UiT

THE ARCTIC
UNIVERSITY
OF NORWAY

# Version control and you

Einar Holsbø
mail: einar@cs.uit.no — web: 3inar.github.io
November 16th, 2016

# Take-home messages

- You're one out of four types of analyst: try to get to the better class

- You should think about the people you collaborate with

- You're a programmer, so you may as well act like one

- Programmers use version control

Thesis of these talks: reproducible code -> reproducible analysis -> reproducible research

# Classification of analysts

1. Does most work in Excel. R, stata, etc. considered "specialist tools"

2. Uses R or other real programming environment, keeps code in files, file names for version control

3. Like #2, with real version control, a programming style guide, peer review

4. Like #3, but more software engineering-savvy. Automated testing, modular code, proper documentation, packages, etc., etc.

http://ellisp.github.io/blog/2016/09/16/version-control

# Classification of analysts

Does most work in Excel. R, stata, etc. considered "specialist tools"

Uses R or other real programming environment, keeps code in files,
le names for version control

2, with real version control, a programming style guide, peer
ew

ke #3, but more software engineering-savvy. Automated testing,
modular code, proper documentation, packages, etc., etc.

http://ellisp.github.io/blog/2016/09/16/version-control

# Classification of analysts

**+++ complexity**

**+++ professionalism**

Does most work in Excel. R, stata, etc. considered "specialist tools"

Uses R or other real programming environment, keeps code files, le names for version control

2, with real version control, a programming style guide er w

ke #3, but more software engineering-savvy. Automated modular code, proper documentation, packages, etc., etc.

http://ellisp.github.io/blog/2016/09/16/version-control

# Class 1: the excel-analyst

## Mistaken Identifiers: Gene name errors can be introduced inadvertently when using Excel in bioinformatics

Barry R Zeeberg[†], Joseph Riss[†], David W Kane, Kimberly J Bussey, Edward Uchio, W Marston Linehan, J Carl Barrett and John N Weinstein ✉

[†] Contributed equally

# Class 2: uses R

# Class 3: version control, style guide, peer review



https://git-scm.com/about

# Class 4: basically a software engineer



https://wiki.teamfortress.com/wiki/File:Engineertaunt1.PNG

# Class 1 & 2: hard to work with

# Class 3 & 4: play well with others

I originally planned to reproduce some of the results on gene expression, but decided not to.

(1) First I had to find out, by guessing and asking, who had written the R scripts. It turned out that several people had been involved. No authors, no dates, no mention of who had written what, several versions writing to the same file, etc.

(2) Some of the authors are no longer employed by ISM, others are on leave, making it a harder job to get answers.

I concluded that it would probably take less time for me to write my own scripts, despite that I have to learn R first.

–Kajsa Møllersen, distributed on nowac mailing list

LACK OF PROVENANCE
+
LACK OF DOCUMENTATION
=
UNPROFESSIONAL
& UNREPRODUCIBLE

I originally planned to reproduce some of the results on gene expression, but decided against it.

(1) First I had to find out, by guessing and asking, who had written the R scripts. It turned out that several people had been involved. No authors, no dates, no mention of who had written what, just a set of R-files with generic names, etc.

(2) Some of the authors are no longer employed by ISM, others are on leave, making it a harder job to get answers.

I concluded that it would probably take less time for me to write my own scripts than to reproduce the old ones from R-list.

–Kajsa Møllersen, distributed on nowac mailing list

```matlab
%EM algorithm

function[cls, probs, mus, sigmas] = EMmachine(X, num_classes)
    % i) init all prog_m, mu_m, sigma_m
    probs = zeros(1, num_classes);
    mus = zeros(1, num_classes);
    sigmas = zeros(1, num_classes);
    cls = 0;
    l = 1;


    [s, samples] = size(X);


    max_iterations = 6000;


    covariance = cov(X);

```

# Your closest collaborator is you six months ago, but you don't reply to emails.

https://github.com/kbroman/datasciquotes

# Version control: Git and friends

- A sane way to keep track of changes to your analysis

- A structured way to collaborate on code

- Manages provenance

# How to Git locally

- git init

- git status

- git add

- git commit

# Git demo

# Git with friends: `push` and `pull`

a file on your computer

v1

a file on your computer

**v1**

another
computer (remote)

a file on your computer

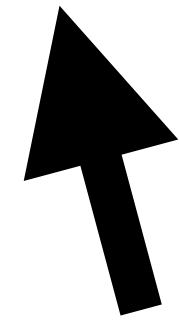**v1**

git push

another
computer (remote)

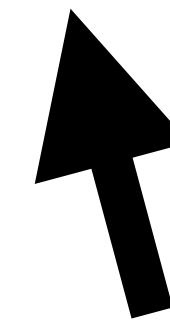a file on your computer

**v1**

**v1**

same file on another
computer (remote)

change + `commit`

**v1**

**v1**

same file on another
computer (remote)

**v2**

**v1**

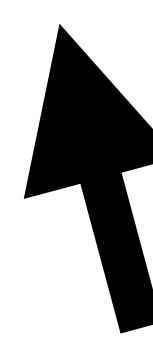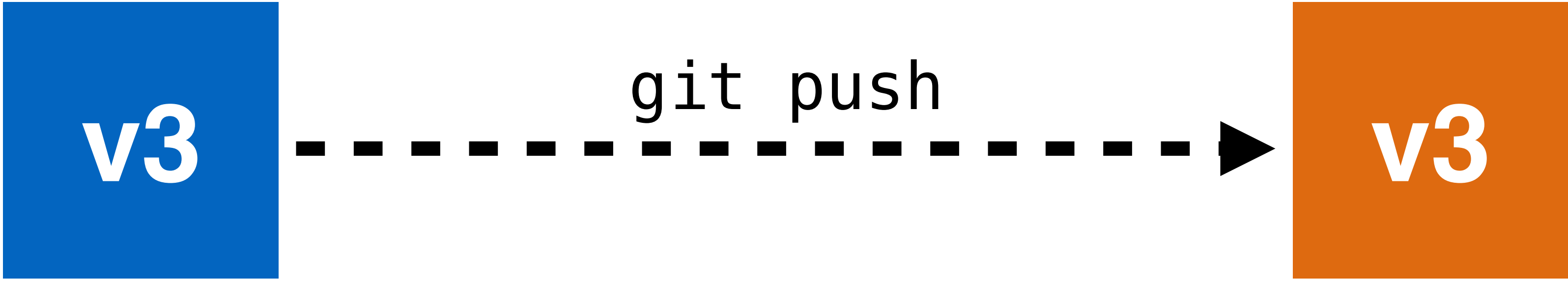same file on another
computer (remote)

change + `commit`

**v2**

**v1**

same file on another
computer (remote)

**v3**

**v1**

same file on another
computer (remote)

**v3** git push ⇢ **v1**

same file on another computer (remote)

v3

v3

Remote is up-to-date

**v3**

**v3**
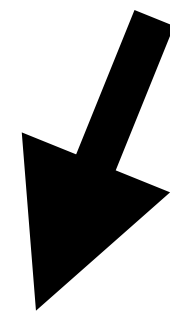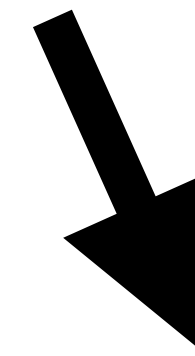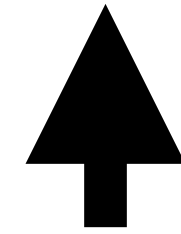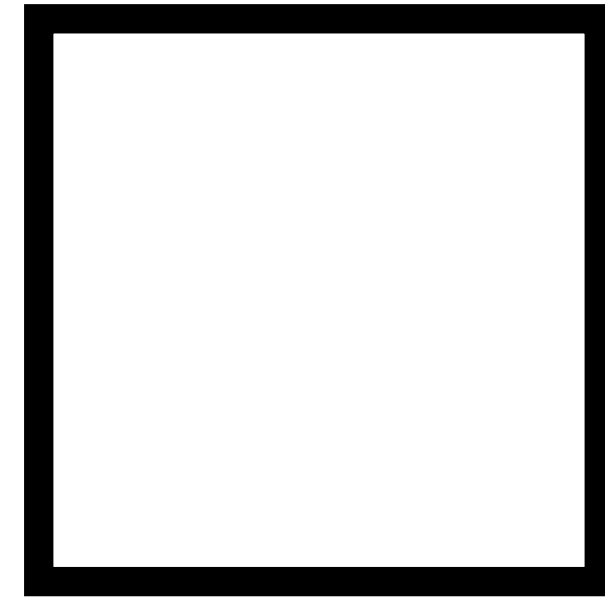
Remote is up-to-date

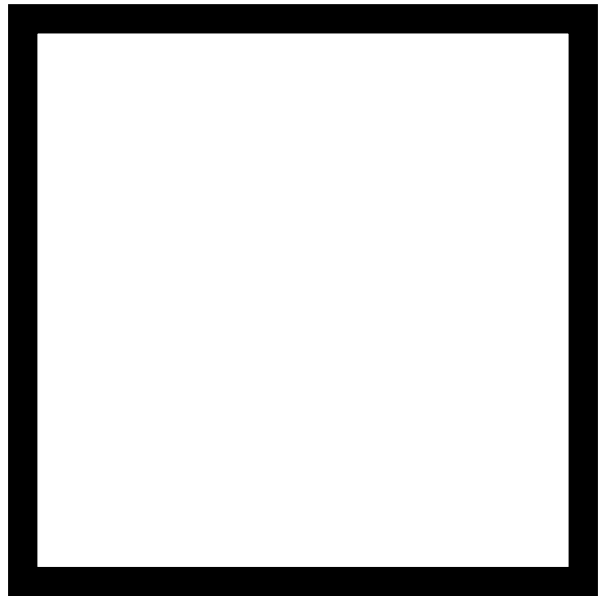Also has version history (`git log`)

a file on your computer

same file on another computer (remote)

**v3**

**v3**

my computer doesn't have this project

**v3**

**v3**

git clone

change + `commit`

v3

v3

v3

v4

v3

v3
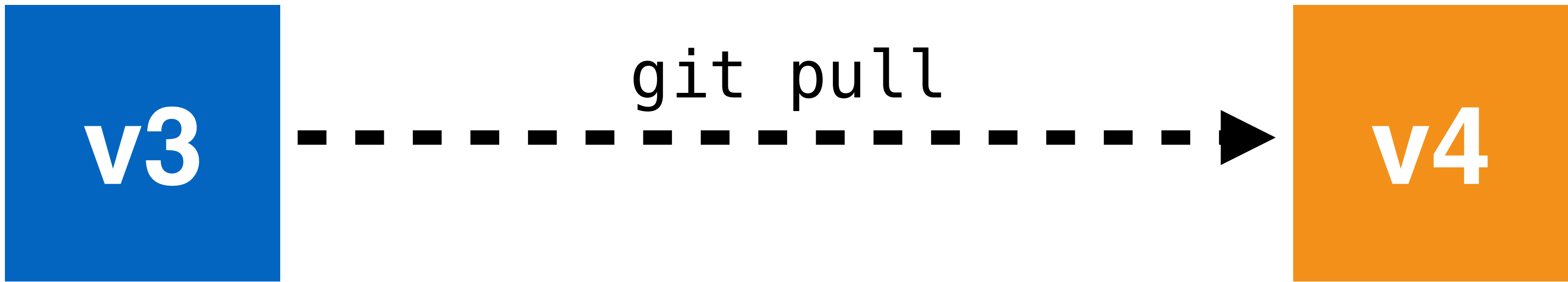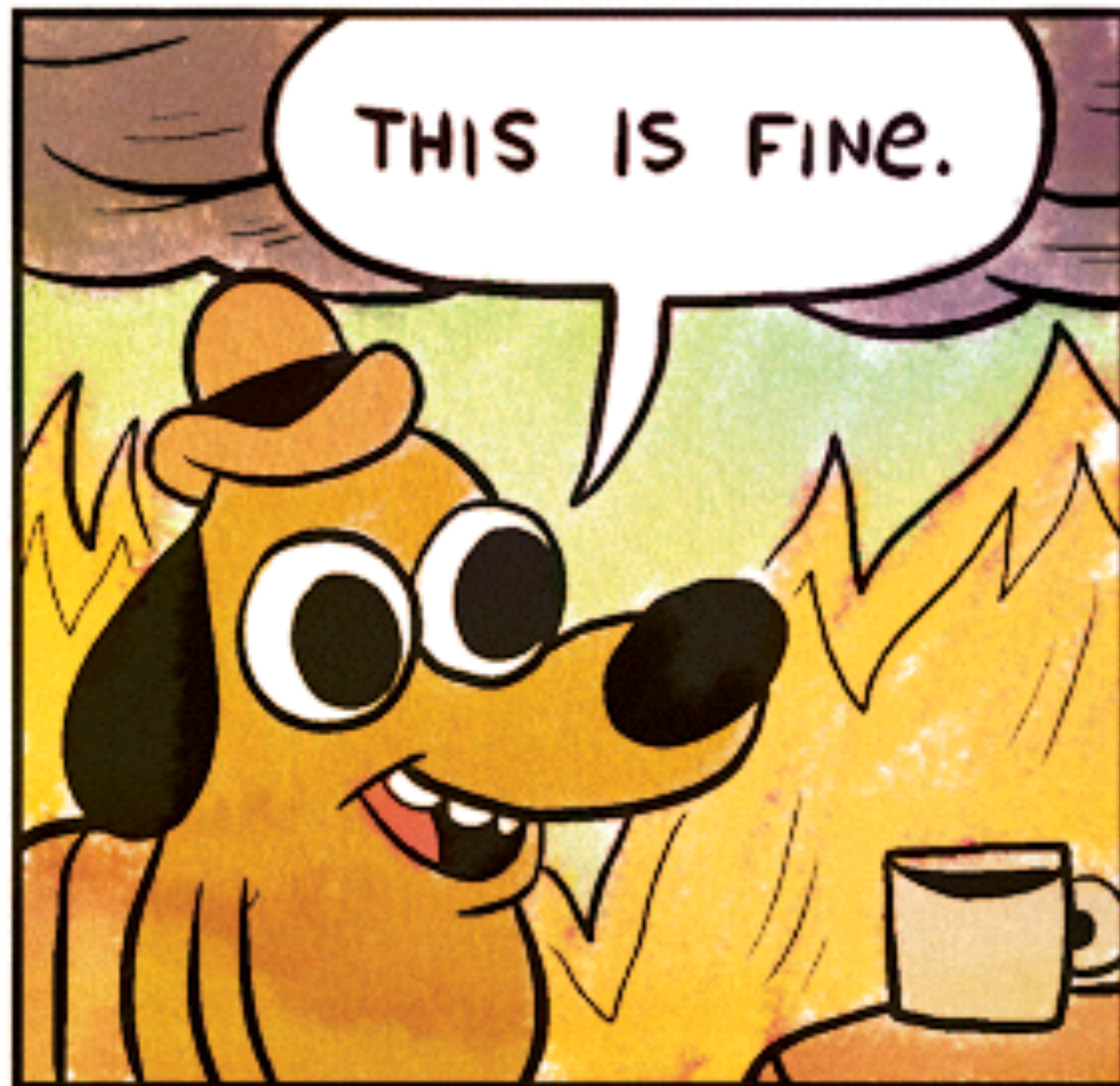
v4

v3

v4

REMOTE = SYNCHRONIZATION PT.

# Local vs remote copy

- Local copy is yours, no-one can see it unless you publish it by `push` to remote

- Remote as synchronization point, get changes by `pull`

- Etiquette note: **never push to remote if your code is broken**

THIS IS FINE.

http://gunshowcomic.com/648

# Some different remotes

- GitHub: 100% the most popular one. It's a very nice website. I host most of my code here. Public with possibility for private repos (free with academic license)

- Ice-git: located at the UiT dep't of computer science. It's private. The NOWAC documentation package lives here.

# Honorable git mentions

- `log`

- `clone`

- `checkout`

- `branch`

- `merge`
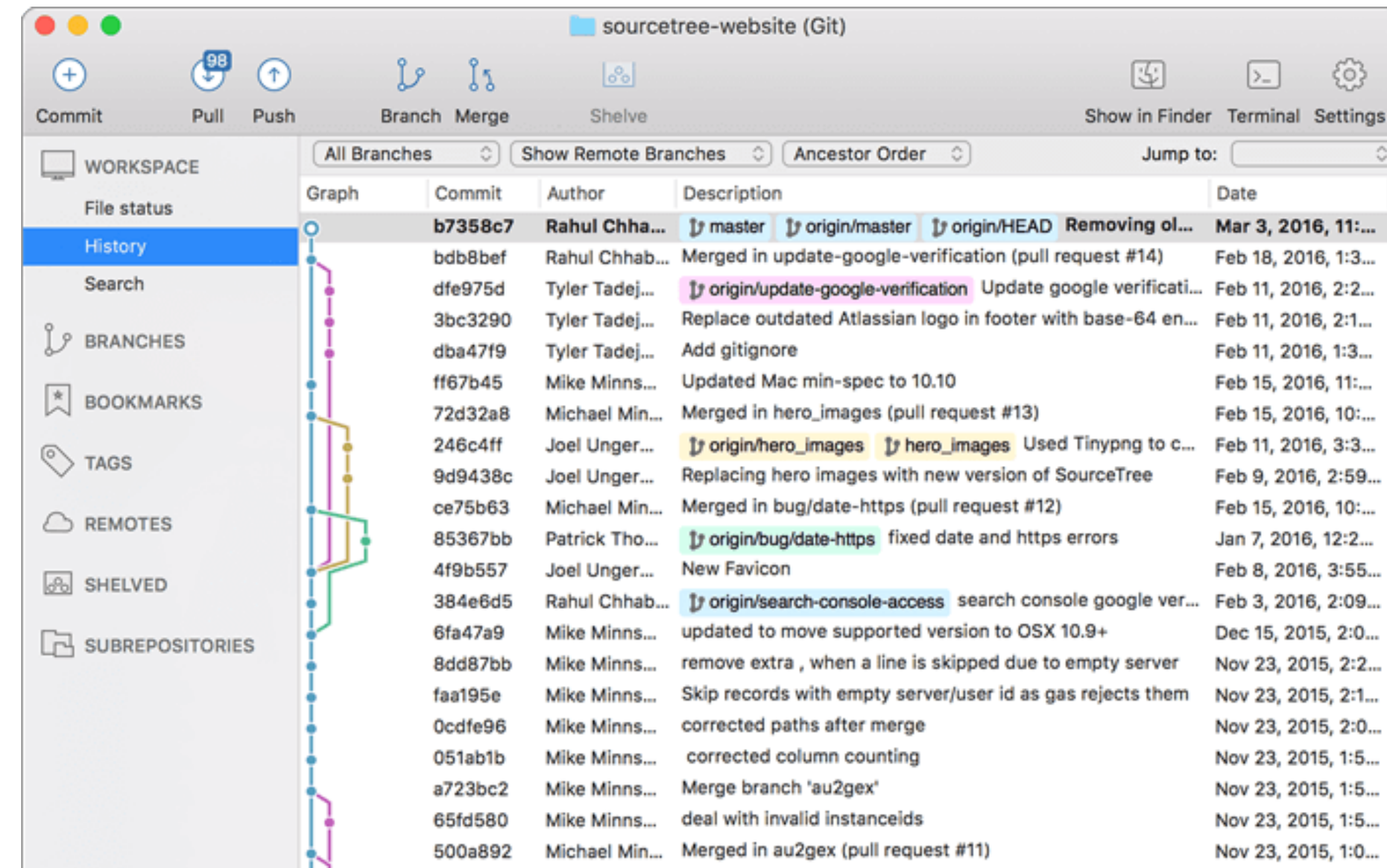
- & so on

# Branching/merging

- For keeping several versions or Branches of your code

- Eg. "development" and "release"

- Makes collaboration easier.

- You should learn about it, but learn the basics 1st.

# https://try.github.io

Learn Git in 15 minutes: Go try all this, including branching, for yourself

# SourceTree: a graphical user interface to git



https://www.sourcetreeapp.com/

# Take-home messages again

- You're one out of four types of analyst: try to get to the better class

- You should think about the people you collaborate with

- You're a programmer, so you may as well act like one

- Programmers use version control

# Next time: Packages, GitHub

# Links/References

- Github is a nice place to keep code: https://github.com/

- Github's interactive guide to git: https://try.github.io

- We use Hadley Wickham's style guide: http://adv-r.had.co.nz/Style.html

- Why you need version control: http://ellisp.github.io/blog/2016/09/16/version-control

# Links/References

- Karl Broman tutorials: http://kbroman.org/pages/tutorials

- A successful Git branching model: http://nvie.com/posts/a-successful-git-branching-model/

# Thank you!

Einar Holsbø
einar@cs.uit.no

Slides available at
3inar.github.io/talks/