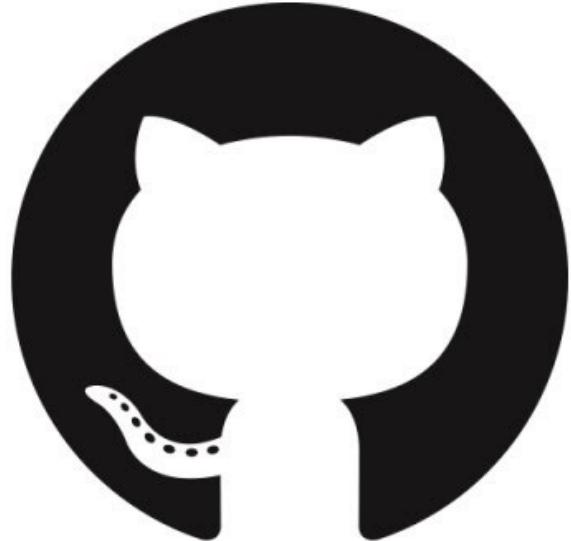


BINF2111 – Introduction to Bioinformatics

Computing

Github – git'ing ready



GitHub

**Richard Allen White III, PhD
RAW Lab**

Lecture 12 – Thursday Sep 26th, 2024

Learning Objectives

- Review quiz/bonus
- Make a github account
- Github terminal authentication
- Github functions
- Quiz 12

Install conda

<https://docs.anaconda.com/anaconda/install/>

Bioawk and conda

```
conda create -n bioawk
```

```
conda activate bioawk
```

```
conda install -c bioconda bioawk
```

```
bioawk -c fastx '{ print $name, length($seq) }' example2.fasta
```

```
bioawk -c fastx '{ print $name, gc($seq) }' example2.fasta
```

```
bioawk -c fastx '{ print ">"$name;print revcomp($seq) }'
```

```
example2.fasta
```

<https://bioinformaticsworkbook.org/Appendix/Unix/bioawk-basics.html#gsc.tab=0>

Bonus 11

- Write a function that will return the number of lines it has in it?

Bonus 11

- Write a function that will return the number of lines it has in it?

```
1 #!/bin/bash
2
3 lines_in_file () {
4     cat $1 | wc -l
5 }
6
7 num_lines=$( lines_in_file $1 )
8
9 echo Your file $1 has $num_lines lines in it.
```

Question

Write a bash script that prints any range of lines you give it from a file.

For example, bash script.sh file.tsv (prints lines 2 to 5). Place code here and make an executable script.

HOW WOULD YOU DO THIS?

Question

Write a bash script that prints any range of lines you give it from a file.

For example, bash script.sh file.tsv (prints lines 2 to 5). Place code here and make an executable script.

HOW WOULD YOU DO THIS?

```
1 #!/bin/sh
2
3 filename="$1"
4 firstline=$2
5 lastline=$3
6 linestoprint=$((lastline-firstline+1))
7
8 tail -n +$firstline "$filename" | head -n $linestoprint
```

Another way?

Question

Write a bash script that prints any range of lines you give it from a file.

For example, bash script.sh file.tsv (prints lines 2 to 5). Place code here and make an executable script.

HOW WOULD YOU DO THIS?

```
1 #!/bin/bash
2
3 range_start=$2
4 range_end=$3
5
6 cat --number $1 | head -$range_end | tail +$range_start
```

Another way?

Question

Write a bash script that prints any range of lines you give it from a file.

For example, bash script.sh file.tsv (prints lines 2 to 5). Place code here and make an executable script.

HOW WOULD YOU DO THIS?

```
1 #!/bin/bash
2
3 filename=$1
4 firstline=$2
5 lastline=$3
6
7 sed -ne "${firstline},${lastline}p;${lastline}q" < ${filename}
```

Question

Oh NO! You tried to save your file from a text editor and it got corrupted (corrupted.fastq). You must now find where the file is corrupted by looking at what is missing in the file?

Write a command or script (better) to find the corrupted lines then print their line numbers.

HOW WOULD YOU DO THIS?

Question

Oh NO! You tried to save your file from a text editor and it got corrupted (`corrupted.fastq`). You must now find where the file is corrupted by looking at what is missing in the file?

Write a command or script (better) to find the corrupted lines then print their line numbers.

HOW WOULD YOU DO THIS?

```
1 #!/bin/bash
2
3 diff $1 $2 --color
```

Another way?

Question

Oh NO! You tried to save your file from a text editor and it got corrupted (corrupted.fastq). You must now find where the file is corrupted by looking at what is missing in the file?

Write a command or script (better) to find the corrupted lines then print their line numbers.

Question

Oh NO! You tried to save your file from a text editor and it got corrupted (`corrupted.fastq`). You must now find where the file is corrupted by looking at what is missing in the file?

Write a command or script (better) to find the corrupted lines then print their line numbers.

```
1 #!/bin/bash
2
3 original=$1
4 corrupted=$2
5 original_count=`grep "@" $original | wc -l`
6 corrupted_count=`grep "@" $corrupted | wc -l`
7
8 Corrupt_check(){
9     diff $original $corrupted --color
10 }
11
12 lines_diff=$(Corrupt_check)
13
14 output=$((original_count-corrupted_count))
15
16 echo -n "Your original file has:"
17 echo $original_count
18 echo -n "Your corrupted file has:"
19 echo $corrupted_count
20 echo -n "Difference in line numbers:"
21 echo $output
22 echo -n "Difference (diff) command found these lines to be different:"
23 echo $lines_diff
```

Question

Oh NO! You tried to save your file from a text editor and it got corrupted (corrupted.fastq). You must now find where the file is corrupted by looking at what is missing in the file?

Write a command or script (better) to find the corrupted lines then print their line numbers.

OUTPUT

Question

Write a bash script or command to find sequences with >4 N's in the sequence (MultiN.fastq), have it print the header of the sequence and the sequence itself.

Question

Write a bash script or command to find sequences with >4 N's in the sequence (MultiN.fastq), have it print the header of the sequence and the sequence itself.

```
1 #!/bin/bash
2
3 input=$1
4
5 N_detector(){
6     grep -B 1 "NNNN" $input
7 }
8
9 output=$(N_detector)
10
11 echo $output
```

Question

Write a bash script or command to find sequences with >4 N's in the sequence (MultiN.fastq), have it print the header of the sequence and the sequence itself.

OUTPU

```
CGGGCCGGCTGGGACTACCTGTTCCTCGATGA|CAAGNCCCCTCCGGCNNCGCGNNGCAC|TCGGTGTGCGCGGANCGCGNNNGCNAAC|GCNCGCNCNCNCNCNNNNNNNNNNNGGAACTATCTGCGGT  
(base) docwhite@system76-pc:~/Desktop/BINF2111/data$ bash lab5_q3.sh MultiN.fastq >lab5_q3.out.txt  
(base) docwhite@system76-pc:~/Desktop/BINF2111/data$ more lab5_q3.out.txt  
@M00618:17:00000000-A31W7:1:1101:15235:1847 1:N:0:21 GTGCAGCTGCTGAAAGTCGTTCACATCTCGCCGTCGCTCGCATCGCCAGCTCGAGCGCGAGCGGCCAGGGNNNGTCGCCGCGAGGCCAGCGCGAGCTCGCCTGCAAGGCACCGCTTATCCACGTCGCCCGAGCGAGCG  
TCGCCGCATCGGAGCCGCTCTCGGTGCGGGGATCTGGCGGCCCTCGACGCCGACCTCTCGAGACTCGGCTGAGGCCGAGGGCG -- @M00618:17:00000000-A31W7:1:1101:17204:1912 1:N:0:21 CGCAGCATCAAGAACGGACCCGATCGGCCAGGGGGCCCTCGATCGCTCGC  
GCTCGANNNNNNNNGCCTCGCTCGACCTCGCCGGCGCCCGTCTCGTGCGACGGGCGCTTCGCGAGCTCGCCGCGCGTGGCGAAGCACATCGAGGTGTGGAGGCCCTCGAACC CGGGAGGGCATCGCCCGAGCCACTCGTAGAGGAAGTACACCGGACCT -- @M00618:17:0000  
00000-A31W7:1:1101:13859:2209 1:N:0:21 GGAGGGAGAGCGGGAGACGACGGGTATCGTGGAGATCGACGAGGGGAGAAAGCGCATCGAGGGCAGGAGGGNNNAAGAGGCCGAGGCGCTGCCCGAGCCCCCGCTACGGGGCTCCGGCCGCGCTCGGGCTCGGCTCGATCGACTCTCGCTTAC  
CGGCTCTCGACTCCCTACTGCTCTCGCCCTCGCTCGTCCGAGCCGACCGCTCGCTCATCGCC -- @M00618:17:00000000-A31W7:1:1101:20727:2816 1:N:0:21 GTCCCCATGATCTCGCCCGCCAGGCCCGATCGGCTCGACCTTCGCCCTCGAGC  
CNGNAACGGAGCCGGGGTTCCAGCGACGCCGGCGTTGAAGGGAGACATGCTCTCGGTGCGGATGAGGAGCTCGTCCACCGCAGGACCGTCGCTCTGAAGGCCCGCCGGAGGGCTGTTCTCAAGAATGCCGGCCCTGCC -- @M00618:17:00000000-A31W7:1:1  
101:23681:4443 1:N:0:21 GCGACGAGTAACCTCGCTCGCTGCTCGAGCTCCCTCGTGTGCGCTGAGCGTGNNNNCGTAGNNNNNNNNNNNNNNNTAGCGNNAGANNCGATCTGNNNNNNTCCANCNNNCNTGGCCGATGAGCGTCAAGACGAGGTGCAAGGGCCGTGAAGNGNTNNNCNGGCGAAGCCG  
GNGCTCCTCNCNGCCGACTCCCNNGNCCTGACGCCGGGGTAGCCCTGGCCGATG -- @M00618:17:00000000-A31W7:1:1101:24313:5033 1:N:0:21 GCCCCGACCTGCCGGCCGCTGCTGCTGAGGGGAGGGCCGAGGCTCGNNNNNNNNNNNNNNNTNTGNNT  
NNCATCGACNNNNNNNGACNNNNNNNNTCCGTGCTGAGAAGGAGGCCGACCCCTCGACGANCGGCGCGGAACCGAGCGNCGCTGGCCTGNCGATGCNNCAGATCTGAAGNGCNAANCTGA  
CTCCCCACGGCCGCTCTACG -- @M00618:17:00000000-A31W7:1:1101:26789:7963 1:N:0:21 CGCGTCGGGGGTGCGGGCCGATGAGACGATTGGGGCGNNNNNNNTGACGGGCTGACCTCTGTCGAGCAGGGGAGGTGACCGCTTCTGGGCCGAACGGGGCGGAAGTCGACGGGATCGGGATGCTGCTGCGACGGGCGGGAGGGAGCGTC  
TTCGGCATCACCGTGTGGCAGCGAATCACCGCTGCTGG -- @M00618:17:00000000-A31W7:1:1101:27801:9262 1:N:0:21 TCGGCCGCTCGCCGCGACGAGCGCCGTCGGCTGCCCTCGCAGCGCTCGGATGCCGNNNNNNNNNNNNNTGGCGACGCGGCCGACCG  
TGTAGTNNNNNTGAAACAGCGGCTCGGGCGACACAGCATGTCGACGACCATGCGAGAGCTCGCCGATCGTGTGCTGCTGAGCTCGGCTGGAGCGCCGAGCTTACGTCGAGGCCGGC -- @M00618:17:00000000-A31W7:1:1101:27772:9473 1:N:0:21 TGCTCAA  
GCTGCTGGCGCTTCGAGTTGGCGCATCGCGCTNNNNNNNNNNNNNNTACACACGCCGACGCCGACGCCGACGCCGACGCCGACGCCGACGCCGACGCCGACGCCGACGCCGACGCCG  
AGGGCGCCGCGATGATCAAACAGAGCT -- @M00618:17:00000000-A31W7:1:1101:28149:9978 1:N:0:21 GATCGCCTGGCATCGATCGGCCGAGATCCGGGAGCGTCTCATCGAGACGNNNNNNNNNNNNNNNNTNGGCNCNGCGANCNNNNNNNGC  
CGGGCGCTGGACTACCTTCTCGATGATCAAGNCTTCCGGCNNGCGNNGACTNCGTGCTGCCGAGCGCGNNNGCNANGCTGNCNCNCNCNNNNNNNNNGACTATCGCGT
```

Question

Write a bash script that counts the number of ATG (starts), Serine (S), Arginine (R), and TAA, TAG, TGA (stops) from the example2.fasta file then converts them into amino acid M, S, R, and * for TAA, TAG, TGA for stop codon.

Remember that ATG encodes for methionine so they only count as start from the beginning of the sequence or the end for the stops. Serine and Arginine can be throughout the sequence.

HOW WOULD YOU DO THIS?

Question

```
1#!/bin/bash
2
3 input=$1
4
5 #Count start codons
6 count_starts(){
7    grep "ATG" $input | wc -l
8 }
9
10 number_starts=$(count_starts)
11
12 #Count stop codons
13 count_stops(){
14    egrep "TAA$|TAG$|TGA$" $input | wc -l
15 }
16
17 number_stops=$(count_stops)
18
19 #Count codons for arginine
20 count_arg(){
21    egrep "CGT|CGC|CGA|CGG|AGA|AGG" $input | wc -l
22 }
23
24 number_arg=$(count_arg)
25
26 #Count codons for serine
27 count_ser(){
28    egrep "AGT|AGC|TCT|TCC|TCA|TCG" $input | wc -l
29 }
30
31 number_ser=$(count_ser)
32
33 #final table
34 echo -n "Number of starts: "
35 echo $number_starts
36 echo -n "Number of stops: "
37 echo $number_stops
38 echo -n "Number of Arginine: "
39 echo $number_arg
40 echo -n "Number of Serine: "
41 echo $number_ser
42
43 #Convert codons
44 convert_cod(){
45    sed 's/^ATG/M/g' $input |
46    sed 's/TAA$/*/g' |
47    sed 's/TAG$/*/g' |
48    sed 's/TGA$/*/g' |
49    sed 's/CGT/S/g' |
50    sed 's/CGC/S/g' |
51    sed 's/CGA/S/g' |
52    sed 's/CGG/S/g' |
53    sed 's/AGA/S/g' |
54    sed 's/AGG/S/g' |
55    sed 's/CGT/R/g' |
56    sed 's/CGC/R/g' |
57    sed 's/CGA/R/g' |
58    sed 's/CGG/R/g' |
59    sed 's/AGA/R/g' |
60    sed 's/AGG/R/g'
61 }
62 convert=$(convert_cod)
63
64 echo "Converting codons to listed amino acids: "
65 echo $convert
```

Question

Write a bash script that counts the number of ATG (starts), Serine (S), Arginine (R), and TAA, TAG, TGA (stops) from the example2.fasta file then converts them into amino acid M, S, R, and * for TAA, TAG, TGA for stop codon.

Remember that ATG encodes for methionine so they only count as start from the beginning of the sequence or the end for the stops. Serine and Arginine can be throughout the sequence.

OUTPU

T

```
ATGCTAAGGCATCNGACAACTGACTAAATAG
(base) docwhite@system76-pc:~/Desktop/BINF2111/data$ bash lab5_q4.sh example2.fasta
Number of starts: 7
Number of stops: 7
Number of Arginine: 7
Number of Serine: 5
Converting codons to listed amino acids:
>chr1_geneA MCTASCTATCTTGACAACTGACTGCC* >chr1_geneB MCTASCTATGTTGGCAACTGACTCCC* >chr1_geneC MCTASCTACCTTGACAACTGACTGGG* >chr1_geneD MAAASCTATCTTGACAACTGACTCCC* >chr1_geneE MCTASCTATCTGATTCTGACTTT* >chr1_geneF MGGGGGGTATCTTGACAACTGACTGCG* >chr1_geneG MCTASCTATCNGACAACTGACTAAA*
```

Question

Write a bash script that prints username, current directory, the location of our desktop directory, and the date/time

Question

Write a bash script that prints username, current directory, the location of our desktop directory, and the date/time

```
1 #!/bin/bash
2
3 echo -n "My user name is: "
4 whoami
5 echo -n "My current directory is: "
6 pwd
7 echo -n "My Desktop directory is: "
8 echo $HOME/Desktop
9 echo -n "The date and time is: "
10 date
```

Question

Write a bash script that prints username, current directory, the location of our desktop directory, and the date/time

OUTPUT

```
(base) docwhite@system76-pc:~/Desktop/BINF2111/data$ bash lab5_q5.sh
My user name is: docwhite
My current directory is: /home/docwhite/Desktop/BINF2111/data
My Desktop directory is: /home/docwhite/Desktop
The date and time is: Mon Oct  4 22:18:06 EDT 2021
```

Github

- What is github?

Github

- What is github?

GitHub, Inc. is a provider of Internet hosting for software development and version control using Git. It offers the distributed version control and source code management functionality of Git, plus its own features.

Founded 2008 (56 million, owned by microsoft)

Github

- Alternatives to github?
 - gitlab (<https://about.gitlab.com/>)
 - bitbucket (<https://bitbucket.org/>)

Github

- Uses of git:
 - Code collaboration
 - Version control
 - Safe storage of code

Github – Create an Account

github.com

- Click “Sign Up” at the top-right
- Use your UNCC email
- 2-5 collaborators is fine
- For Features I suggest “collaborative coding”
- Free account is OK, It can also be upgraded for free:
 - https://education.github.com/discount_requests/student_application
 - Select: “University of North Carolina at Charlotte”
 - Make sure to use your @uncc.edu address for this

Github – Example Markdown (Readme)

```
# Metaome Stats: Calculating denovo assembly statistics from metaomes
```

```
=====
```

```
## Installing
```

```
- pip installation <br /> (LINE BREAKS)
```

```
`pip install MetaomeStats`
```

```
- conda installation <br /> (LINE BREAKS)
```

```
`conda install -c bioconda MetaomeStats`
```

```
- source (github)
```

```
```bash
```

```
git clone https://github.com/raw-lab/metaome_stats
```

```
cd metaome_stats
```

```
python setup.py
```

```
```
```

```
=====
```

```
## Usage examples
```

```
```bash
```

# Github – terminal authentication

- Github no longer allows password authentication. There are two options to authenticate
  - ssh authentication (more involved setup)
  - Personal Access Token (Easier setup)
  - <https://docs.github.com/en/authentication/keeping-your-account-and-data-secure/creating-a-personal-access-token>
  - When git asks you for a password, use the created token instead of your password

The command on the terminal “git config --global credential.helper store” will tell git to store your password, it saves it as plain text (no encryption) so it is not secure. This is the reason for the Personal Access Token, giving limited access to anyone that may steal this token. SSH is much better, we may set that up on a future date.

# Github – Basic Functions

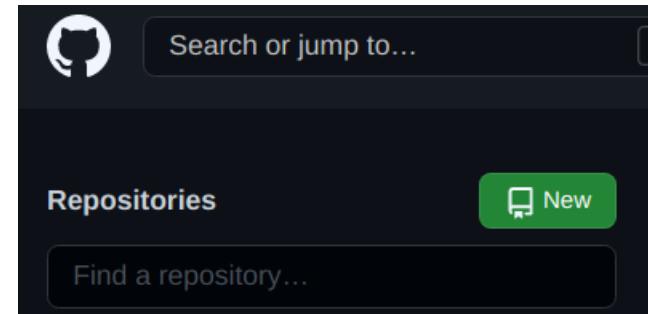
- Pulling data from github
  - git clone <https://github.com/....>
  - git pull (for updating repo after someone else makes changes)
- Git needs to authenticate before you can push, or pull from a private repo.
- Push to github (assuming you are in root directory of the code)
  - git add .
  - git commit –m “Description of changes”
  - git push

# Github – Basic Functions

- Show changes
  - git status
  - git fetch (checks server for updates without ‘pulling’ them)
- Create a new branch (for working on fixes/features)
  - git checkout –b <branch name> (i.e. dev)
- Switch between branches
  - git checkout <name>
- List local and remote branches
  - git branch -a

# Github – Practice: Hello World

- Create a new repository on github
  - Repository name: hello\_world
  - Select public and add the README file
- Now to clone it in your home directory
  - On the terminal type ‘cd’ to make sure you are in your home
  - git clone [https://github.com/<username>hello\\_world](https://github.com/<username>hello_world)
  - Type ‘ls’ to confirm the new folder “hello\_world” is there
- Congratulations on creating your first github repo!



# Github – Practice: Making Changes

- Change the README file:
  - cd hello\_world
  - ls
- You should see README.md in the directory. With your favorite text editor change the readme file (vim, nano, gedit, echo...)
  - Add the line ‘## created by <your name>” and save the file
- Push your changes: (type ‘git status’ after each command)
  - git add README.md
  - git commit –m “Added author to README”
  - git push
- Now on your browser at [https://github.com/<username>/hello\\_world](https://github.com/<username>/hello_world) you will see your updated README file

# Quiz 12

- On canvas now

# Bonus 12

- Write a bash script that converts this into a tsv file (specific files and all files)?

With this file:

Rat,steven,bear,Xi  
Olf,thor,flower,Ton  
Lazy,aaron,larry,jose  
old,thor,flower,Ton