# Collaborative Workflow for SAP

# When you think workflow

Workflow includes:

- Editor you use to write code

- Name of your home directory

- Git client you use

Products include:

- Raw data

- Packages used for analysis

- Functions you wrote

- Figures and code used to make those figures

- Code someone needs to run to reproduce your results

**Don't hard code anything about workflow into products**

# Openness as a spectrum

66 Open Science is the principle and practice of making research products and processes available to all, while respecting diverse cultures, maintaining security and privacy, and fostering collaborations, reproducibility, and equity.

- Varying degrees of openness

- Practice peer-review

- It's okay to share imperfect code

yourself    collaborators    SAP team    science centers    public

# Tools for Openness

- Use the tools we have to make your work as open as possible

- Google Drive - good for sharing data among collaborators

    - Use R code to pull directly from Google drive into R or onto personal computer

- Github repositories - can keep private until published

- Open science community

    - Openscapes training

    - NMFS R User Group

# Benefits of automated workflows

- Save time and energy

  - Can put your energy towards developing models and not the tedious details

- Make things uniform

- Makes onboarding easier

- Help future us

- Helps you avoid mistakes and easier to fix them

**Automate as much as possible!**

# Why SAP should adopt an open and automated workflow

- Anyone can download a repository and reproduce our results

- Makes it easier during WPSAR reviews to run reviewer requests

- Easier onboarding when joining new projects

- Helps share the responsibility among the team

- Share knowledge within the group

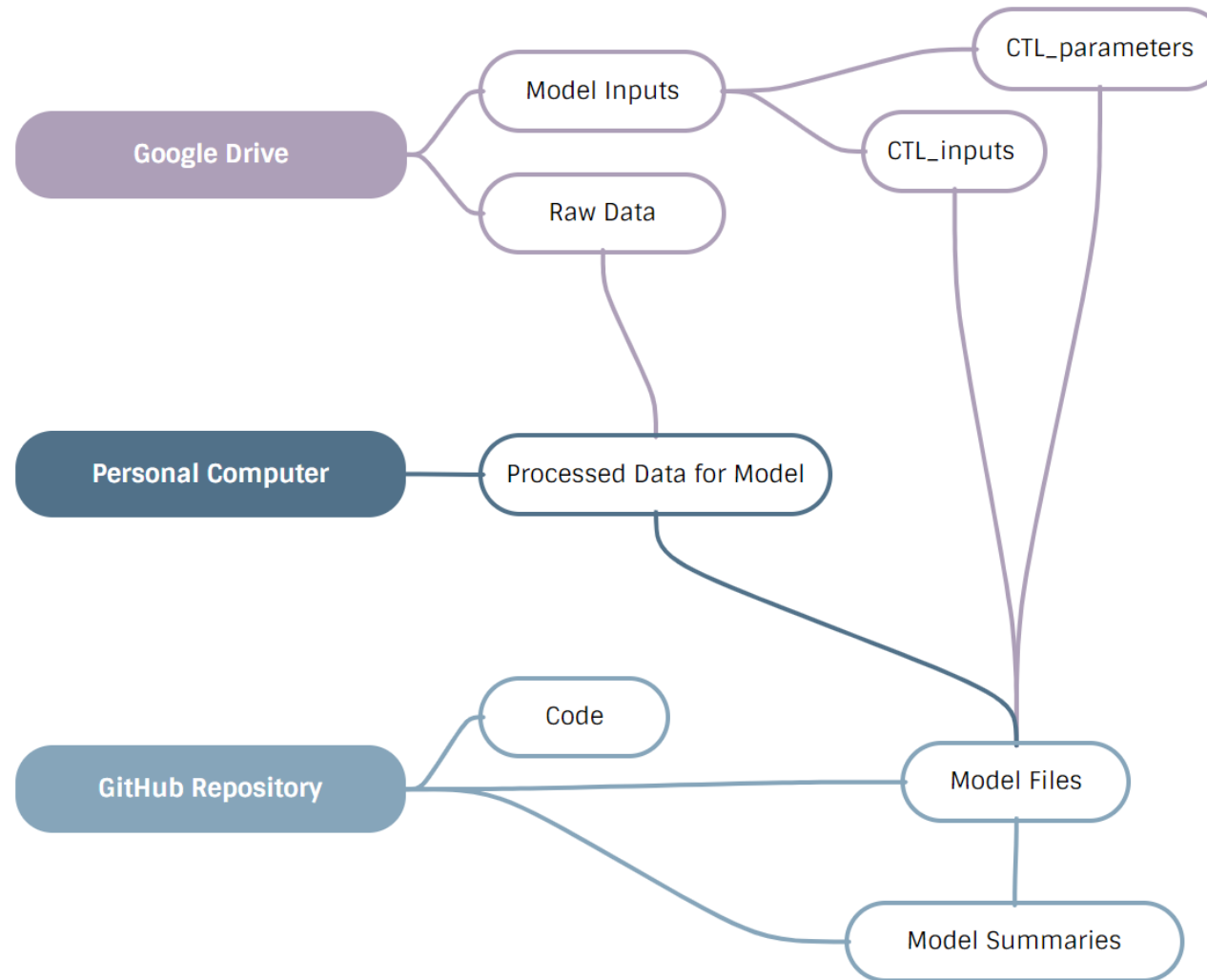# Challenges of American Samoa Bottomfish Assessment

## Problems

- Develop 9 individual models that have the same structure but different data and parameter values.

- Develop models simultaneously and update parameter values without having to upload and download new versions of a spreadsheet.

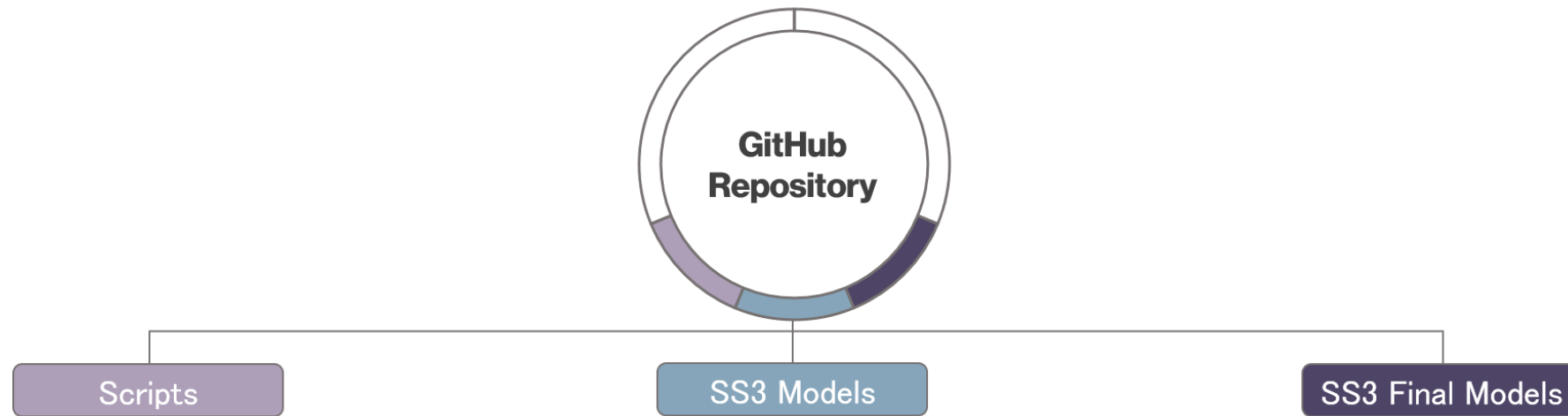- Share quick summaries of each model version easily with people.

## Solutions

- Wrote custom functions to build the model input files and include the correct data and parameter values for each species.

- Set up 2 Google Sheets (inputs and parameters) that we could change and pull the information directly into R with real-time updated information.

- Included pdf reports of summary plots and diagnostics for each model version with 4 input files needed to reproduce those plots.

# Our Workflow

# Repository Structure

# Repository Structure



GitHub Repository

Scripts
— 01_Run data scripts.R
— 02_Run SS models.R

01_Data scripts
— CPUE scripts
— Catch scripts
— Size scripts

02_SS scripts
— Build SS Models
— Run Forecasts
— Run Bootstraps
— Run Alternate Models

03_Report scripts
— Create Figures
— Create Tables

SS3 Models

SS3 Final Models

# Repository Structure



**GitHub Repository**

**Scripts**
- 01_Run data scripts.R
- 02_Run SS models.R

**01_Data scripts**
- CPUE scripts
- Catch scripts
- Size scripts

**02_SS scripts**
- Build SS Models
- Run Forecasts
- Run Bootstraps
- Run Alternate Models

**03_Report scripts**
- Create Figures
- Create Tables

**SS3 Models**

APVI

ETCO

LUKA

PRZO

Template Models

APRU
- Model 1
- Model 2

CALU

LERU

PRFL

VALO

**SS3 Final Models**

# Repository Structure

# Step 01: Retrieve raw data

```r
1   library(googledrive)
2   library(googlesheets4)
3
4   ########## DOWNLOAD DATA FROM GOOGLE DRIVE ##############
5   # Check latest data from Google Drive but only download if its more recent than on local repo
6   a <- drive_reveal(
7     drive_ls(
8     path="https://drive.google.com/drive/u/1/folders/1pnH38cupmDU4O_KkKDhYWee_p4sTSD6u",
9     pattern="Data"),
10    what = "modified_time")
11  a <- arrange(a, by = desc(modified_time))[1,] # Select most recent "Data" zip file
12
13  if(dir.exists(file.path(here(..=1),"Data"))){
14          Date.CurrentFolder <- as_datetime(
15              file.info(paste0(file.path(here(..=1)),"/Data"))$mtime)
16  } else {
17    Date.CurrentFolder <- "1900-01-01 01:01:01 UTC"
18    }
19
20  Date.GoogleFolder <- as_datetime(map_chr(a$drive_resource, "modifiedTime"))
21
22  if(Date.CurrentFolder<Date.GoogleFolder){
23     drive_download(file=a$id,
24                 overwrite = TRUE,
25                 path = file.path(here(..=1),a$name))
26
```

# Step 02: Process data

```r
1   ########## PROCESS CATCH, CPUE, AND SIZE DATA ###############
2   set.seed(123)
3   source(paste0(here(..=1),"/Scripts/01_Data scripts/01_CPUE_BBS_InitPrep.r"));    rm(list=ls())
4   source(paste0(here(..=1),"/Scripts/01_Data scripts/02_CPUE_BBS_PropTable.r"));   rm(list=ls())
5   source(paste0(here(..=1),"/Scripts/01_Data scripts/03a_CPUE_BBS_Wind.r"));        rm(list=ls())
6   source(paste0(here(..=1),"/Scripts/01_Data scripts/03b_CPUE_BBS_PCA.r"));         rm(list=ls())
7   source(paste0(here(..=1),"/Scripts/01_Data scripts/04_CPUE_BBS_FinalPrep.r"));   rm(list=ls())
8   source(paste0(here(..=1),"/Scripts/01_Data scripts/06_CATCH_BBS_FinalPrep.r")); rm(list=ls())
9   set.seed(123)
10  source(paste0(here(..=1),"/Scripts/01_Data scripts/07_CATCH_SBS_PropTable.r")); rm(list=ls())
11  source(paste0(here(..=1),"/Scripts/01_Data scripts/08_CATCH_SBS_FinalPrep.r")); rm(list=ls())
12  source(paste0(here(..=1),"/Scripts/01_Data scripts/09_CATCH_Final.r"));          rm(list=ls())
13  source(paste0(here(..=1),"/Scripts/01_Data scripts/10_SIZE.r"));                 rm(list=ls())
14
15  ############### RUN CPUE STANDARDIZATION#########################
16  # Run CPUE standardization and export indices for input into SS
17  source(paste0(here(..=1),"/Scripts/01_Data scripts/05_CPUE_BBS_Standardize_Function3.r"))
18  source(paste0(here(..=1),"/Scripts/01_Data scripts/05_CPUE_BBS_Standardize_Function2.r"))
19
20  # Run CPUE standardization for all species, areas combined in a loop
21  root_dir <- root_dir <- this.path::here(.. = 1)
22  Species.List <- c("APRU","APVI","CALU","ETCO","LERU","LUKA","PRFL","PRZO","VALO")
23  for(i in 1:length(Species.List)){
24      Standardize_CPUE3(Sp=Species.List[i],Interaction=T,minYr=1988,maxYr=2015)
25  }
```

# Step 03: Set input arguments

## Arguments we rarely changed

```
1  Lt        <-vector("list",9) # Species options
2
3  Lt[[1]]<-list("APRU",          #Name
4                "SW_Then",        #M
5                "SW_BBS_BIOS",    #Growth
6                "Kamikawa",       #Length-Weight
7                "SW_BBS_BIOS",    #Maturity
8                F,                #Use InitF?
9                c(0.5,1.6),       #Range of R0 prof
10               0.29,             #Btarg. value
11               T,                #Include superyea
12               list(c(2019,2020)),#Blocks of supery
13               c(2.5,5.5,0.2))    #Proj catch range
```

## Arguments we changed frequently

```
1  DirName     <- "65_Base" # Name of directory
2  runmodels   <- T    # Run ss.exe
3  printreport <- T    # Run ss_diags report
4  Create_species_report_figs <- F # Produce formatte
5                                  #figures and table
6                                  #word document
7  N_boot      <- 0    # Bootstrap on/off
8  N_foreyrs   <- 0    # Nyears for forecast
9  RD          <- F    # Run diagnostics
10 ProfRes     <- .1 # R0 profile resolution
11 profile     <- "SR_LN(R0)" # Parameter to profile
12 Begin       <- c(1967,1986)[1] #Model start year
13 DeleteForecastFiles <- T #Remove extra files
```

# Step 04: Build models and run

- The core function is `Build_All_SS()`

- Wrapper function for modular functions:

  - `Build_Data()` , `Build_Control()` , `Build_Starter()` , `Build_Forecast()`

- Incorporate parameter information stored in Google Sheets

```
1   if(readGoogle == T){
2       ctl.params <- read_sheet("1XvzGtPls8hnHHGk7nmVwhggom4Y1Zp-gOHNw4ncUs8E",
3                             sheet=species)
4       ctl.inputs <- read_sheet("11lPJV7Ub9eoGbYjoPNRpcpeWUM5RYl4W65rHHFcQ9fQ",
5                             sheet=scenario)
6
7   }else{
8     ctl.inputs <- readxl::read_excel(file.path(root_dir, "Data", "CTL_inputs.xlsx"),
9                             sheet=scenario)
10    ctl.params <- readxl::read_excel(file.path(root_dir, "Data", "CTL_parameters.xlsx"),
11                         sheet=species)
12  }
```

# Modular code

- Created "smaller" functions for each step
- Used the main function as a wrapper so could run everything from 1 script and easy to turn functions on and off
  - bootstraps
  - diagnostics
  - forecasting
  - alternate models
  - summary reports
  - creating specific plots and tables

📁 01_Data scripts

📁 02_SS scripts

📁 03_Report scripts

📁 04_Other scripts

📄 01_Run data scripts.r

📄 02_Run SS models.R

# Step 05: Produce model diagnostics reports

- `r4ss::SS_plots()`
- Created a custom qmd with model diagnostic plots and tables
  - runs test, likelihood profiles, retrospectives, jitters, etc.
- Created html and pdf versions
  - PDFs were synced with the repository so anyone could look at the diagnostics for a certain model without cloning the repository onto their computer

# Step 06: Generate formatted tables and figures

Quarto document produced a .docx with figures and tables formatted for the report

```r
library(openxlsx)
library(flextable)

set_flextable_defaults(
  font.size = 11, font.family = "Arial",
  border.color = "black", font.color = "black")
border <- officer::fp_border()

refpoints <- read.xlsx(file.path("./", "scr", "01_tables.xlsx"),
                       sheet = "02_RefPoints")
refpoints.ft <- flextable(refpoints)

refpoints.ft <- set_header_labels(refpoints.ft, REF_POINT = "Reference point", ALL = "Value")
refpoints.ft <- hline_bottom(refpoints.ft, border = border)
refpoints.ft <- hline_top(refpoints.ft, border = border, part = "header")
refpoints.ft <- hline_bottom(refpoints.ft, border = border, part = "header")

refpoints.ft <- compose(refpoints.ft, i = 1, j = 1, part = "body",
                        value = as_paragraph(as_i("F"), as_sub("MSY"),
                                             " (yr", as_sup("-1"), ")"))

refpoints.ft <- compose(refpoints.ft, i = 4, j = 1, part = "body",
                        value = as_paragraph(as_i("SSB"), as_sub("MSST")))
```

# Formatted tables for report

| | A | B |
|---|---|---|
| 1 | REF_POINT | ALL |
| 2 | Fmsy | 0.136 (0.127 - 0.146) |
| 3 | F2021 | 0.001 (0 - 0.002) |
| 4 | F2021/Fmsy | 0.006 (0.003 - 0.014) |
| 5 | SSBmsst | 4.6 (2.47 - 8.55) |
| 6 | SSB2021 | 14.2 (7.06 - 26.88) |
| 7 | SSB2021/SSBmsst | 3.12 (2.56 - 3.8) |
| 8 | MSY | 2.16 (0.89 - 3.43) |
| 9 | Catch2019-2021 | 0.51 (0.11 - 0.9) |
| 10 | SPRmsy | 0.36 (0.36 - 0.36) |
| 11 | SPR2021 | 0.99 (0.98 - 1) |
| 12 | | |
| 13 | | |
| 14 | | |

01_Quants    **02_RefPoints**    03_Sensitivies    (+)

| Reference point | Value |
|---|---|
| $F_{MSY}$ (yr$^{-1}$) | 0.136 (0.127 - 0.146) |
| $F_{2021}$ (yr$^{-1}$) | 0.001 (0 - 0.002) |
| $F_{2021}/F_{MSY}$ | 0.006 (0.003 - 0.014) |
| $SSB_{MSST}$ | 4.6 (2.47 - 8.55) |
| $SSB_{2021}$ (mt) | 14.2 (7.06 - 26.88) |
| $SSB_{2021}/SSB_{MSST}$ | 3.12 (2.56 - 3.8) |
| MSY (mt) | 2.16 (0.89 - 3.43) |
| Catch$_{2019-2021}$ (mt) | 0.51 (0.11 - 0.9) |
| $SPR_{MSY}$ | 0.36 (0.36 - 0.36) |
| $SPR_{2021}$ | 0.99 (0.98 - 1) |

# Lessons learned in our process

- Use packages that work with all platforms of R

- Make functions flexible enough to give you options but don't worry about coding every single possibility

- Develop as you go, streamline as you go

- Make code modular

# Small changes you can implement now

- More comments, or use software such as Quarto or RMarkdown for literate programming

- Read other people's code

- Make sure you are using relative paths and including all the packages you used in that script

- When naming scripts, give them informative names and try sequencing them

- Get in the habit of syncing code with repository regularly

# Discussion points

- Is this system useful for single-species assessments?
    - I think its useful for any assessment with more than 1 person working on it
- How can we help each other with setting these systems up?
    - coworking time in small groups
    - set up GitHub repository templates for a common set-up
- How can this be useful for international assessments?
- What about reference manager software?
- What are some good tools people use for project management?
    - GitHub Projects