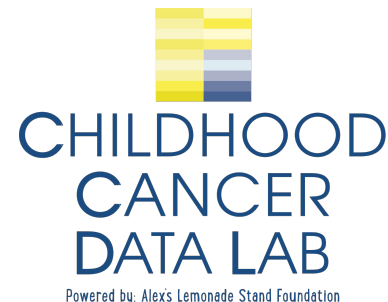
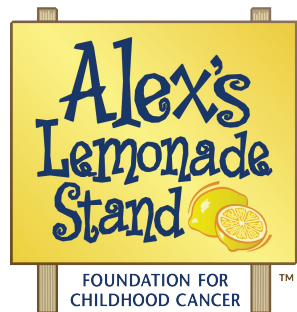




Welcome to the July 2020 Virtual CCDL Training Workshop!

July 27-31, 2020
Childhood Cancer Data Lab
<https://alexslemonade.github.io/2020-july-training/>



Meet your instructors



JOSH

Joshua Shapiro

Data Scientist @ the CCDL

PhD Ecology & Evolution, *UChicago*

Postdoc Integrative Genomics, *Princeton*

Research interests:

- **Evolutionary Genomics**



jashapiro

Meet your instructors



JACLYN

Jaclyn Taroni

Principal Data Scientist @ the CCDL

PhD Genetics *Dartmouth*

Postdoc Computational Biology *UPenn*

Research interests:

- **Transcriptomics in rare, complex diseases**
- **Unsupervised pattern extraction**



jaclyn-taroni

Meet your instructors



CANDACE

Candace Savonen

Biological Data Analyst @ the CCDL

Masters Neuroscience at *Michigan State University*

Research interests:

- **Neurogenomics**
- **Single-cell transcriptomics**



cansavvy

Meet your instructors



CHANTE

Chante Bethell

Biological Data Analyst @ the CCDL

Bachelor's in Bioinformatics from *Rowan University*

Research interests:

- **Functional motifs in the proteome**



cbethell

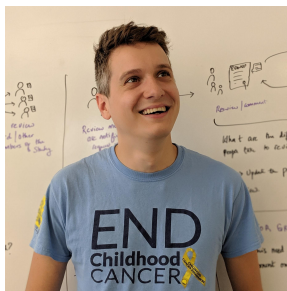
Other staff you may see



STEVEN
Steven Foltz

Postdoctoral Research Fellow
@ CCDL

- Interested in cancer genomics and tumor evolution
- Passionate about data science, visualization, and teaching



KURT
Kurt Wheeler

Data Engineer
@ CCDL

- Builds scalable systems
- Manages servers



DEEPA
Deepa Prasad

User Experience Designer
@ CCDL

- Talks to researchers about their needs and frustrations
- Designs usable software

Tell us about you!

- What's your name?
- What are you studying?
- What was the last movie, TV show, book, or album that you loved?





Code of Conduct



Be kind, have fun

We value the involvement of everyone in the community. We are committed to creating a friendly and respectful place for learning, teaching, and contributing.

- Use welcoming and inclusive language
- Be respectful of different viewpoints and experiences
- Gracefully accept constructive criticism
- Focus on what is best for the community
- Show courtesy and respect towards other community members

Read the full Code of Conduct here:

<https://alexslemonade.github.io/2020-july-training/code-of-conduct.html>



If you at any time feel harassed or treated inappropriately, please contact
ccd1@alexslemonade.org.

Monday

Workshop Intro

Single-cell RNA-seq

Technology overview
QC & normalization

Consultations

Exercise notebooks
Your own data

Wednesday

Pathway Analysis

Overrepresentation
GSEA

Consultations

Exercise notebooks
Your own data

Friday

Consultations

Your own data
Exercise notebooks

Presentations

Tuesday

Single-cell RNA-seq

Droplet analysis
Dimension reduction

Consultations

Exercise notebooks
Your own data

Thursday

Machine Learning

Clustering & heatmaps
PLIER

Consultations

Exercise notebooks
Your own data

Full schedule: <https://alexslemonade.github.io/2020-july-training/workshop/SCHEDULE.html>

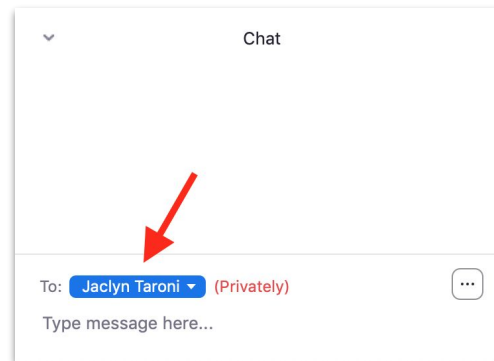


Virtual Training Procedures



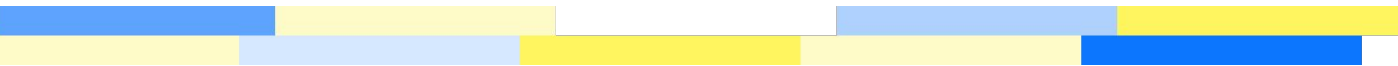
General Zoom Etiquette

- Keep your microphone muted
- Type questions in the Chat window (directed to the host)
 - Click on the “Chat” button at the bottom of your window to open the chat.
- Use the Zoom status buttons to tell us how you are doing!
 - Click on the “Participants” button at the bottom of your zoom window to see these buttons

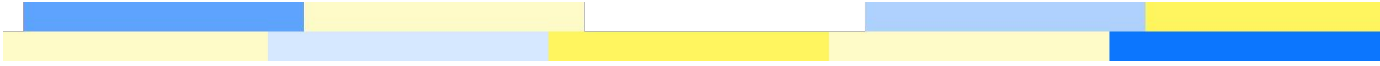


Using Slack & asking for help

- Use the **#2020-july-training** Slack channel
- Post public questions, get help with errors and debugging, make comments, and help your fellow participants!
 - Use threads to keep related content together
- Help us (and others) help you!
 - <https://alexslemonade.github.io/2020-july-training/workshop/posting-errors-guidelines.html>
 - If asking for help with an error, include the error message
 - Include what you tried, and code as appropriate
 - Use text, not screenshots (and take advantage of Slack's formatting tools)



What you will learn (and what you won't)




What you will learn

We will introduce you to the R programming language, R Notebooks, and some reproducible research practices.

We cover pipelines for the quality control, processing, and downstream analysis of bulk and RNA-seq data almost entirely through hands-on exercises.

We generally elect to go *broad* and not *deep*.

Our overarching goals: To prepare you to perform “frontline” analyses of your own data, to get you more comfortable reading documentation/learning new methods on your own, and to give you tools to collaborate more effectively with analysts when needed



What you won't learn


We don't address experimental design (e.g., how many replicates you need).

We won't compare tools (e.g., edgeR vs. DESeq2 for differential gene expression).

We won't cover every feature (or assumption) of the tools we do present.

You may not be able to perform every analysis you need to perform for your own work, particularly for complex experimental designs.

We present analysis as a series of *linear steps*. In practice, it's **not**. It's important to consult analysis experts when you need to and to keep track of and report what you've done.



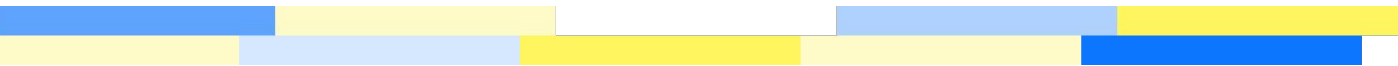
How do we pick what we teach?

We want methods to be or to have:

- Useful for a wide range of experimental designs, sample sizes
- Easy to use, well-documented, and consistently updated
- Solid tutorials, a sizeable user base, and responsive authors/maintainers

We have a preference for methods that integrate easily into a single workflow that can be run on a laptop (and our own personal biases as scientists).





Schedule



Daily Schedule Outline

Instruction

Full group

Lectures

- Introduce concepts and background
- Demonstrate usage
- Answer general questions

Breakout

Small groups

Start exercise notebooks

- Split up into Zoom breakout rooms
- Ask questions of instructors and other participants

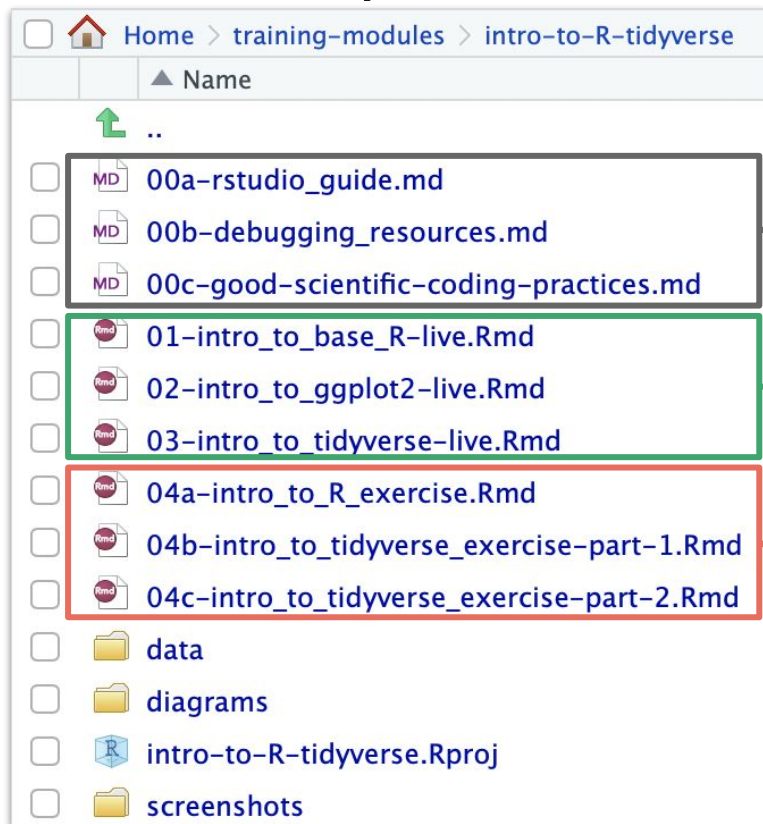
Consultation Period

Exercise notebooks

Your own data

- Practice what you have learned
- Work on exercises (at your own pace, or with others)
- Work with your own data

Module Layout



This is a reference document.
We will not go through this.

We'll walk through these notebooks
together, step-by-step

You will practice what you have
learned. We're here to help!

Module cheatsheets cover key functions

<https://github.com/AlexsLemonade/training-modules/tree/2020-july/module-cheatsheets>

dplyr

Read the `dplyr` package documentation [here](#).

A vignette on the usage of the `dplyr` package can be found [here](#).

Library/Package	Piece of code	What it's called	What it does
<code>dplyr</code>	<code>%>%</code>	Pipe operator	Funnels a <code>data.frame</code> through tidyverse operations
<code>dplyr</code>	<code>filter()</code>	Filter	Returns a subset of rows matching the conditions of the specified logical argument
<code>dplyr</code>	<code>arrange()</code>	Arrange	Reorders rows in ascending order. <code>arrange(desc())</code> would reorder rows in descending order.
<code>dplyr</code>	<code>select()</code>	Select	Selects columns that match the specified argument
<code>dplyr</code>	<code>mutate()</code>	Mutate	Adds a new column that is a function of existing columns
<code>dplyr</code>	<code>summarise()</code>	Summarise	Summarises multiple values in an object into a single value. This function can be used with other functions to retrieve a single output value for the grouped values. <code>summarize</code> and <code>summarise</code> are synonyms in this package.
<code>dplyr</code>	<code>rename()</code>	Rename	Renames designated columns while keeping all variables of the <code>data.frame</code>
<code>dplyr</code>	<code>group_by()</code>	Group By	Groups data into rows that contain the same specified value(s)
<code>dplyr</code>	<code>inner_join()</code>	Inner Join	Joins data from two data frames, retaining only the rows that are in both datasets.

Friday

Your own projects Exercise notebooks

Spend Friday working with your own data, getting assistance as needed from CCDL staff and each other.

Presentations

Present what you worked on during the consultation times to the group!

Communication during instruction



- I have an **urgent question** that needs an answer before moving on:
 - **Raise Hand** or **Chat** with the room host
- I'm **stuck with an error** and can't proceed with the hands-on exercise
 - **Chat** with meeting host: Request 1:1 and you will be placed in a breakout room with a CCDL staff member



- I have an **general question** that does not need an answer right away.
 - **Post** in #2020-july-training
- I'm having trouble **logging in** to RStudio Server
 - **Direct Message** a CCDL staff member (not the current host or instructor)

Trouble logging into Zoom and Slack? **Email** training@ccdatalab.org

Communication at other times (consultation periods)

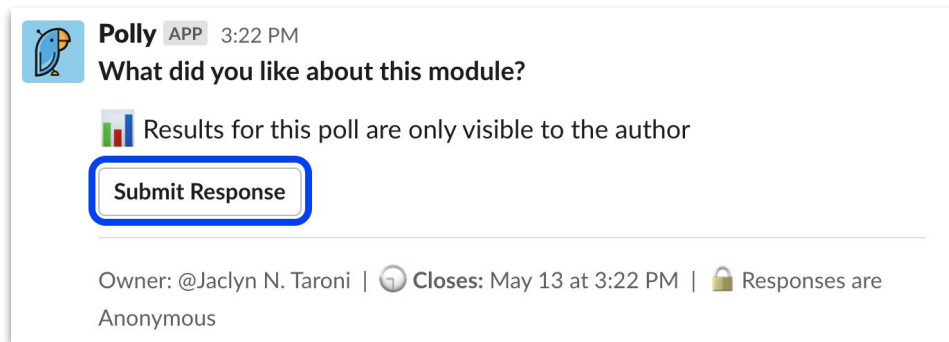


- I have questions about **previous instruction or exercise notebooks**
 - **Post** in #2020-july-training
 - If you need to share your screen, we will set up a 1:1 or group Zoom call
- I would like to be paired up with other participants
 - **Post** in #2020-july-training; we can set you up in a Zoom breakout room
- I have a question that is **highly specific to my data**
 - **Direct Message** a CCDL staff member
- I'm having trouble **logging in** to RStudio Server
 - **Direct Message** a CCDL staff member

Trouble logging into Zoom and Slack? **Email** training@ccdatalab.org

We want your feedback!

At the end of each module,
we will post a few questions
in the Slack channel.



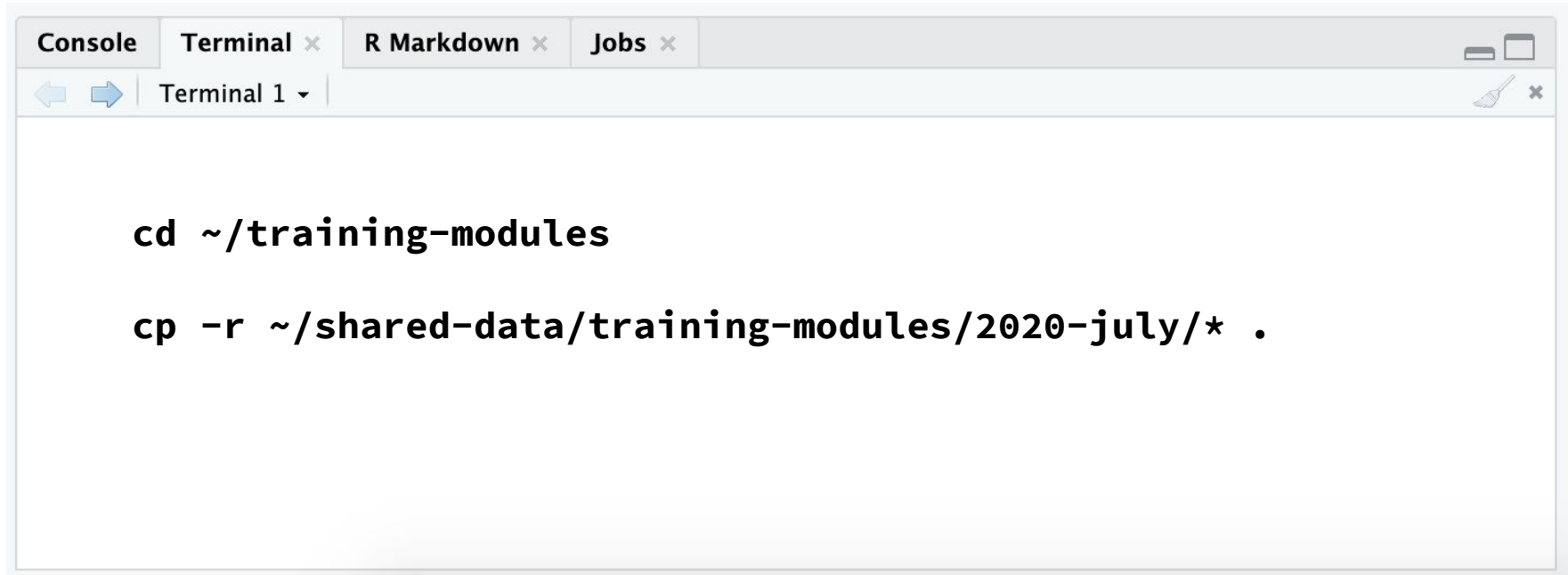
The screenshot shows a Slack poll interface. At the top, there is a header bar with a penguin icon, the name 'Polly', the label 'APP', and the time '3:22 PM'. Below this, the poll question 'What did you like about this module?' is displayed. A small bar chart icon is followed by the text 'Results for this poll are only visible to the author'. A prominent blue-outlined button labeled 'Submit Response' is centered below the text. At the bottom of the interface, a footer line contains the text 'Owner: @Jaclyn N. Taroni | ⌚ Closes: May 13 at 3:22 PM | 🔒 Responses are Anonymous'.

- The most difficult or confusing point of the module ("muddiest point")
We will post additional material answering your questions the next day
Responses to this question will appear in the channel anonymously
- What did you like about the module?
- How we can improve the module?
These responses will be collected anonymously (and not posted)

Get the modules for this workshop

Login to <http://rstudio.ccdatalab.org>

Enter the following commands in the **Terminal**:



The image shows a screenshot of the RStudio interface. At the top, there are four tabs: 'Console', 'Terminal', 'R Markdown', and 'Jobs'. The 'Terminal' tab is selected and active. Below the tabs, the terminal window is titled 'Terminal 1'. It contains two lines of text, which are commands to be entered in the terminal:

```
cd ~/training-modules  
  
cp -r ~/shared-data/training-modules/2020-july/* .
```

The terminal window has a light gray background and a white border. The commands are in a monospaced font. The first command is `cd ~/training-modules` and the second command is `cp -r ~/shared-data/training-modules/2020-july/* .`. There is a small icon of a broom and a close button in the top right corner of the terminal window.