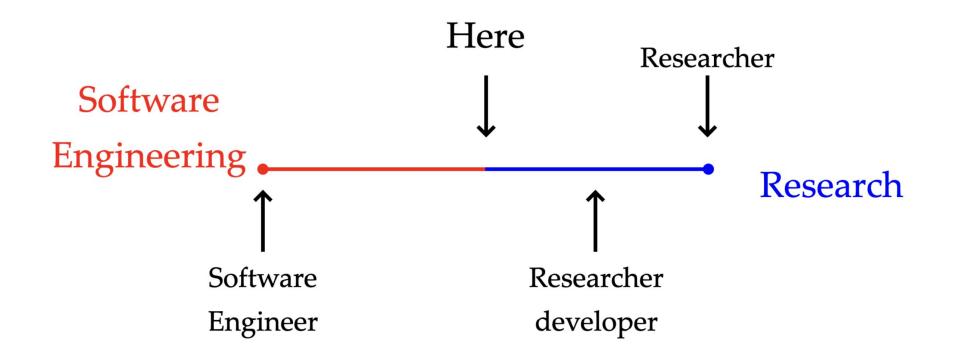
- Considers the *lifecycle* of software
- Writes software for stakeholders
- Takes team ethic into account
- Applies a process to understanding, designing, building, releasing, and maintaining software

"Programmers tend to start coding right away. Sometimes this works."

- Eric Larsen, 2018

Where are you?





Beyond building a 'sequence of instructions'

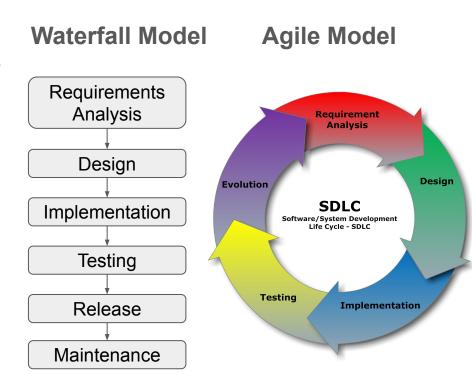


Software is far more than that...

Outcome of a development process

But also...

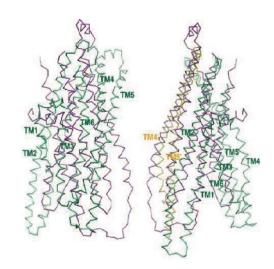
- Architecture
- Implementation of algorithms
- Data model
- Programming paradigm
- Documentation
- Best practices and conventions ...



Testing



- Humans are fallible! Our software will contain defects
 - o In requirements, design, as well as code
 - 1-10-150 hours to fix in design/development/production
- Validation: are we building the right product?
- Verification: are we building the product right?
 - Manual testing, unit testing, automated testing, code reviews
- Highly-cited papers published on multidrug resistance transporters between 2001 - 2010
- Results couldn't be reproduced 5 retractions
- Caused by error in an internal software utility
 - Flipped two columns of data, inverting electron-density map used to derive protein structure



"I didn't question it then. Obviously now I check it all the time."

- Geoffrey Chang[3]

Platform support?



... Density functional theory nuclear magnetic resonance calculations established the relative configurations of 1 and 2 and revealed that the calculated shifts depended on the operating system when using the "Willoughby-Hoye" Python scripts to streamline the processing of the output files, a previously unrecognized flaw that could lead to incorrect conclusions.

 Due to different sorting of file names on different operating systems



Organic Letters, October 8 2019 https://doi.org/10.1021/acs.orglett.9b03216

Optimisation



"Three orders of magnitude in **machine speed** and three orders of magnitude in **algorithmic speed** add up to six orders of magnitude in solving power. A model that might have taken a year to solve 10 years ago can now solve in less than 30 seconds."

- Robert Bixby, review of linear programming solvers from 1987-2002
- Faster code, faster results!
- Understanding trade-offs
 - Maintainability, accuracy
- When & where to optimise?
 - o 80/20 rule, code profiling

Amdahl's Law:

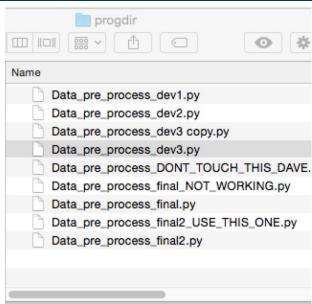
Time to result = Develop Time (D) + Time to Run (R)

As effort is put into reducing R, overall time required to get new result is dominated by writing, testing, maintaining, installing, configuring software.

Code management & collaboration



- Version control provides a full history of your project's software and other assets
- Makes for easy:
 - Backups
 - Collaboration
 - Recovering from dead-ends
- What should be in version control?
 - Code, documentation, tests, test data, analysis scripts
 - o Reports, papers, etc.
- Packaging and deployment



"If you're not using version control, whatever else you might be doing with a computer, it's not science."

- Greg Wilson, SWC

Other key points



These skills will save you time

Always assume others will use and develop your software

Be clear on requirements and assume they will change

 Funders are increasingly expecting software outputs to be sustainable and reusable

More on software engineering



Facts and Fallacies of Software Engineering



Robert L. Glass Foreword by Alan M. Davis

Robert L Glass, Addison-Wesley Professional

Group mini-project



For week 2, in same groups as those for Zoom...

Design and implement Python library to specify & solve a Pharmacokinetic model, which should ideally have the following functionality:

- pip installable
- github repository, with issues + PRs that fully document development process
- unit testing with good test coverage
- fully documented, e.g. README, API documentation, OS license
- continuous integration for automated testing/doc generation
- Ability to specify form of the PK model
- Users can specify protocol independently from the model
- Ability to solve for drug quantity in each compartment over time
- Ability to visualise the solution of a model, compare two different solutions
- Something else? Feel free to suggest alternative features!

Group projects: SABS students



- 1. Automated Data Extraction to Future-Proof Therapeutic and Natural Antibody Databases [UCB, Roche, GSK]
 - Antibody informatics databases manually curated & can't harvest literature/other sources automatically
 - Goal to increase usability and automation of their updating systems, possibly leveraging NLP
 - Develop frameworks to maximally automate updating process for each database
- 2. "Drug Discovery Game" SABS R³ software project [Roche]
 - Develop game software to help learn how medicinal chemists make compound progression decisions
 - o Inspired by "Drug Discovery Game", similar to Mastermind game; successful guesses cost play money
 - Medicinal variant: basis of possible substituents/functional groups attached to molecular scaffold

Group projects: SABS students



3. Cost effectiveness of HPV vaccination in Asia-Pacific Region [Roche]

- o Builds on web-based infrastructure of 2020-2021 SABS students to simulate epidemics
- Dynamic age-stratified epidemiological-economic modelling codebase for HPV in Laos
- o Professionally-engineered open-source select disease inputs, run model across scenarios
- Display metrics representing cost-effectiveness analysis, interactive exploration of models for policymakers

4. Computer Vision-based Clinical Imaging Quality Control [GE Healthcare]

- Clinical trial imaging data is multisite, but centrally analyse and variable quality
- o 2020-2021 SABS:R3 project open, extensible framework to ingest image metadata w.r.t. Rules
- Develop plug-in to extend QC to the pixel data using computer vision characterisation techniques
- QC tasks may include checking for "burned-in" annotation, anatomical region/completeness
- Possible extensions include auto-masking of identifiable info, extensions to QC for other checks

Group project: DTP, NERC students



Tree Generation Desktop Application

- Develop a desktop GUI application to generate and visualise branching tree structures using L-systems
- Investigate GUI design patterns such as Model View Controller, and methods for unit testing GUI applications
- Deploy to a stand-alone executable suitable for Linux, Mac OS X and Windows

OR

Suggest a suitable group software engineering project, perhaps relevant to your own research interests

General daily teaching structure

16:30 Finish



09:30-10:00	Welcome ar	nd Q&A	Main Zoom room
00.00 10.00	VVCICOTTIC at	ila Qa/ t	

12:30-13:30 Lunch (advisory only!) Advisory (instructors' may be unavailable)

13:30-14:00 Q&A session Main Zoom room

Advisory (can keep working if you like)

Each Zoom room will have 'roving' demonstrators Also - there will be a "common room" for breaks!

A few infrastructure things...



- Ensure you have your full name set in the *Participants* list
 - Participants -> hover over your entry, select More then Rename
- Please mute when not talking
- Need help?
 - In first instance, ask demonstrators for help by raising hand in Zoom room
 - Particularly if there are a lot of people asking questions demonstrators can answer them in order
 - Or if quiet, just ask them (or your Zoom buddies!)
 - o If demonstrator not in room, ask on Slack
- Move Zoom rooms (i.e. to Common Room) at break times if you like
- Be sure to tick off the exercises you complete in Canvas as you go!

A few last things!



- Usually, this is a F2F course
 - Online training is challenging!
- Please bear with us!
- Remember to tick off exercises in Canvas!
- Please fill in the after-course survey!
- Download, import new VM: https://bit.ly/NewSABSVM
 - Those who have done this already: who has had any issues?



Say hi to your Zoom demonstrators & neighbours!

References



[1] "It's impossible to conduct research without software, say 7 out of 10 UK researchers", http://www.software.ac.uk/blog/2014-12-04-its-impossible-conduct-research-without-software-say-7-out-10-uk-researchers

[2] "An investigation of the funding invested into software-reliant research", https://github.com/softwaresaved/software in grants GTR

[3] "Retractions unsettle structural bio",

https://www.the-scientist.com/daily-news/retractions-unsettle-structural-bio-46891