

# DSCI351-351m-451: Class 01a, ODS Tool Chain (CWRU, Pitt, UCF)

Profs: R. H. French, Pawan Tripath, P. Leu, M. Li, K. Davis

TAs: Redad Medhi, Olatunde Akanbi

April 16, 2025

## Contents

1.1.1.1	Class Readings, Assignments, Syllabus Topics . . . . .	1
1.1.1.1.1	Reading, Lab Exercises, SemProjects . . . . .	1
1.1.1.1.2	Textbooks . . . . .	2
1.1.1.1.3	Syllabus . . . . .	2
1.1.1.1.4	Prof. Pawan Tripath will present in class Thursday, on SemProjs . .	2
1.1.1.2	The Lab Exercises (LEs) . . . . .	2
1.1.1.3	Where we are at present in Class . . . . .	4
1.1.1.4	Markov HPC and Open Data Science (ODS) Compute Engines . . . . .	5
1.1.1.5	What we need to do now . . . . .	7
1.1.1.5.1	So go make accounts, using your case.edu email address . . . . .	9
1.1.1.6	Your Open Data Science Tool Chain . . . . .	9
1.1.1.6.1	Its all about a Data Science Tool Chain . . . . .	9
1.1.1.6.2	Online Git Server Communities . . . . .	9
1.1.1.6.3	Slack, another component of Agile Software Development . . . . .	11
1.1.1.7	Your Online Data Science Portfolio . . . . .	11
1.1.1.7.1	Sign up for a Stack Exchange Account . . . . .	12
1.1.1.7.2	Efficiently browse you SX sites . . . . .	12
1.1.1.7.3	An Example, Emeline Liu . . . . .	12
1.1.1.8	Links . . . . .	12

### 1.1.1.1 Class Readings, Assignments, Syllabus Topics

#### 1.1.1.1.1 Reading, Lab Exercises, SemProjects

- Readings:
  - For today:
  - For next class: Peng R Programming, pages 4-33
- Laboratory Exercises:
  - LE1 : Given out Thursday
  - LE1 : Is Due Tuesday Sept. 10th
- Office Hours: (Class Canvas Calendar for Zoom Link)
  - Mondays @ 5:00 PM to 6:00 PM, Redad Medhi
  - Wednesday @ 3:00 PM to 5:00 PM, Olatunde Akanbi
  - **Office Hours are on Zoom, and recorded**
- Semester Projects
  - DSCI 451 Students Biweekly Update Due

- DSCI 451 Students
  - \* Next **Report Out #4 (Full Report) is Due December 12th**
- All DSCI 351/351M/451 Students:
  - \* **Peer Grading of Report Out # is Due**
- SemProject Office Hours: (Class Canvas Calendar for Zoom Link)
  - \* Prof. Pawan Tripathi guides the SemProjs for DSCI451 students
    - @Abhishek Daundkar [aad157@case.edu](mailto:aad157@case.edu)
  - \* SemProject Office Hours. (twice a week during the first month)
    - Tuesdays @ 3:00 PM
    - Thursdays @ 2:00 PM
- Exams
  - Final: Monday December 16, 2024, 12:00PM - 3:00PM, Nord 356 or remote

#### 1.1.1.1.2 Textbooks

- Introduction to R and Data Science
  - For R, Coding, Inferential Statistics
    - \* Peng: R Programming for Data Science
    - \* Peng: Exploratory Data Analysis with R

Textbooks for this class

- OIS = Diez, Barr, Cetinkaya-Runde: Open Intro Stat v4
- R4DS = Wickham, Grolemund: R for Data Science

Textbooks for DSCI353/353M/453, And in your Repo now

- ISLR = James, Witten, Hastie, Tibshirani: Intro to Statistical Learning with R 2nd Ed.
- ESL = Trevor Hastie, Tibshirani, Friedman: Elements of Statistical Learning
- DLwR = Chollet, Allaire: Deep Learning with R

Magazine Articles about Deep Learning

- DL1 to DL13 are “Deep Learning” articles in 3-readings/2-articles/

#### 1.1.1.1.3 Syllabus

#### 1.1.1.1.4 Prof. Pawan Tripath will present in class Thursday, on SemProjs

- To give more information on the Semester Projects for DSCI453 students
  - This includes 3 Reports Outs by 453 Students
  - That **all students will view and do peer grading of**
  - The SemProj TA is Abhishek Daundakar
    - \* White 540, [aad157@case.edu](mailto:aad157@case.edu)

#### 1.1.1.2 The Lab Exercises (LEs)

- Each LE is worth
  - LE1,2 are 7 points
  - LE3-7 are 10 points
    - \* (except LE0 = 0 points)

So 64 points are in the Lab Exercises

- So these are important and critical to learning
- You will need to start on the early
  - This is why you are given two weeks to do them
- You turn in both the .Rmd and the .pdf file

Day:Date	Foundation	Practicum	Reading	Due
w01a:Tu:8/27/24	ODS Tool Chain	R, Rstudio, Git		
w01b:Th:8/29/24	Knuth-Literate Prog.	Bash, Git, Slack, Agile	PRP4-33	LE1
w02a:Tu:9/3/24	ODS Setup	RIntroR	PRP35-64	
w02b:Th:9/5/24	How Git Operates	OIS:Intro2R	OIS1,2	
w02Pr:Fr:9/6/24			PRP65-93	451 Update1
w03a:Tu:9/10/24	Intro to Data Science	Data Analytic Style	PRP94-116	LE2 LE1 Due
w03b:Th:9/12/24	Intro to Data	R, Tidy Example	OIS4	
w04a:Tu:9/17/24	Summarizing Data	Rmd Paths Loops Tidy	EDA1-31	
w04b:Th:9/19/24	Rand. Var. Normal Dist.	CapMinder Tidy EDA	R4DS1-3	LE3 LE2 Due
w04Pr:Fr:9/20/24			EDA32-58	451 Update2
w05a:Tu:9/24/24	OIS4 Rand. Var.	PET Degr EDA, Hyper-spec	OIS5	
w05b:Th:9/26/24	OIS5 Found. of Infer.	Anscombe's Quartets	R4DS4-6	
w05Pr:Fr:9/27/24				451 RepOut1
w06a:Tu:10/1/24	Pred., Algorithm, Model	Multivar Corr. Plot	R4DS7-8	
w06b:Th:10/3/24	Tidy Data	Summary Stats	R4DS9-16	LE4 LE3 Due
w06Pr:Fr:10/4/24				451 Update3
w07a:Tu:10/8/24	Midterm Rev.	Penguin EDA	OIS6.1-2	PeerRv1 Due
w07b:Th:10/10/24	HypoTest, Infer. Recap	Sampling		
w08a:Tu:10/15/24	Programming & Coding	Code Packaging	R4DS17-21	LE4 Due
w08a:Th:10/17/24	<b>MIDTERM</b>	<b>EXAM</b>		
w08Pr:Fr:10/18/24				451 Update4
Tu:10/21/24	<b>CWRU</b>	<b>FALL BREAK</b>		
w09b:Th:10/24/24	Cat. Inf. 1 & 2 propor.	Indep. Test, 2-way tables	OIS6.3-4	LE5
w09Pr:Fr:10/25/24				451 RepOut2
w10a:Tu:10/29/24	Goodness of Fit, $\chi^2$ test	t-tests 1&2 means	OIS7.1-4	
w10b:Th:10/31/24	Num. Infer, Cont. Tables	Stat. Power		
w10Pr:Fr:11/1/24				451 Update5
w11a:Tu:11/5/24	Sample & Effect Size	Stat. Power GGmap	OIS8	PeerRv2 Due
w11b:Th:11/7/24	Regr Part 1, Test & Train	Curse of Dimen.	ISLR1,2,1,2	LE6 LE5 Due
w12a:Tu:11/12/24	Regr. Outliers	Regr Part 2, GIS	OIS9	
w12b:Th:11/14/24	Mult.Regr., Var. Select	Regr. Diagnostics		
w12Pr:Fr:11/15/24				451 Update6
w13a:Tu:11/19/24	Log. Regr.	Mult. Regression	ISLR3.1	LE7 LE6 due
w13b:Th:11/21/24	Statistical learning	Log. Regr.	ISLR4.1-3	
w14a:Tu:11/26/24	Classificat.	Caret, Broom 4 modeling	ISLR3.2	
w14Pr:We:11/27/24	Sup. Unsup. Lrning			451 RepOut3
Th:10/28,29	<b>CWRU</b>	<b>THANKSGIVING BRK</b>		
w15a:Tu:12/3/24	Final Exam Review	Caret, Broom 4 modeling	Khalil.2020	PeerRv3 Due
w15b:Th:12/5/24	Big Data Analytics	Dist. Comp., Hadoop	Fr.Br.2020	LE7 due
<b>Friday 12/13/23</b>	<b>SemProj</b>	<b>Final Report</b>		<b>SemProj4 due</b>
<b>Monday 12/16/23</b>	<b>FINAL EXAM</b>	<b>12:00-3:00pm</b>	Nord 356	or remote

Table 1: DSCI351-451 Weekly Syllabus. w01a is week 1, class a. w01b is week 1 class b. w02Pr is DSCI451 SemProj. Readings are defined by book and chapters, sections in Peng R Prog. (PRPx.y), Peng Exp. Data An. (EDAx.y), R for Data Sci. (R4DSx.y), Open Intro Stats (OISx.y) & Intro. to Stat. Learn. with R (ISLRx.y).

Figure 1: DSCI351-351M-451 Syllabus

- We grade on the .pdf file in Canvas
- We expect good code styling
  - That matches the Google/Rstudio R Style Guide
  - Since this aides collaboration

## **Structure of DSCI 351/451 Course**

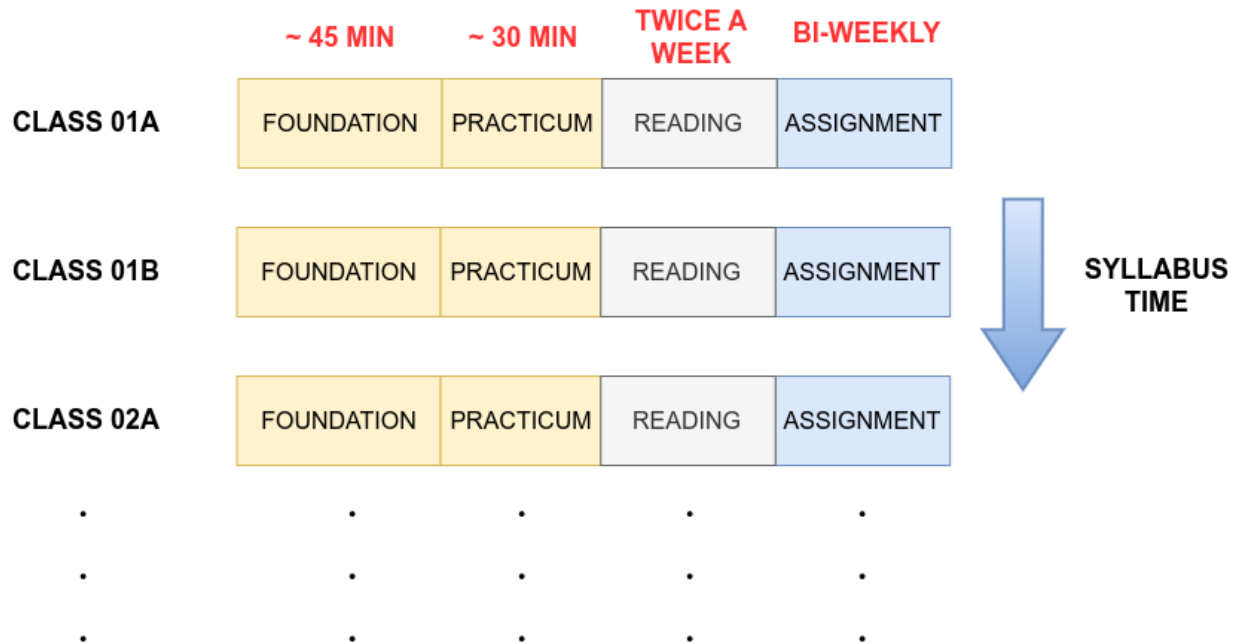


Figure 2: getting started

### **1.1.1.3 Where we are at present in Class** And getting familiar with data science

- So as of today,

We need to make all elements for the ODS tools chain working for you

- You have logged into your CaseID email at <http://webmail.case.edu>
  - And have setup Duo for Two Factor Authentication (2FA)
- You have joined the DSCI Slack
  - At <https://cwru-dsci.slack.com>
  - Using your [CaseID@case.edu](mailto:CaseID@case.edu) email
- You setup a bitbucket.org account
  - using your CaseID email account
  - And have setup your Bitbucket “App Password”
- You have “forked” the 24f-dsci351-451-prof “prof” repo
  - And have change “prof” to your caseID
  - And made your fork in the CWRU-DSCI team
- You have configured your git server
  - on both Markov, in your /home/CaseID/Git folder
    - \* and on ODS Desktop, in your H:/Git folder
    - \* and on your personal notebook computer, in a Git folder you make
  - And these configurations define your name and email
    - \* `git config --global user.name "[name]"`
    - \* `git config --global user.email "[email address]"`

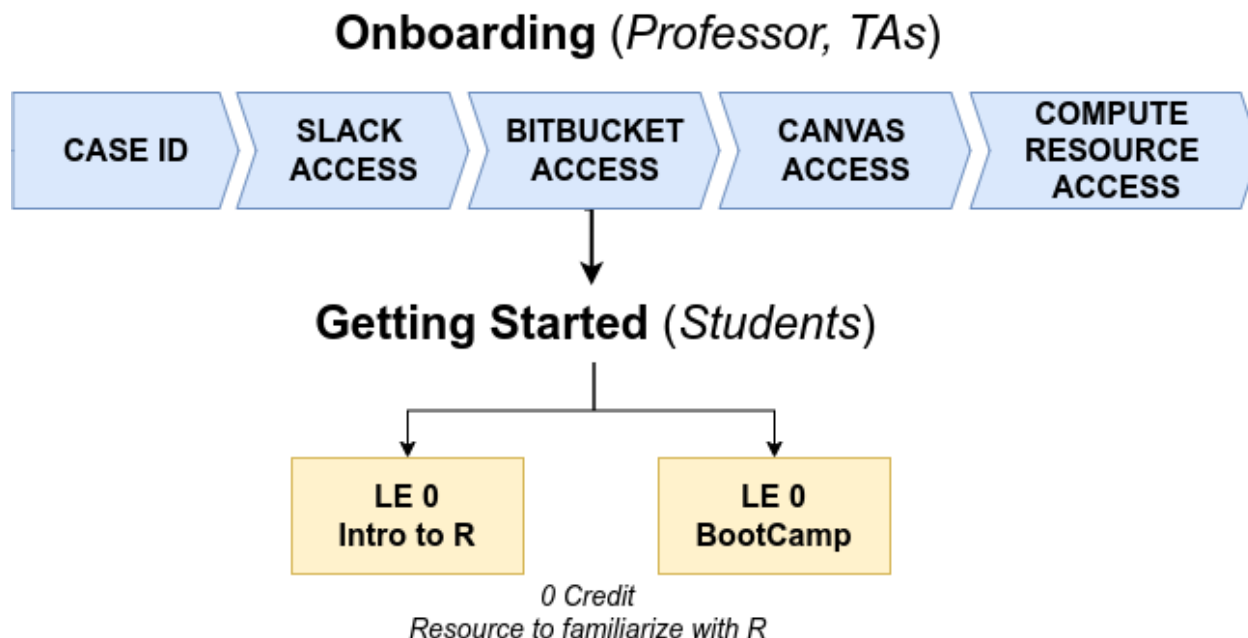


Figure 3: getting started

- Then you want to clone your personal course repo to 3 places
  - Markov/OnDemand: git clone... to /home/CaseID/Git/
  - ODS Desktop/MyApps: git clone... to H:/Git/
  - On your own computer to Git folder (to enable easy reading pdf)

If not, reach out to the TAs ( Redad Medhi, Hein Htet Aung, Nicole Lipa )

- Using the <http://cwru-dsci.slack.com>
  - Which you can join directly using your [CaseID@case.edu](mailto:CaseID@case.edu) email address
- Defining where you issue is
- And we'll fix it

#### 1.1.1.4 Markov HPC and Open Data Science (ODS) Compute Engines

- You can do data analysis on your notebook computer
  - You can setup your own notebook
    - \* For data science using R or Python
    - \* Full instructions are in the class syllabus
      - Section 11
    - \* For Linux, Mac's or Windows Operating Systems
    - \* But Many times you'll need more compute power than your notebook
      - Such as GPUs (Graphics Processing Units) to accelerate computations

But its useful to learn about a variety of Compute Resources

- In Class we'll use
  - Markov Data Science Cluster
    - \* A high performance computing cluster
    - \* via <http://ondemand.case.edu>

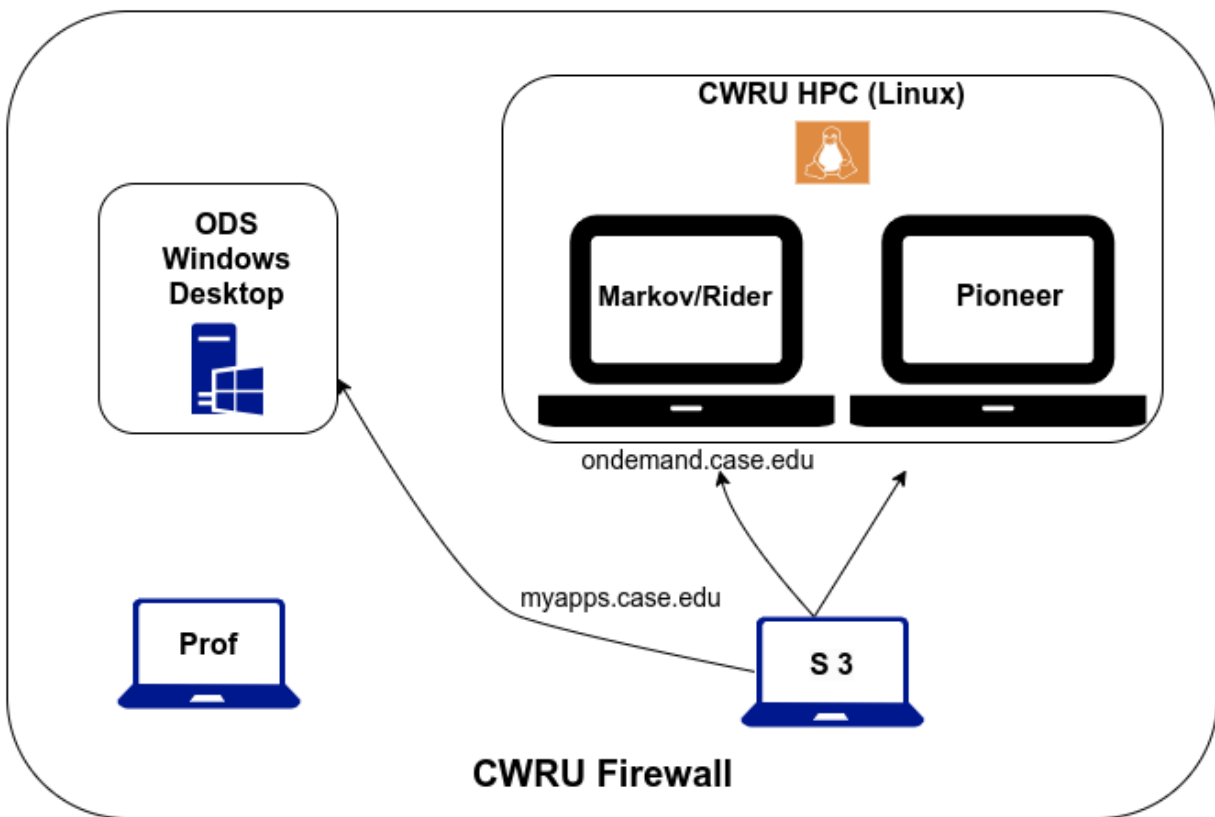


Figure 4: ODS Infrastructure

- or Open Data Science Desktops
  - \* A Win10 cloud desktop
  - \* via <http://myapps.case.edu>

These are all configured the same

- Independent of the Operating System
- They have R with Rstudio IDE (Integrated Development Environment)
- Git for code versioning
- LaTeX for publication quality report generation
- And also Python3 with VS Codium or PyCharm IDE

The two cloud computing systems: Markov HPC Cluster & ODS Win10 Desktop

- Markov Data Science HPC Compute Cluster, via OnDemand
  - Log in to <http://ondemand.case.edu>
  - Using your CaseID and password
  - Launch the Rstudio Server (rxf131)
    - \* Which runs R version 4.4.1
  - You can also get an XFCE graphical desktop on Markov

CWRU HPC provides Markov

- CWRU's HPC (High Performance Computing) Markov Cluster
  - This runs RedHat Linux version 7
  - Has 4400 CPU cores
  - Has 100,000 GPU cores
  - Up to a terabyte of Ram
- And has a new Data Science Cluster, named [Markov.case.edu](http://Markov.case.edu)
  - With a Hadoop Cluster for distributed computing
  - And dedicated GPUs
- You'll get accounts on CWRU HPC
- And use <http://ondemand.case.edu>
  - To login to Markov and get a Rstudio Server (rxf131) session
  - Or a xfce graphical desktop session
    - \* for simple file operations or a browser
- You also have access to the ODS Win10 Desktops
  - These are cloud Windows computers
    - \* That you log into from a Browser
    - \* login to <http://myapps.case.edu>
    - \* With your CaseID and password
  - The ODS VDIs are Windows 10 computers
  - The ODS VDIs don't have GPUs

Not for class, but for your own data science projects.

And you can use Google's Collaboratory](<https://colab.research.google.com/notebooks/welcome.ipynb>)

- For Jupyter Notebooks
- Running Python3
- Doesn't support R language yet
- Free GPUs and TPUs (Tensor Processing Unit)

#### 1.1.1.5 What we need to do now

- Setup our Markov and Open Data Science (ODS) Computers
  1. For Markov Data Science Cluster

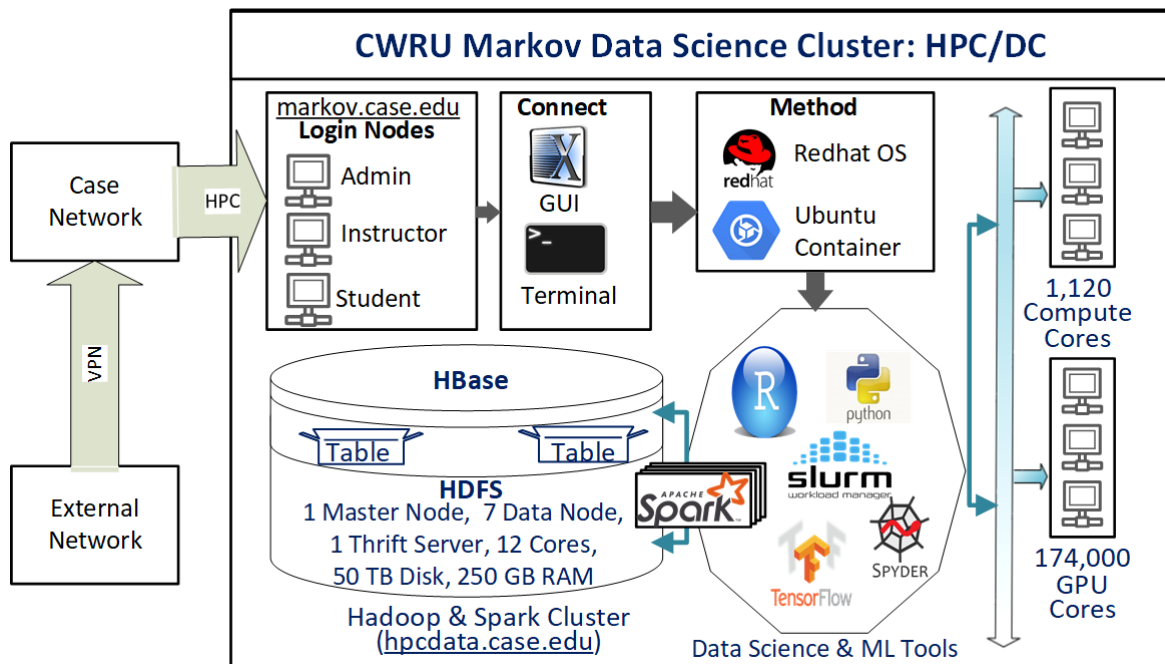


Figure 5: Markov Cluster

- login to <http://ondemand.case.edu> with your CaseID account
  - Launch the SDLE Rstudio Server (rxf131)
  - Check your “Library Paths”
    - \* in the R console
    - \* run `.libPaths()`
    - \* And the first directory MUST be
    - \* `“/home/rxf131/ondemand/ubuntu2004/r4” “/usr/local/lib/R/site-library”`
  - otherwise refer to the file in the root directory of your repo
    - \* named `FixRstudioServer-R-libPaths.txt`
    - \* and run the `“source(‘/home/rxf131/ondemand/share/config/r-lib-path-fix.R’)`
    - \* In the R console
    - \* then check your `.libPaths()` again
  - On Markov, launch XFCE Desktop (rxf131)
    - \* make a Git folder under `/home/CaseID/`
    - \* Login to DSCI Slack in your firefox browser on XFCE desktop
2. For the ODS Desktop
- login to <http://myapps.case.edu> with your CaseID account
  - Drag icons of to your desktop
    - \* for R, Rstudio, Git Bash, VScodium, PyCharm, Jupyter Notebook, Slack
3. Setup Git
- make `/home/caseID/Git` folder on Markov
    - \* git config your name and email of your git server
  - make `H:\Git` folder on ODS Desktop
    - \* git config your name and email of your git server
4. Git Fork the Class “Prof” Repo



- In your Bitbucket Account
- 5. Git Clone your Fork of the Class Repo
- 6. When in Rstudio (on Markov or ODS)
  - Its ESSENTIAL that you open the .Rproj file in the upper right corner
  - this tells Rstudio where your root directory of your project is.
- 7. Setup Bitbucket account
- 8. Setup [DSCI Slack Account](#)
- 9. Setup StackExchange account

#### 1.1.1.5.1 So go make accounts, using your case.edu email address

- Most students have already been invited
  - Pitt, UCF, UTRGV students have been issued CaseIDs
  - That you will use for logging in to
    - \* case email: at <http://webmail.case.edu>
    - \* Markov
    - \* ODS Desktop
    - \* DSCI Slack
    - \* CWRU Canvas
- Our DSCI Slack class channel
  - [CWRU Data Science Slack](#)
  - This is [an invite link to CWRU DSCI Slack](#)
- For you cloud Git server
  - [Bitbucket.org](#)
- A [Stack Exchange account](#)

#### 1.1.1.6 Your Open Data Science Tool Chain

##### 1.1.1.6.1 Its all about a Data Science Tool Chain

- Use R and build on the communities foundation
- Use Rstudio as a comfy environment
- Share your Open Data and Open Source Code
- Produce Reproducible Science with Rmarkdown
  - Use [Creative Commons Licenses](#)
  - Or other [Open Source Licenses](#)
  - Such as the [Gnu Public License: GPL](#)
  - Or one of my favorites, [the Apache License](#)

Pilot your Data Science studies using available data

- Find available data sets
- Before starting the costly process of making data

Use Git repositories

- For Code Version Control
- For Collaboration
- For Open Science sharing

##### 1.1.1.6.2 Online Git Server Communities

- We use [BitBucket Account](#)
  - In class, for our class code repositories

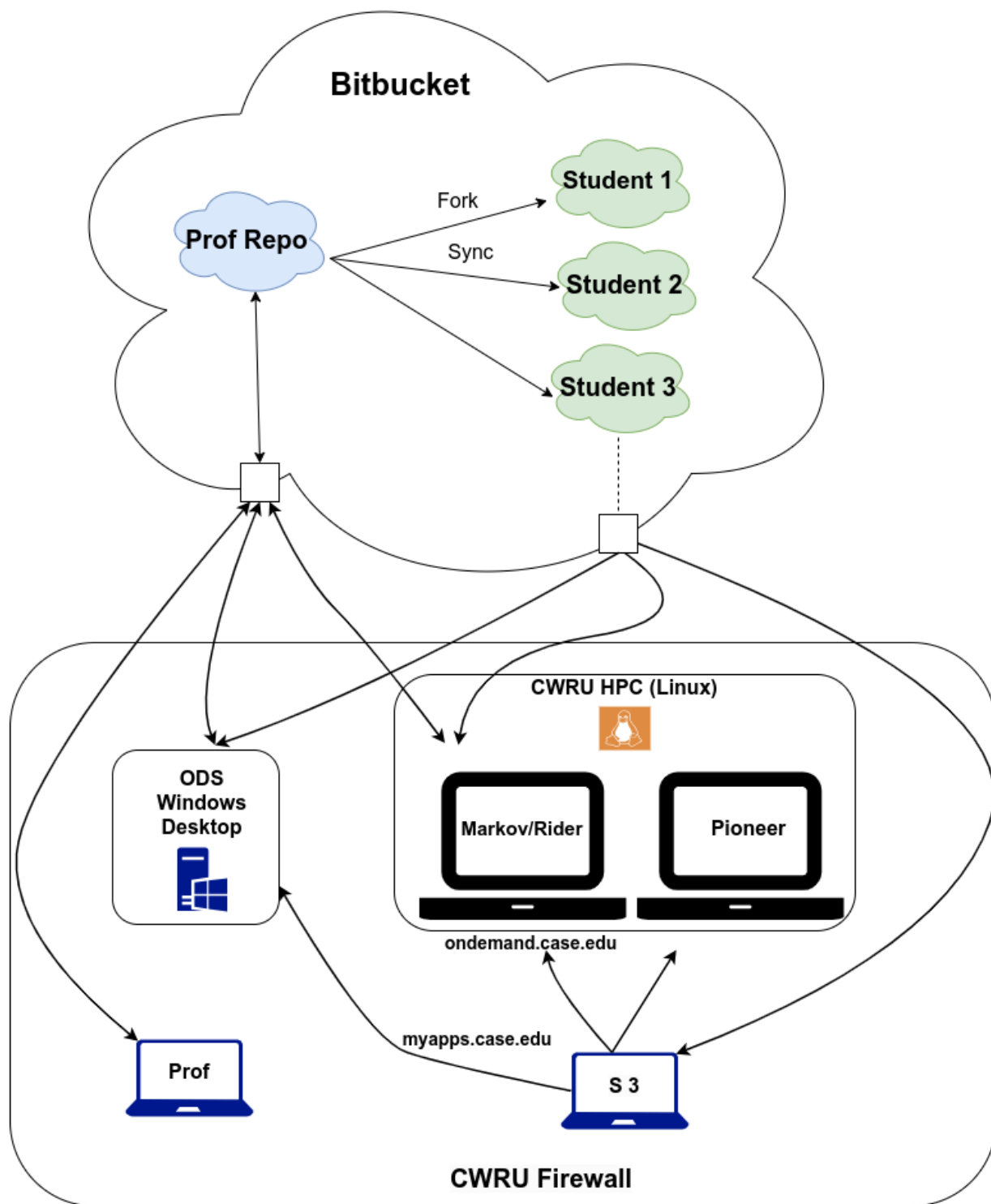


Figure 6: ODS Infrastructure

- These are private repositories
- You'll probably also want a [GitHub](#) account.
  - Many Rprojects are there, and
  - you can fork their repo's as inspect the code very easily.

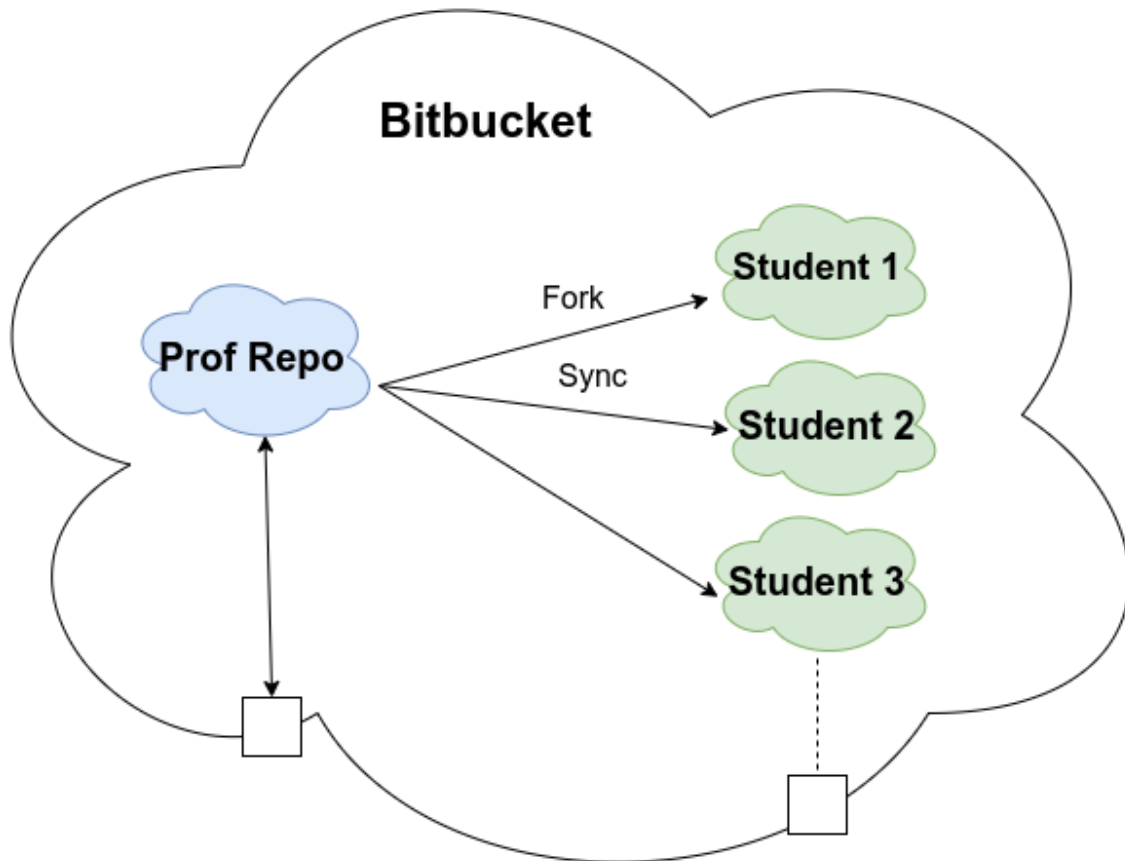


Figure 7: ODS Infrastructure

#### 1.1.1.6.3 Slack, another component of Agile Software Development

- [Slack.com](#)
  - We have a [CWRU DSCI Slack room](#)
  - There is Slack app for phones
  - And client for computers, its on vdi.
  - Slack client available for windows, mac and Linux
- an online collaboration tool

#### 1.1.1.7 Your Online Data Science Portfolio

- Doing open, reproducible data science
- Lets you share a portfolio of codes and projects
- Cite it in your resume
- Build a community of supporters and collaborators

#### 1.1.1.7.1 Sign up for a Stack Exchange Account

- Stack Exchange, Stack Overflow
  - are a Q&A community focused on many topics.

Stack Overflow allows you to search by tag

- `r` and `rmarkdown` are useful tags for Stack Overflow (SO)

[Stack Exchange's Tour of Stack Overflow](#)

Specific Stack Exchange websites

- for [SX Data Science](#)
- for [SX Statistics on Cross Validated](#)
- for [SX Open Data](#)

#### 1.1.1.7.2 Efficiently browse you SX sites

- Google (but more random)
- [The Stack Exchange apps](#)
- Using an [RSS Feed reader such as Feedly](#) is a good way

#### 1.1.1.7.3 An Example, Emeline Liu

- [emelineliu.com](#)
  - This website, which runs off of [Github Pages](#) and [Jekyll](#), is my latest project.
  - Right now, I'm using [Poole](#) as a foundation for my website/blog.

#### 1.1.1.8 Links

- <http://www.r-project.org>
- Rory Winston, for the [Learning R Intro](#)
- StackExchange <http://stackexchange.com/sites>
- Slack <http://slack.com>
- CWRU-DSCI Slack
- [emelineliu.com](#)
- [Github Pages](#)
- [Kaggle.com](#)
- [Colaboratory](#)