

# An Overview of GitHub for Data Analysis Projects

Emerging Technologies Information Sessions 2025 Webinar Series

Modeling and Analysis Session

January 8, 2025

---

## Ben Staton

Quantitative Fisheries Scientist

Columbia River Inter-Tribal Fish Commission



bstaton1



bstaton@critfc.org

*Disclaimer: This overview is much more “what and why” of GitHub  
See upcoming PNAMP FMWG Webinar Series for “How-To”*

# Why Version Control?

*Especially for Data Analysis Projects?*

Avoid this →

"FINAL".doc



FINAL.doc!



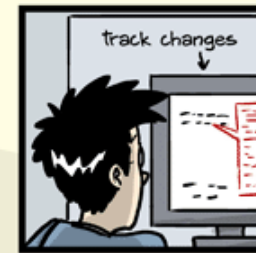
FINAL\_rev.2.doc



FINAL\_rev.6.COMMENTS.doc



FINAL\_rev.8.comments5.  
CORRECTIONS.doc



FINAL\_rev.18.comments7.  
corrections9.MORE.30.doc



FINAL\_rev.22.comments49.  
corrections.10.#@\$%WHYDID  
ICOMETOGRADSCHOOL????.doc



JORGE CHAM © 2012

<https://phdcomics.com/comics/archive.php?comid=1531>

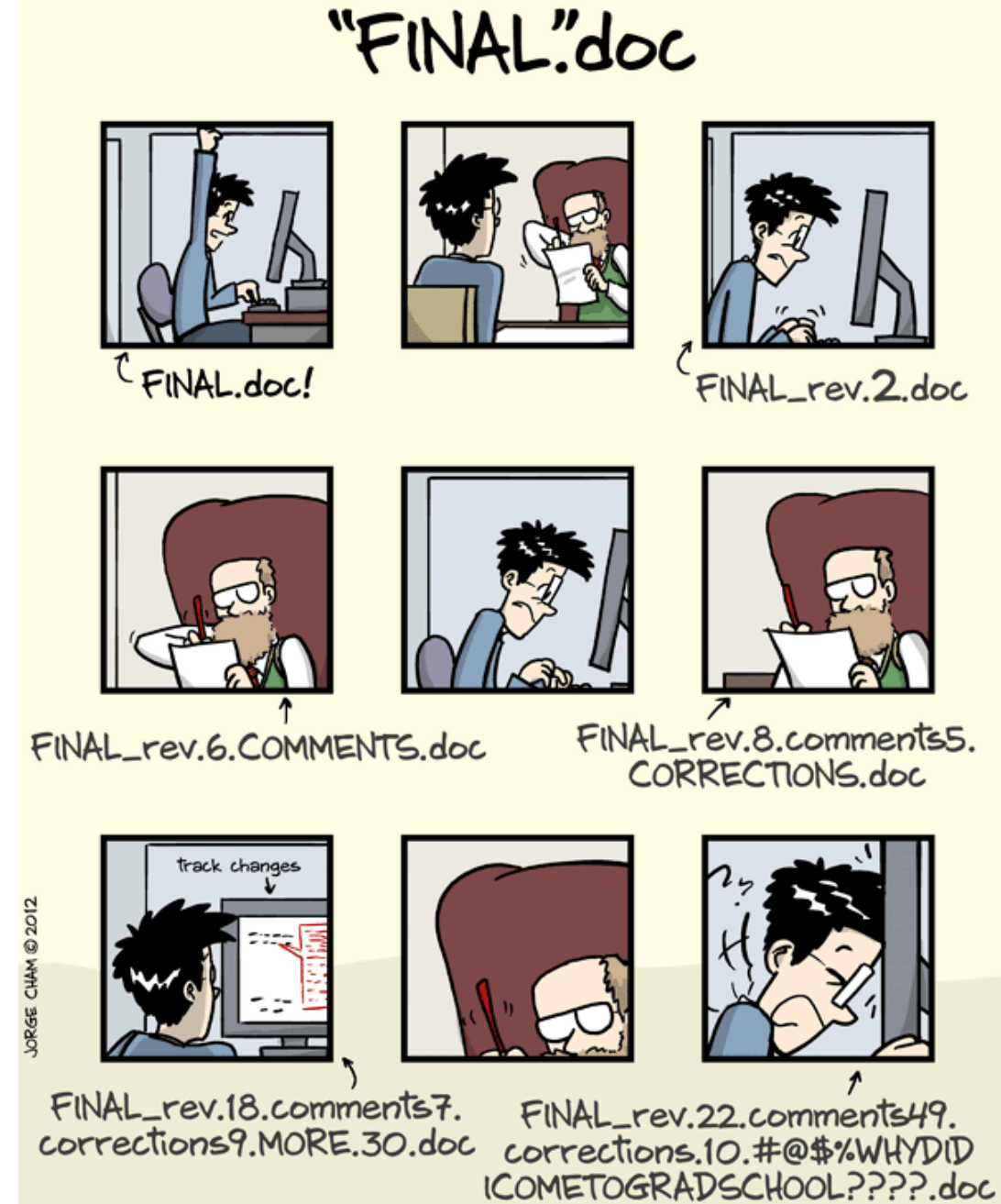
WWW.PHDCOMICS.COM

# Why Version Control?

*Especially for Data Analysis Projects?*

Avoid this →

- Data analyses have more files
  - Code & data files
  - Projects are directories



<https://phdcomics.com/comics/archive.php?comid=1531>

WWW.PHDCOMICS.COM

# Why Version Control?

*Especially for Data Analysis Projects?*

Avoid this →

- Data analyses have more files
  - Code & data files
  - Projects are directories

- version\_01
- version\_02
- version\_03
- version\_04
- version\_05
- version\_06
- version\_07
- version\_08
- version\_09
- version\_10

"FINAL".doc



FINAL.doc!



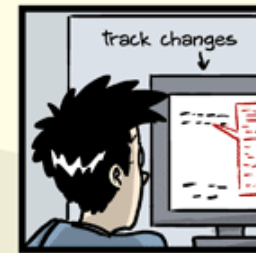
FINAL\_rev.2.doc



FINAL\_rev.6.COMMENTS.doc



FINAL\_rev.8.comments5.  
CORRECTIONS.doc



FINAL\_rev.18.comments7.  
corrections9.MORE.30.doc



FINAL\_rev.22.comments49.  
corrections.10.#@\$%WHYDID  
ICOMETOGRADSCHOOL????.doc



JORGE CHAM © 2012

<https://phdcomics.com/comics/archive.php?comid=1531>

WWW.PHDCOMICS.COM

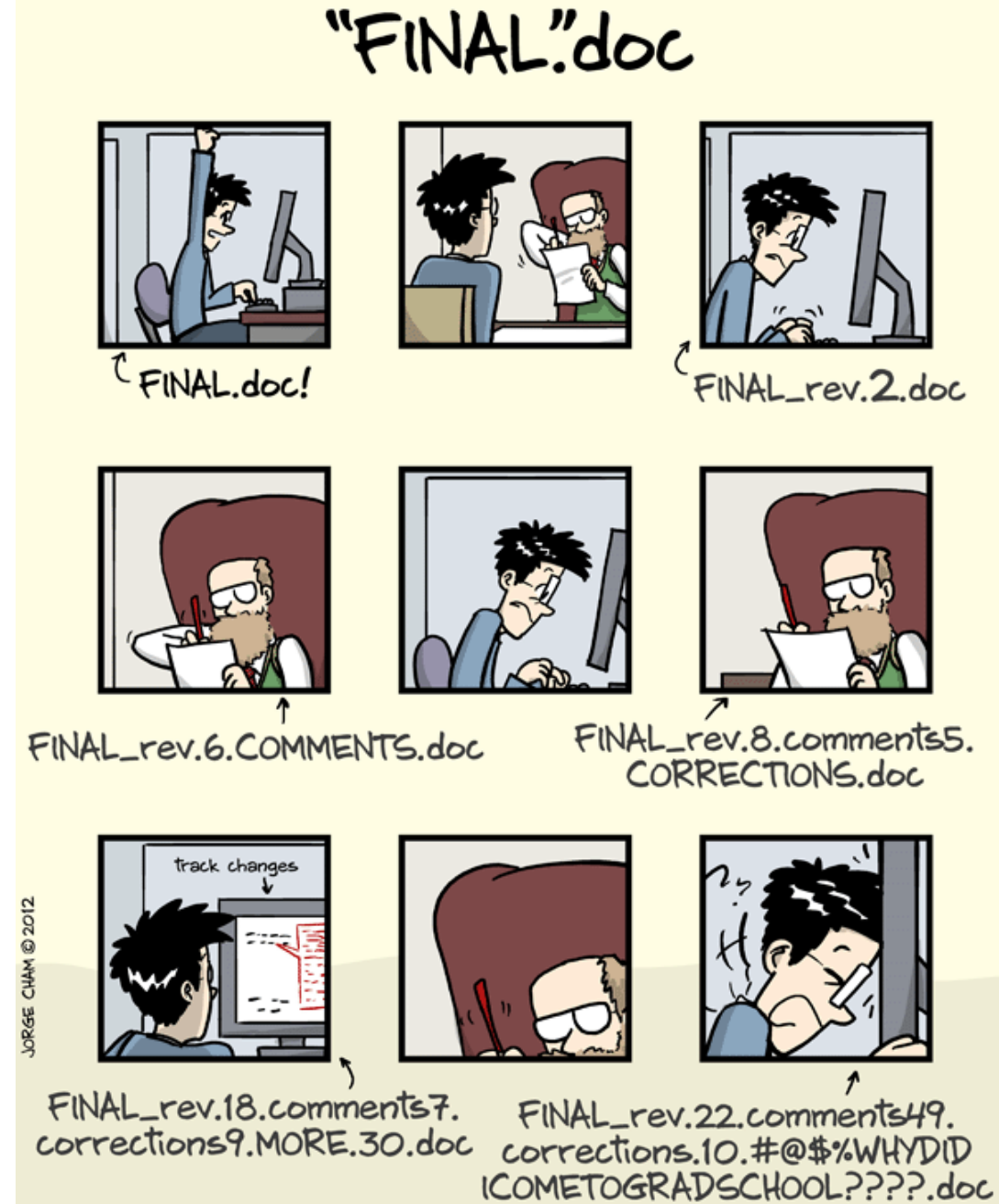
# Why Version Control?

*Especially for Data Analysis Projects?*

Avoid this →

- Data analyses have more files
  - Code & data files
  - Projects are directories →
- What if you need to (*easily*):
  - Try out something new?
  - See the history of changes?
  - Compare differences in versions?
  - Roll back to an earlier version?

- version\_01
- version\_02
- version\_03
- version\_04
- version\_05
- version\_06
- version\_07
- version\_08
- version\_09
- version\_10



<https://phdcomics.com/comics/archive.php?comid=1531>

WWW.PHDCOMICS.COM

# Why Version Control?

*Especially for Data Analysis Projects?*

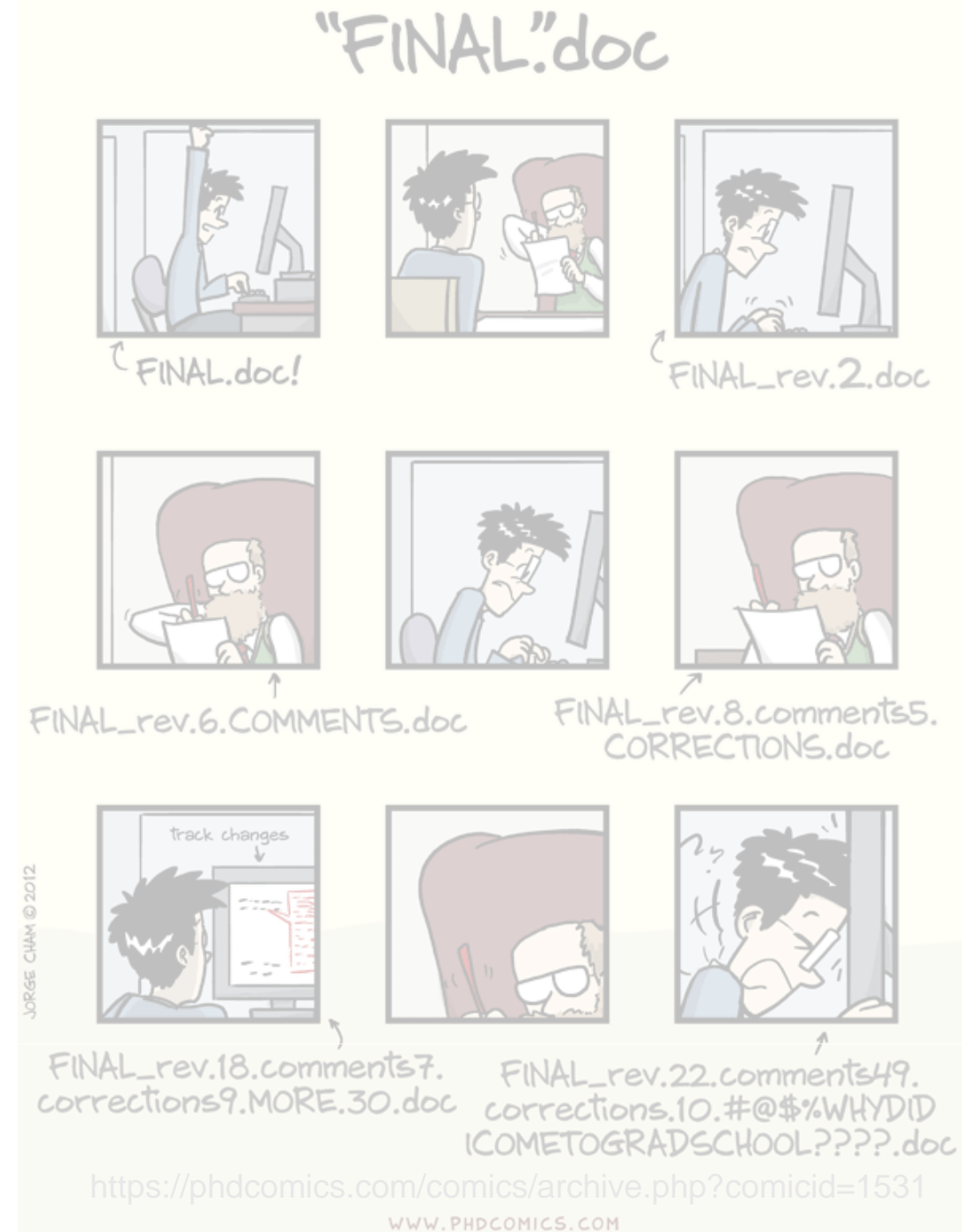
Avoid this →

- Data analyses have more files
  - Code & data files
  - Projects are directories
- What if you need to (*easily*):

- Try out something new?
- See the history of changes?
- Compare differences in versions?
- Roll back to an earlier version?

- version\_01
- version\_02
- version\_03
- version\_04
- version\_05
- version\_06
- version\_07
- version\_08
- version\_09
- version\_10

**This is what version control systems (VCS) are for (and more!)**



# Why GitHub?

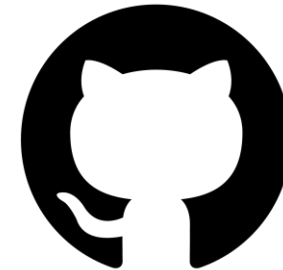
*As the Version Control Platform?*

---

- Old school git (Unix-like CLI) is hard
  - But this is where the VCS tools live
- GitHub offers:
  - All the VCS benefits of Git
  - Desktop GUI to interact with Git (*point/click*)
  - Online hosting (*centralized, most current version*)
  - Collaborative development
  - Project management
  - Project documentation
  - Free & pro versions (*no need for pro generally*)



*VCS software “under the hood”*



**GitHub**

*User-friendly, GUIs,  
collaboration, remote hosting*

*Other options: e.g., GitLab*

# Why GitHub?

*As the Version Control Platform?*

---

- Old school git (Unix-like CLI) is hard
  - But this is where the VCS tools live
- GitHub offers:
  - All the VCS benefits of Git
  - Desktop GUI to interact with Git (*point/click*)
  - Online hosting (*centralized, most current version*)
  - Collaborative development
  - Project management
  - Project documentation
  - Free & pro versions (*no need for pro generally*)

Tired of sending code/data back and forth,  
each time with different file names and  
describing what you changed?



*VCS software “under the hood”*



**GitHub**

*User-friendly, GUIs,  
collaboration, remote hosting*

*Other options: e.g., GitLab*



# History of Changes

*Who Made the Change, When, and What was the Purpose?*

Commits on May 11, 2023

Merge pull request [#195](#) from bstaton1/constant-prespawn-surv

bstaton1 authored on May 11, 2023

Remove time-varying respawn mort from sim-vs-obs

bstaton1 committed on May 11, 2023

Remove time-varying respawn mort from all output plots

bstaton1 committed on May 11, 2023

Update initial value generator

bstaton1 committed on May 11, 2023

Remove time-varying pre-spawn mortality in fitting script

bstaton1 committed on May 11, 2023

Remove time-varying pre-spawn mortality in JAGS model

bstaton1 committed on May 11, 2023

Merge pull request [#194](#) from bstaton1/time-varying-fecundity

bstaton1 authored on May 11, 2023



# History of Changes

*Who Made the Change, When, and What was the Purpose?*

Commits on May 11, 2023

Merge pull request #195 from bstaton1/constant-prespawn-surv

bstaton1 authored on May 11, 2023

Remove time-varying respawn mort from sim-vs-obs

bstaton1 committed on May 11, 2023

Remove time-varying respawn mort from all output plots

bstaton1 committed on May 11, 2023

**Update initial value generator**

bstaton1 committed on May 11, 2023

Remove time-varying pre-spawn mortality in fitting script

bstaton1 committed on May 11, 2023

Remove time-varying pre-spawn mortality in JAGS model

bstaton1 committed on May 11, 2023

Merge pull request #194 from bstaton1/time-varying-fecundity

bstaton1 authored on May 11, 2023

*Click to see more details*



# File Diffs

## Easily View File Changes

### Original Version

```
77 -   prespawn_surv = x_carcass_spawned/x_carcass_total
78 -   prespawn_surv[is.na(prespawn_surv)] = 0.75
79 -   prespawn_surv[prespawn_surv < 0.5] = 0.5
80
81     # get the age/origin compositon of the total adult return
82     # assume 60% pHOS and 20%, 60%, 20% of age 3, 4, and 5
@@ -96,7 +94,7 @@ get_E_obs = function(jags_data, fill_missing = FALSE, append_sim_yrs = FALSE) {
96
97     # remove prespawn morts
98     age_origin_survived = lapply(1:nj, function(j) {
99 -     apply(age_origin_above_weir[:,j], 2, function(a) a *
    prespawn_surv[:,j])
100     })
101     age_origin_survived = do.call(abind, append(age_origin_survived,
list(along = 3)))
```

### Changed Version

```
77 +   prespawn_surv = colSums(x_carcass_spawned, na.rm =
    TRUE)/colSums(x_carcass_total, na.rm = TRUE)
78
79     # get the age/origin compositon of the total adult return
80     # assume 60% pHOS and 20%, 60%, 20% of age 3, 4, and 5
@@ -96,7 +94,7 @@ get_E_obs = function(jags_data, fill_missing = FALSE, append_sim_yrs = FALSE) {
94
95     # remove prespawn morts
96     age_origin_survived = lapply(1:nj, function(j) {
97 +     apply(age_origin_above_weir[:,j], 2, function(a) a *
    prespawn_surv[j])
98     })
99     age_origin_survived = do.call(abind, append(age_origin_survived,
list(along = 3)))
```

# What do you need to get started?

*Beyond Patience...*

---

- GitHub account (*free*; [info](#))
- GitHub Desktop installed ([info](#))
- A beginner's knowledge of markdown syntax ([info](#))
- Working knowledge of terminology...

Term	Description
<b>1. Repository</b>	A directory (folder) with files under version control ( <i>aka "repo"</i> )
<b>2.</b>	
<b>3.</b>	
<b>4.</b>	
<b>5.</b>	
<b>6.</b>	
<b>7.</b>	
<b>8.</b>	
<b>9.</b>	
<b>10.</b>	
<b>11.</b>	
<b>12.</b>	
<b>13.</b>	

<b>Term</b>	<b>Description</b>
<b>1. Repository</b>	A directory (folder) with files under version control ( <i>aka "repo"</i> )
<b>2. Local (Repository)</b>	A repository found on your hard drive
<b>3.</b>	
<b>4.</b>	
<b>5.</b>	
<b>6.</b>	
<b>7.</b>	
<b>8.</b>	
<b>9.</b>	
<b>10.</b>	
<b>11.</b>	
<b>12.</b>	
<b>13.</b>	

<b>Term</b>	<b>Description</b>
<b>1. Repository</b>	A directory (folder) with files under version control ( <i>aka "repo"</i> )
<b>2. Local (Repository)</b>	A repository found on your hard drive
<b>3. Remote (Repository)</b>	A repository found online
<b>4.</b>	
<b>5.</b>	
<b>6.</b>	
<b>7.</b>	
<b>8.</b>	
<b>9.</b>	
<b>10.</b>	
<b>11.</b>	
<b>12.</b>	
<b>13.</b>	

<b>Term</b>	<b>Description</b>
<b>1. Repository</b>	A directory (folder) with files under version control ( <i>aka "repo"</i> )
<b>2. Local (Repository)</b>	A repository found on your hard drive
<b>3. Remote (Repository)</b>	A repository found online
<b>4. Clone</b>	Create a local repo from remote ( <i>one time only</i> )
<b>5.</b>	
<b>6.</b>	
<b>7.</b>	
<b>8.</b>	
<b>9.</b>	
<b>10.</b>	
<b>11.</b>	
<b>12.</b>	
<b>13.</b>	



<b>Term</b>	<b>Description</b>
<b>1. Repository</b>	A directory (folder) with files under version control ( <i>aka "repo"</i> )
<b>2. Local (Repository)</b>	A repository found on your hard drive
<b>3. Remote (Repository)</b>	A repository found online
<b>4. Clone</b>	Create a local repo from remote ( <i>one time only</i> )
<b>5. Commit</b>	A <u>small</u> set of changes that accomplish an "atomic task"
<b>6.</b>	
<b>7.</b>	
<b>8.</b>	
<b>9.</b>	
<b>10.</b>	
<b>11.</b>	
<b>12.</b>	
<b>13.</b>	

# Commit Messages


*They Matter*


	COMMENT	DATE
○	CREATED MAIN LOOP & TIMING CONTROL	14 HOURS AGO
○	ENABLED CONFIG FILE PARSING	9 HOURS AGO
○	MISC BUGFIXES	5 HOURS AGO
○	CODE ADDITIONS/EDITS	4 HOURS AGO
○	MORE CODE	4 HOURS AGO
○	HERE HAVE CODE	4 HOURS AGO
○	AAAAAAAAA	3 HOURS AGO
○	ADKFJSLKDFJSDKLFJ	3 HOURS AGO
○	MY HANDS ARE TYPING WORDS	2 HOURS AGO
○	HAAAAAAAAAANDS	2 HOURS AGO

AS A PROJECT DRAGS ON, MY GIT COMMIT MESSAGES GET LESS AND LESS INFORMATIVE.

<b>Term</b>	<b>Description</b>
<b>1. Repository</b>	A directory (folder) with files under version control ( <i>aka "repo"</i> )
<b>2. Local (Repository)</b>	A repository found on your hard drive
<b>3. Remote (Repository)</b>	A repository found online
<b>4. Clone</b>	Create a local repo from remote ( <i>one time only</i> )
<b>5. Commit</b>	A <u>small</u> set of changes that accomplish an "atomic task"
<b>6. Fetch</b>	Check if local is up-to-date with remote
<b>7.</b>	
<b>8.</b>	
<b>9.</b>	
<b>10.</b>	
<b>11.</b>	
<b>12.</b>	
<b>13.</b>	

<b>Term</b>	<b>Description</b>
<b>1. Repository</b>	A directory (folder) with files under version control ( <i>aka "repo"</i> )
<b>2. Local (Repository)</b>	A repository found on your hard drive
<b>3. Remote (Repository)</b>	A repository found online
<b>4. Clone</b>	Create a local repo from remote ( <i>one time only</i> )
<b>5. Commit</b>	A <u>small</u> set of changes that accomplish an "atomic task"
<b>6. Fetch</b>	Check if local is up-to-date with remote
<b>7. Pull ↓</b>	Sync commits to local from remote ( <i>if local is behind remote</i> )
<b>8. Push ↑</b>	Sync commits from local to remote ( <i>if remote is behind local</i> )
<b>9.</b>	
<b>10.</b>	
<b>11.</b>	
<b>12.</b>	
<b>13.</b>	

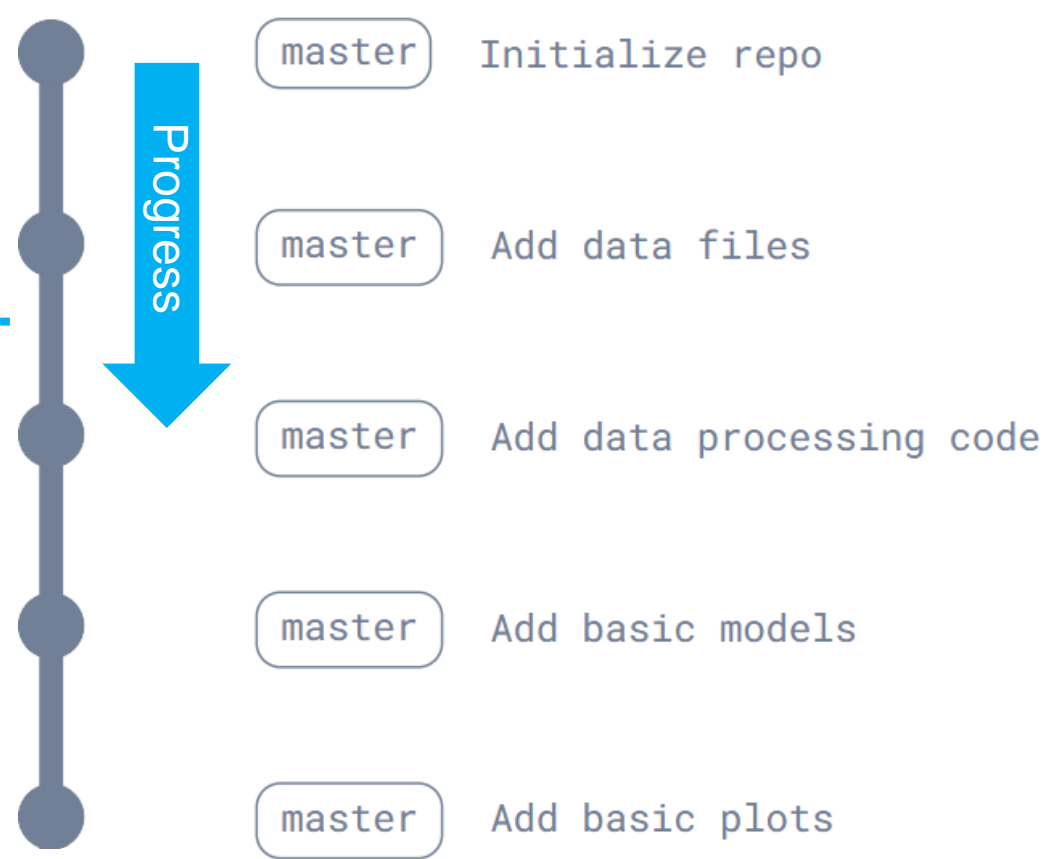
Term	Description
1. Repository	A directory (folder) with files under version control ( <i>aka "repo"</i> )
2. Local (Repository)	A repository found on your hard drive
3. Remote (Repository)	A repository found online
4. Clone	Create a local repo from remote ( <i>one time only</i> )
5. Commit	A <u>small</u> set of changes that accomplish an "atomic task"
6. Fetch	Check if local is up-to-date with remote
7. Pull ↓	Sync commits to local from remote ( <i>if local is behind remote</i> )
8. Push ↑	Sync commits from local to remote ( <i>if remote is behind local</i> )
9. Branch 	An isolated version of the repo ( <i>doesn't affect other branches</i> )
10.	
11.	
12.	
13.	


<b>Term</b>	<b>Description</b>
<b>1. Repository</b>	A directory (folder) with files under version control ( <i>aka "repo"</i> )
<b>2. Local (Repository)</b>	A repository found on your hard drive
<b>3. Remote (Repository)</b>	A repository found online
<b>4. Clone</b>	Create a local repo from remote ( <i>one time only</i> )
<b>5. Commit</b>	A <u>small</u> set of changes that accomplish an "atomic task"
<b>6. Fetch</b>	Check if local is up-to-date with remote
<b>7. Pull</b> ↓	Sync commits to local from remote ( <i>if local is behind remote</i> )
<b>8. Push</b> ↑	Sync commits from local to remote ( <i>if remote is behind local</i> )
<b>9. Branch</b> 	An isolated version of the repo ( <i>doesn't affect other branches</i> )
<b>10. Main Branch</b>	A special branch storing the "best" version
<b>11.</b>	
<b>12.</b>	
<b>13.</b>	

# A “Git Graph”

*Visualizes Development History*

- Commits constitute progress
- “Master” = “Main” branch



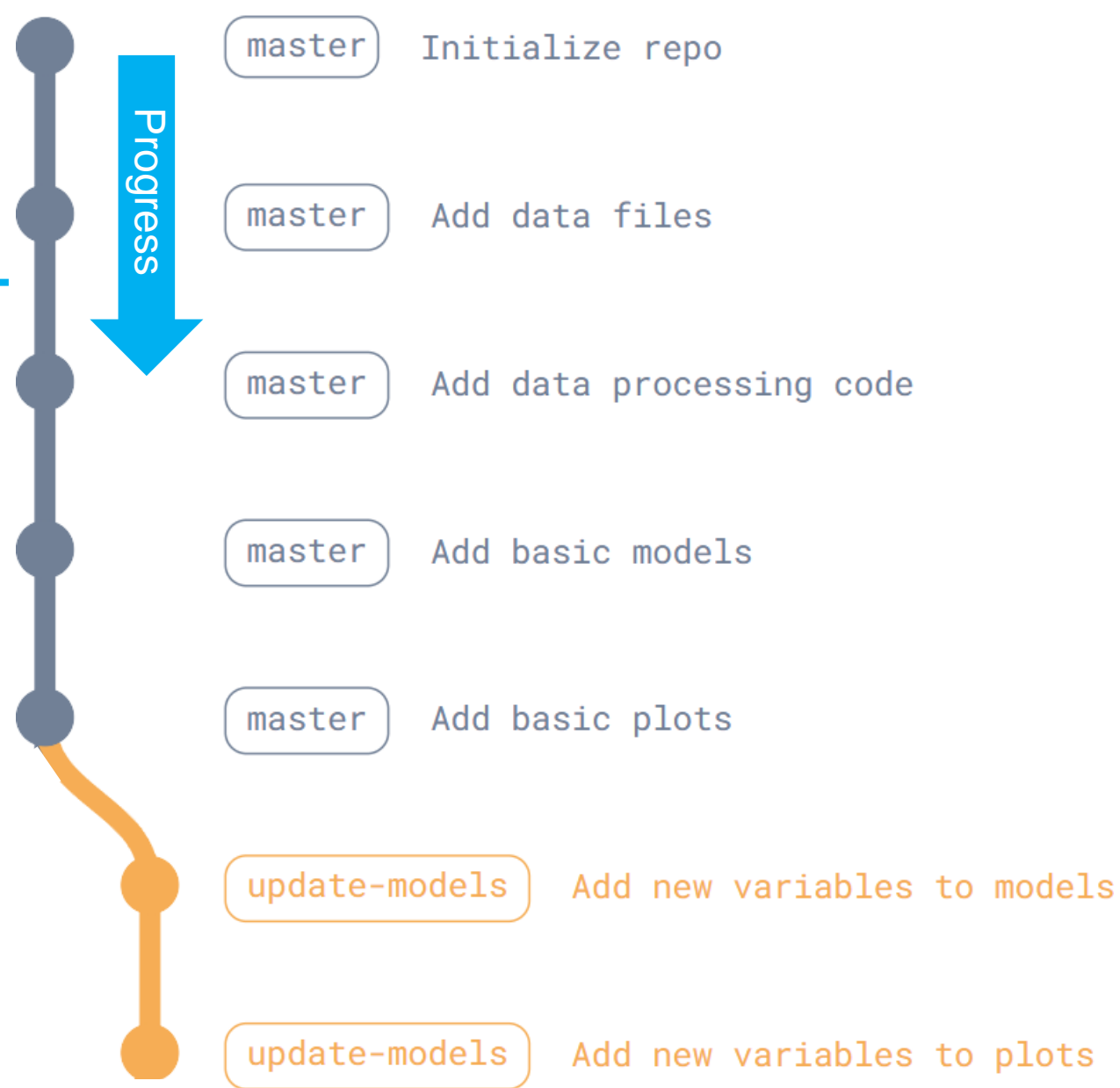
<b>Term</b>	<b>Description</b>
<b>1. Repository</b>	A directory (folder) with files under version control ( <i>aka "repo"</i> )
<b>2. Local (Repository)</b>	A repository found on your hard drive
<b>3. Remote (Repository)</b>	A repository found online
<b>4. Clone</b>	Create a local repo from remote ( <i>one time only</i> )
<b>5. Commit</b>	A <u>small</u> set of changes that accomplish an "atomic task"
<b>6. Fetch</b>	Check if local is up-to-date with remote
<b>7. Pull</b> ↓	Sync commits to local from remote ( <i>if local is behind remote</i> )
<b>8. Push</b> ↑	Sync commits from local to remote ( <i>if remote is behind local</i> )
<b>9. Branch</b> 	An isolated version of the repo ( <i>doesn't affect other branches</i> )
<b>10. Main Branch</b>	A special branch storing the "best" version
<b>11. Feature Branch</b>	A branch storing a version under development
<b>12.</b>	
<b>13.</b>	





# A “Git Graph”

*Visualizes Development History*

- Commits constitute progress
- “Master” = “Main” branch
- Create a branch and add commits
- Feature branch has all upstream commits on main branch
- But main remains unchanged

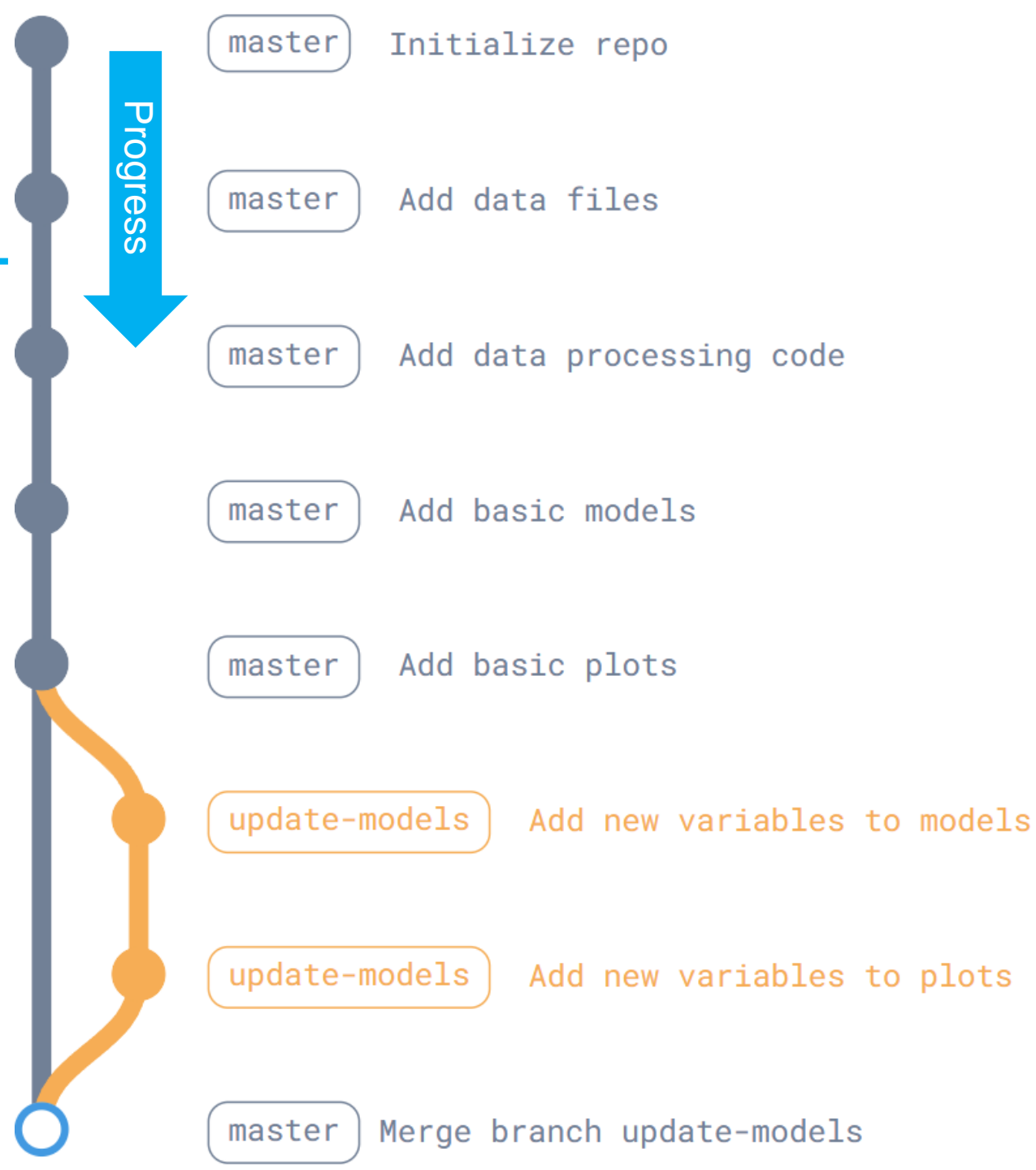




Term	Description
1. Repository	A directory (folder) with files under version control ( <i>aka "repo"</i> )
2. Local (Repository)	A repository found on your hard drive
3. Remote (Repository)	A repository found online
4. Clone	Create a local repo from remote ( <i>one time only</i> )
5. Commit	A <u>small</u> set of changes that accomplish an "atomic task"
6. Fetch	Check if local is up-to-date with remote
7. Pull ↓	Sync commits to local from remote ( <i>if local is behind remote</i> )
8. Push ↑	Sync commits from local to remote ( <i>if remote is behind local</i> )
9. Branch 	An isolated version of the repo ( <i>doesn't affect other branches</i> )
10. Main Branch	A special branch storing the "best" version
11. Feature Branch	A branch storing a version under development
12. Merge 	Add commits from one branch to another
13.	

# A “Git Graph”

*Visualizes Development History*

- Commits constitute progress
- “Master” = “Main” branch
- Create a branch and add commits
- Feature branch has all upstream commits on main branch
- But main remains unchanged
- If feature commits are beneficial, add them to main with merge
- After merge, delete feature branch



<b>Term</b>	<b>Description</b>
<b>1. Repository</b>	A directory (folder) with files under version control ( <i>aka "repo"</i> )
<b>2. Local (Repository)</b>	A repository found on your hard drive
<b>3. Remote (Repository)</b>	A repository found online
<b>4. Clone</b>	Create a local repo from remote ( <i>one time only</i> )
<b>5. Commit</b>	A <u>small</u> set of changes that accomplish an "atomic task"
<b>6. Fetch</b>	Check if local is up-to-date with remote
<b>7. Pull</b> ↓	Sync commits to local from remote ( <i>if local is behind remote</i> )
<b>8. Push</b> ↑	Sync commits from local to remote ( <i>if remote is behind local</i> )
<b>9. Branch</b> 	An isolated version of the repo ( <i>doesn't affect other branches</i> )
<b>10. Main Branch</b>	A special branch storing the "best" version
<b>11. Feature Branch</b>	A branch storing a version under development
<b>12. Merge</b> 	Add commits from one branch to another
<b>13. Pull (merge) Request</b>	Perform a merge in a documented/collaborative fashion

# Miscellaneous Topics

*Time Allowing*

*These Topics are More Advanced*

*Will Require Exploring Some on Your Own*

# Issue Tracker

*Organize Ideas, Next Steps, Discuss Problems*

---

- Online feature only
- “# system” allows cross-referencing
- Can use as something of a notebook
- Comment threads
- Can group issues, tag people
- Uses markdown syntax
- Open vs. closed issues
- [Info](#)

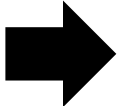
☐ ○ 2 Open ✓ 111 Closed

- ☐ ✓ Add output-plots section for adult survival from estuary to successful spawning  
#209 by bstaton1 was closed on Dec 6, 2023
- ☐ ✓ Implement a higher observation CV on mean length data  
#205 by bstaton1 was closed on Jul 11, 2023
- ☐ ✓ New section in output plots: assumed-known quantities **priority: low**  
#204 by bstaton1 was closed on Jul 11, 2023 6 tasks done
- ☐ ✓ Evaluate use of non-stochastic pre-spawn survival terms **priority: high**  
#203 by bstaton1 was closed on Jul 11, 2023
- ☐ ✓ Revise bio data processing for age data  
#198 by gibsonpp was closed on May 18, 2023
- ☐ ✓ Revise bio data processing for weir data **priority: high**  
#193 by gibsonpp was closed on May 16, 2023
- ☐ ✓ Process Noise Confounded: Pre-spawn & Egg-to-Parr Survival **priority: high**  
#192 by bstaton1 was closed on May 11, 2023

# .gitignore

*Exclude some files from tracking*

---

- E.g.,
  - Large output files
  - Things that get re-written often
- Example .gitignore file 
- Place in main repo directory
  - Subdirectories can have separate ones
- Edits to ignored files cannot be committed, files cannot be pushed

```
# Files associated with R sessions/projects
.Rproj.user/
.Rhistory
.RData
.Ruserdata

# Saved R objects (e.g., large model output)
*.rds

# Temporary files created by Rmarkdown
*.utf8.md
*.knit.md

# HTML files (generated from Rmarkdown)
*.html

# Figure files (can be recreated from source)
*.pdf
*.png
*.jpg
*.jpeg

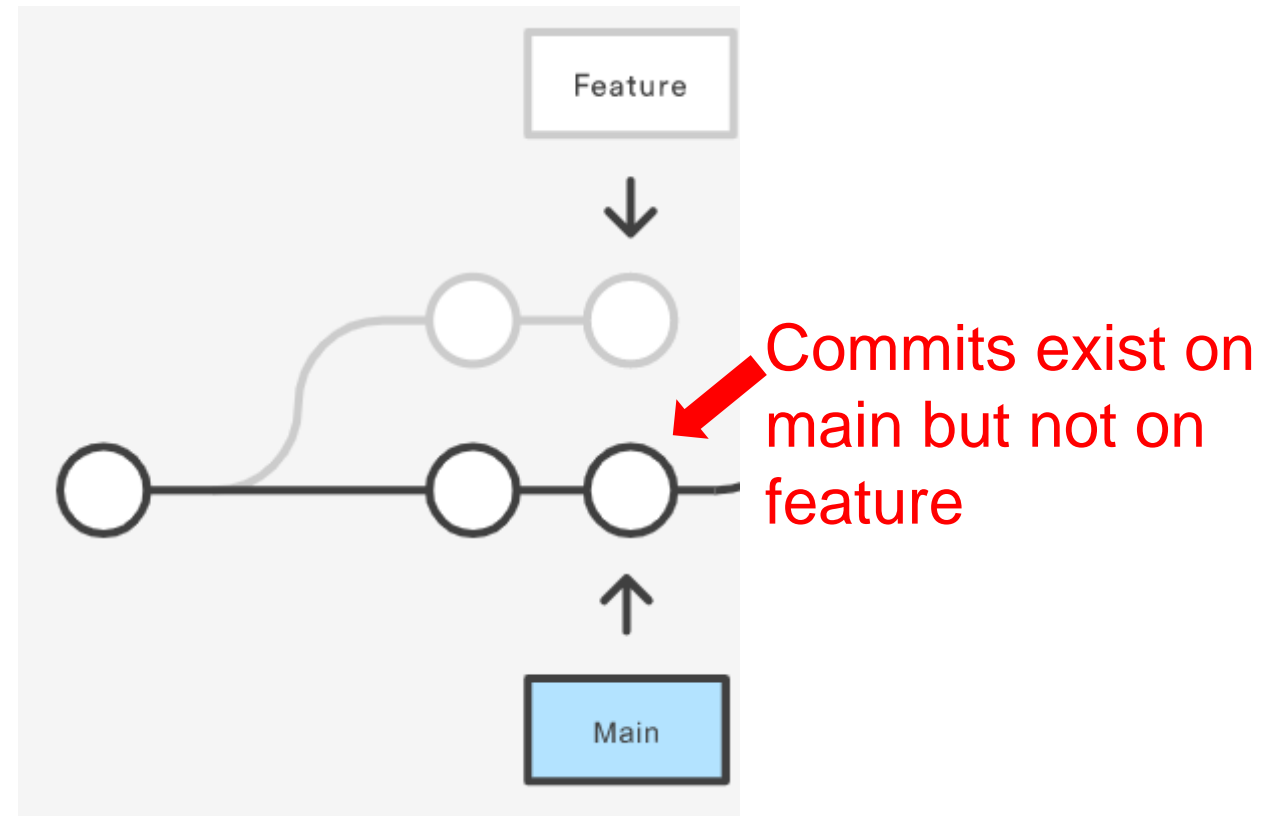
# Anything in scratch/ subdirectory
scratch/*
```

# Merge Conflicts

## *When Things Go Wrong*

- Occurs when merging branches with inconsistent commit history

In these cases, merging feature into main will often cause merge conflicts



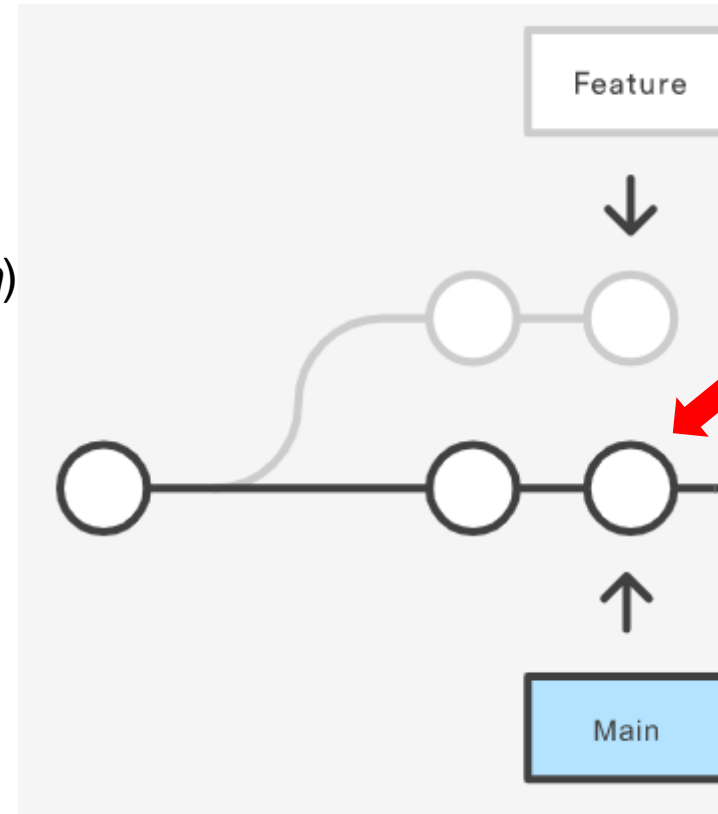


# Merge Conflicts

## *When Things Go Wrong*

- Occurs when merging branches with inconsistent commit history
- Avoid by:
  - Don't commit to main (*esp. with active branches*)
  - 1 branch = 1 person
  - Otherwise, communicate & pull (*often*)
  - Rebase (*details [here](#)*)

In these cases, merging feature into main will often cause merge conflicts



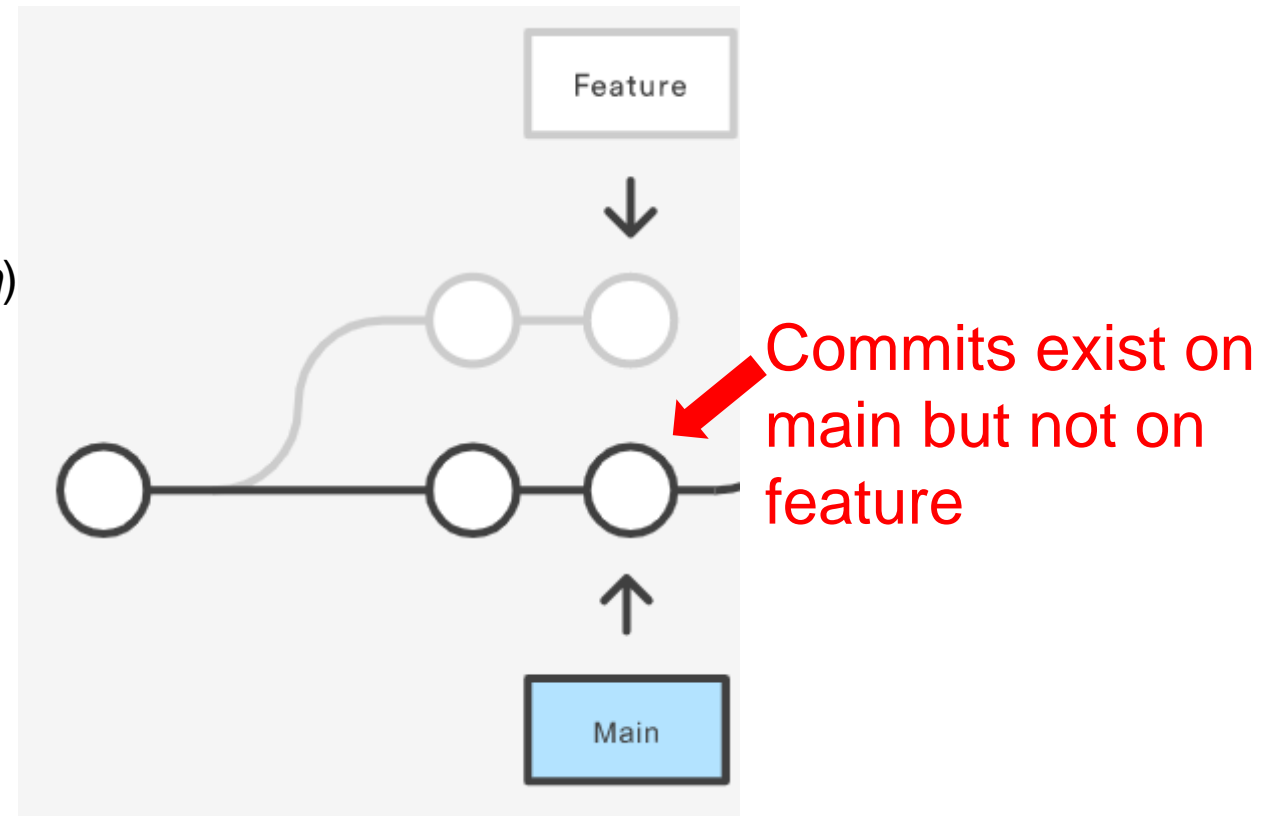
Commits exist on main but not on feature

# Merge Conflicts

## *When Things Go Wrong*

- Occurs when merging branches with inconsistent commit history
- Avoid by:
  - Don't commit to main (*esp. with active branches*)
  - 1 branch = 1 person
  - Otherwise, communicate & pull (*often*)
  - Rebase (*details [here](#)*)
- Resolve by selecting which version of conflicting file(s) to keep after merge

In these cases, merging feature into main will often cause merge conflicts



# R Packages

*Have an R Package? Host it on GitHub!*

---

- R packages are a clean way of sharing code that will be called repeatedly by multiple users
- You could
  - Host it on CRAN (non-trivial)
  - Share the compiled package source files
  - Host it on GitHub!
- Users install via 1-2 lines of code:

```
# install remotes package
install.packages("remotes")

# install package from GitHub
remotes::install_github("user/packageName")
```



Free ebook: <https://r-pkgs.org/>

# Cite Your Code

*Need to make a version permanent?*

---

- E.g., for journal article supplementary material
- [Zenodo](#): archival service

## Example of Staton et al. (2021)

Accounting for uncertainty when estimating drivers of imperfect detection: An integrated approach illustrated with snorkel surveys for riverine fishes

Article [10.1016/j.fishres.2021.106209](https://doi.org/10.1016/j.fishres.2021.106209)

Code [10.5281/zenodo.3928691](https://doi.org/10.5281/zenodo.3928691)



# Cite Your Code

*Need to make a version permanent?*

---

- E.g., for journal article supplementary material
- [Zenodo](#): archival service
- Steps ([more details](#)):
  1. Link Zenodo to GitHub account
  2. Make your GitHub repo public
  3. Sync Zenodo & turn on tracking
  4. Publish a GitHub [release](#)
  5. Each release will receive a DOI, archiving that exact version

## Example of Staton et al. (2021)

Accounting for uncertainty when estimating drivers of imperfect detection: An integrated approach illustrated with snorkel surveys for riverine fishes

Article [10.1016/j.fishres.2021.106209](https://doi.org/10.1016/j.fishres.2021.106209)

Code [10.5281/zenodo.3928691](https://doi.org/10.5281/zenodo.3928691)



# Practice Makes Perfect!



(And gives confidence that you won't mess something up)

*But don't be too perfect.*

*Time spent documenting changes is time spent not making changes (progress).*

**Strongly** *advise practicing on a toy analysis before applying to real new or existing analyses. Get the feel for how this works.*

***Upcoming PNAMP FWMG Webinar  
A "How-To" Tutorial***

<b>Term</b>	<b>Description</b>
<b>1. Repository</b>	A directory (folder) with files under version control ( <i>aka "repo"</i> )
<b>2. Local (Repository)</b>	A repository found on your hard drive
<b>3. Remote (Repository)</b>	A repository found online
<b>4. Clone</b>	Create a local repo from remote ( <i>one time only</i> )
<b>5. Commit</b>	A <u>small</u> set of changes that accomplish an "atomic task"
<b>6. Fetch</b>	Check if local is up-to-date with remote
<b>7. Pull ↓</b>	Sync commits to local from remote ( <i>if local is behind remote</i> )
<b>8. Push ↑</b>	Sync commits from local to remote ( <i>if remote is behind local</i> )
<b>9. Branch </b>	An isolated version of the repo ( <i>doesn't affect other branches</i> )
<b>10. Main Branch</b>	A special branch storing the "best" version
<b>11. Feature Branch</b>	A branch storing a version under development
<b>12. Merge </b>	Add commits from one branch to another
<b>13. Pull (merge) Request</b>	Perform a merge in a documented/collaborative fashion