

**Mansuri Naushin Parveen Sagir Ahmed (24310041)**

**24310041**

## **Lab Assignment - 2**

1. Open a new file called notes.txt in vi.

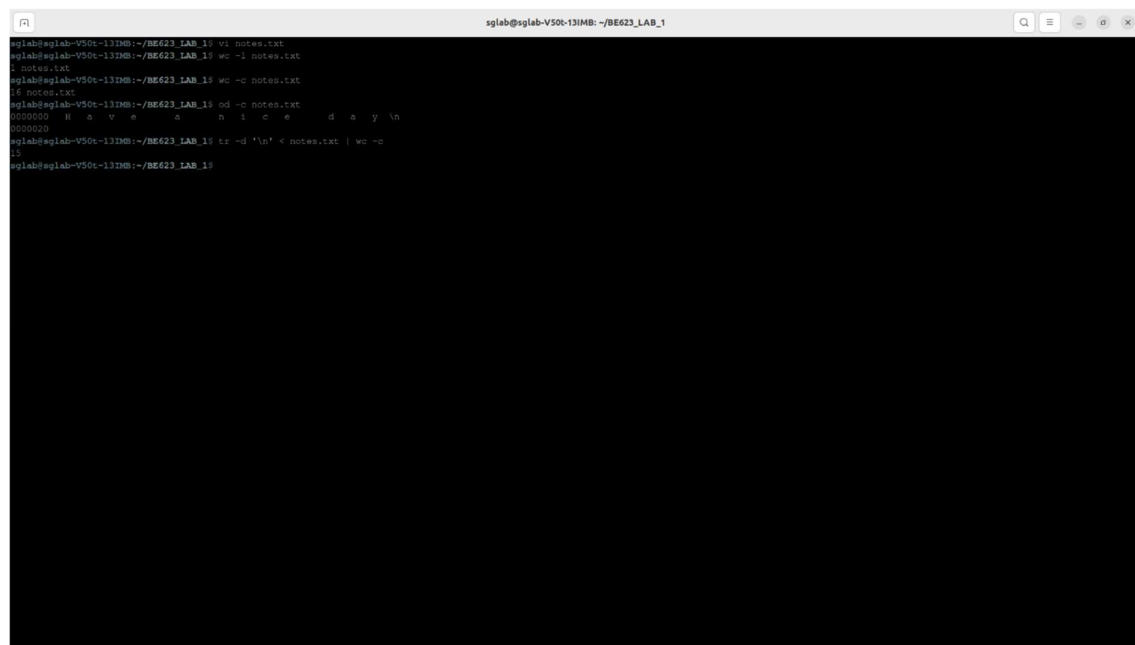
- Insert exactly one line of text:

Have a nice day

(Make sure there is no trailing space at the end.)

- Save and exit.

Verify that the file contains exactly one line and 15 characters.



```
sglab@sglab-V50t-131MB: ~/BE623_LAB_1
sglab@sglab-V50t-131MB:~/BE623_LAB_1$ vi notes.txt
sglab@sglab-V50t-131MB:~/BE623_LAB_1$ wc -l notes.txt
1 notes.txt
sglab@sglab-V50t-131MB:~/BE623_LAB_1$ wc -c notes.txt
16 notes.txt
sglab@sglab-V50t-131MB:~/BE623_LAB_1$ od -c notes.txt
000000  H a v e   a   n i c e   d a y  \n
000010
sglab@sglab-V50t-131MB:~/BE623_LAB_1$ tr -d '\n' < notes.txt | wc -c
15
sglab@sglab-V50t-131MB:~/BE623_LAB_1$
```

During Part 1, I faced an issue where `wc -m` showed 16 characters instead of 15. To solve this, I used two extra Linux commands (beyond the lab sheet):

`od -c notes.txt` → shows hidden characters like newline (`\n`) or carriage return (`\r`). This confirmed that the extra count was due to a newline, and then I used `tr -d '\n' < notes.txt | wc -m` to count characters without the newline (Ref: used chat GPT).

```
sglab@sglab-V50t-13IMB: ~/BE623_LAB_1
have a nice day

*notes.txt* 1L, 16B
1,15 All
```

Q.2. Display the last four lines of sequence.fasta without opening the file in an editor.

Q.3. In sequence5.fasta, print all header lines (lines starting with >).

```
sglab@sglab-V50t-13IMB: ~/BE623_LAB_1
sglab@sglab-V50t-13IMB:~/BE623_LAB_1$ vi notes.txt
sglab@sglab-V50t-13IMB:~/BE623_LAB_1$ tail -n 4 sequence.fasta
TAACTACTGATAAGTTACAAAAGTGTTCCTATCTAAAGGGCAATACAGCCCTAGACTCTCCAGGTAT
TTGACTCTGTCAGCAAAAAGGAAATTGAGGAAATAGAGCAAGCTATTTCAGAGGCAACTATATCACA
TAGACACCCCG

sglab@sglab-V50t-13IMB:~/BE623_LAB_1$ grep "^>" sequence5.fasta
>shr
>clock
>hif1a
>hif2a
>hif3a
>npsa1
>npsa2
>npsa3
>npsa4
>sim1
>sim2
>arnt1
>bm11
sglab@sglab-V50t-13IMB:~/BE623_LAB_1$ ls
1 notes.txt protein.fasta sequence1.fasta sequence2.fasta sequence3.fasta sequence4.fasta sequence5.fasta sequence.fasta
sglab@sglab-V50t-13IMB:~/BE623_LAB_1$
```

Q.4. Find all matches in sequence5.fasta where A is followed by any single character and then G.

```
sglab@sglab-V50t-131MB: ~/BE623_LAB_1
sglab@sglab-V50t-131MB:~/BE623_LAB_1$ grep "A.G" sequence5.fasta
HFRTHKKLOFTFIGCDAKGRIVLVGYTEAELCTRGSGVQFIHAADMLYCAESHIMIKTGESGMIVFRLLT
DAARSRRSQETVLYQLAHTLPARGVSAHLDKASIMRLTISYLRMHRLCAAGEWQVQAGGEPLDACYL
HALEGFVWVLTADGGDMAYLSENVSKHLGLSQLELIGHSIFDFHPCDQELQDALTPPTERCFSLRMKST
REKSRHAARSRRGKENLEFFELAKLLPLFGAISSQLDKASIVRLSVTYLRLRAFALGAPFWGLRAAGP
AGLAGRRGPAALVSEFEQHLGGHILQSLDGFVFAINQEGKFLYISETVSIYLGLSQVENTGSSVFDYI
HPGDHSEVLEQLGLVQERSFFVVRKSTILTRGLHVLASGYKVIHVIGRLRALGLVALGHTLFPAPLAELE
WLCAGGFWACQVATVAGSRSPGEHVLVWVHVLAGGQFI
DASKARDDQIMAEIRALRELPAEADKVLVYLNIMSLACIITRKGVFFAGGTPAGPTGLLSAGELED
IYAALPGFLVITAEGNLLYLSSEVSEHLGSMVDVAGGGSIXDIDFADHLTVRQOLTLTRLTFRCF
EKSKNAARTREKENSEFYELAKLLPLPSAITSQLDKASIIRLITSYLMRVVFFEGLOAGWHSSRTSP
EERSFFLRMKCVLAKRNAGLTCGGYKVIHCSGYLKIRNVGLVAVGHSLLPSAVTEIKLHNNMFMRASL
EKSKNAARTREKENSEFYELAKLLPLPSAITSQLDKASIIRLITSYLMRAVFFEGLOAGQPSRAGP
EERSFFLRMKCVLAKRNAGLTCGGYKVIHCSGYLKIRIVGLVAVGQSLPSAITEIKLYNNMFMRASL
ELKHLILEAADGFLFIVSCETGRVVVSDSVTVLNPQSEWFGSTLYDQVHPDDVDKLRQLSTSRMCM
GSRASFICRMCGSGSEPHFVVVHCTGYKAKFCVLVAGRLQVTSNPCTDMGNVCOPTFISRHNIEGIF
DELKHLILEAADGFLFVVGCDRGKTLFVSESVFKILNYSQNDLIGQSLFEDYLHPKDIAKVRELSSRLC
SGARRSFFCRMKCNRPFRSFCITISTGYLNSLSCVAGRLSHSVVFPVNGEIRVKSMEYVSRHAIDG
sglab@sglab-V50t-131MB:~/BE623_LAB_1$
```

Q.5. Find all matches in sequence5.fasta where P is followed by any character except A, then L.

```
sglab@sglab-V50t-131MB: ~/BE623_LAB_1
sglab@sglab-V50t-131MB:~/BE623_LAB_1$ grep "P[^A]L" sequence5.fasta
QLHWQIPFENSPLMERCFICRLCLDNSSGFLAMNFQKLVLPPQLALFAIATPLQPPSILEIRTKNF
WRMKCTVTNRGRTVNLKSAITWKVLHCTGQKVYVEPLLGCLIMCEPIQHPSHMDIFLDSKTFLSRHSNDM
LTSRGRTLNKAAATWKVLNCSGHMRAYEPPLCCLVLICEAIPHGSLPPLGRGAFLSRHSMDMKFTYCD
FTQMLLEALDGFIIAVTIDGSIIVVSDSITPLLGHLPSDVMDQNLNLFPEQHSSEVYKILSEYLSKSDS
ELKHLILEAADGFLFIVSCETGRVVVSDSVTVLNPQSEWFGSTLYDQVHPDDVDKLRQLSTSRMCM
sglab@sglab-V50t-131MB:~/BE623_LAB_1$
```

6. Print all lines in sequence5.fasta that have exactly 2 consecutive Vs anywhere in the line.

```
sglab@sglab-V50t-131MB: ~/BE623_LAB_1
sglab@sglab-V50t-131MB:~/BE623_LAB_1$ grep "VV" sequence5.fasta
AANFRGELNLSGEFLGALNGFVLVWTTDAIVFASTIGQYLFGQSDVHQSVYELIHTEBRAEFQR
IMLQTHYIITYHQNNRSEPFIVCTHIVVGYAEVRAE
TVIYNTKNSQFCQICVCHVTVVSGIIOHDL
CMDNLYLKALEGFIAVWTDGDMIFLSENISKFMGLTQVELTGHISIFDTHPCDHEEIRENLSTERDFF
RFTYCDRITELIGYHPEELLGRSAYEFYHALDSENMTKSHQNLCTKGQVSGQYRMIAKHGGYVWLETQ
DRIAEVAGYSDDLIGCSAYEYIHALDSDAVSEKSIHTLLSKQAVTGGYRFLARSGGYLWTQQTAVVSG
QTHYIITYHQNNRSEPFIVCTHIVVGYADVRE
GVVHFGDGVEMASQGLMTERSFIRMKSTLTKRGVHIKSSGKYVHITGRLLRLMGLVVAHALPFFTIT
ISESVLIYLGFERSELLCKSWYGLLHPEDLAHASAQHYRLLAESGDIOAEMVVALQARTGGWAMIYCLLY
EKSKNAARTREKENSEFEYELAKLLPLPSAITSQLOKASIIRLTTSYLMKRVVFEGLGEAWGHSRTSP
LDNVRELGLSHLQTLDOFIVWADPKIMYISETASVHLGLSQVELTGNIIYEYIHPADHENTAVLTA
LDGVANELGSHLQTLDOFIVWADPKIMYISETASVHLGLSQVELTGNIIYEYIHPADHENTAVLTA
RYATVWNSRSRSPHCIVSNVYVLEIYEYEL
ELKHLILEAADGFLFVSCETGRVWVSDSVTPVLNQPQSENFGSTLYDQVHPDDVDKLEQISTSRMCM
RSRRSFICRMACGSEPEHVWVHCTGYIKAKCLVAIGRLQVTSFNPCTDMSNVQPTFESISHNIEGIF
TFVDHRCVATVGVQPELLGKNIVEFCHPEDQLLRDSFOQVVLKQGVLSVMFRFRSKNQEWLMMRTSS
DELKHLILEAADGFLFVWCDRGRKILFVSESVFKILNYSQNDLIGQSLFDYLPKDIKAKVKEQLSSSRLC
SGARRSFFCRMKNRPRKSFCTIHTSTGYLKNLSCLVAIGRLHSHVVPQPVNGEIRKVSMEVYSRHAIDG
RWFSEFNNFWTKEVEYIVSTNTVVL
sglab@sglab-V50t-131MB:~/BE623_LAB_1$
```

7. Print all lines in sequence5.fasta that contain either AA or DD.

```
sglab@sglab-V50t-131MB: ~/BE623_LAB_1
sglab@sglab-V50t-131MB:~/BE623_LAB_1$ grep -E "AA|DD" sequence5.fasta
AANFRGELNLSGEFLGALNGFVLVWTTDAIVFASTIGQYLFGQSDVHQSVYELIHTEBRAEFQR
RFTYCDRITELIGYHPEELLGRSAYEFYHALDSENMTKSHQNLCTKGQVSGQYRMIAKHGGYVWLETQ
DRIAEVAGYSDDLIGCSAYEYIHALDSDAVSEKSIHTLLSKQAVTGGYRFLARSGGYLWTQQTAVVSG
QTHYIITYHQNNRSEPFIVCTHIVVGYADVRE
GVVHFGDGVEMASQGLMTERSFIRMKSTLTKRGVHIKSSGKYVHITGRLLRLMGLVVAHALPFFTIT
ISESVLIYLGFERSELLCKSWYGLLHPEDLAHASAQHYRLLAESGDIOAEMVVALQARTGGWAMIYCLLY
EKSKNAARTREKENSEFEYELAKLLPLPSAITSQLOKASIIRLTTSYLMKRVVFEGLGEAWGHSRTSP
LDNVRELGLSHLQTLDOFIVWADPKIMYISETASVHLGLSQVELTGNIIYEYIHPADHENTAVLTA
LDGVANELGSHLQTLDOFIVWADPKIMYISETASVHLGLSQVELTGNIIYEYIHPADHENTAVLTA
RYATVWNSRSRSPHCIVSNVYVLEIYEYEL
ELKHLILEAADGFLFVSCETGRVWVSDSVTPVLNQPQSENFGSTLYDQVHPDDVDKLEQISTSRMCM
RSRRSFICRMACGSEPEHVWVHCTGYIKAKCLVAIGRLQVTSFNPCTDMSNVQPTFESISHNIEGIF
TFVDHRCVATVGVQPELLGKNIVEFCHPEDQLLRDSFOQVVLKQGVLSVMFRFRSKNQEWLMMRTSS
DELKHLILEAADGFLFVWCDRGRKILFVSESVFKILNYSQNDLIGQSLFDYLPKDIKAKVKEQLSSSRLC
SGARRSFFCRMKNRPRKSFCTIHTSTGYLKNLSCLVAIGRLHSHVVPQPVNGEIRKVSMEVYSRHAIDG
RWFSEFNNFWTKEVEYIVSTNTVVL
sglab@sglab-V50t-131MB:~/BE623_LAB_1$
```

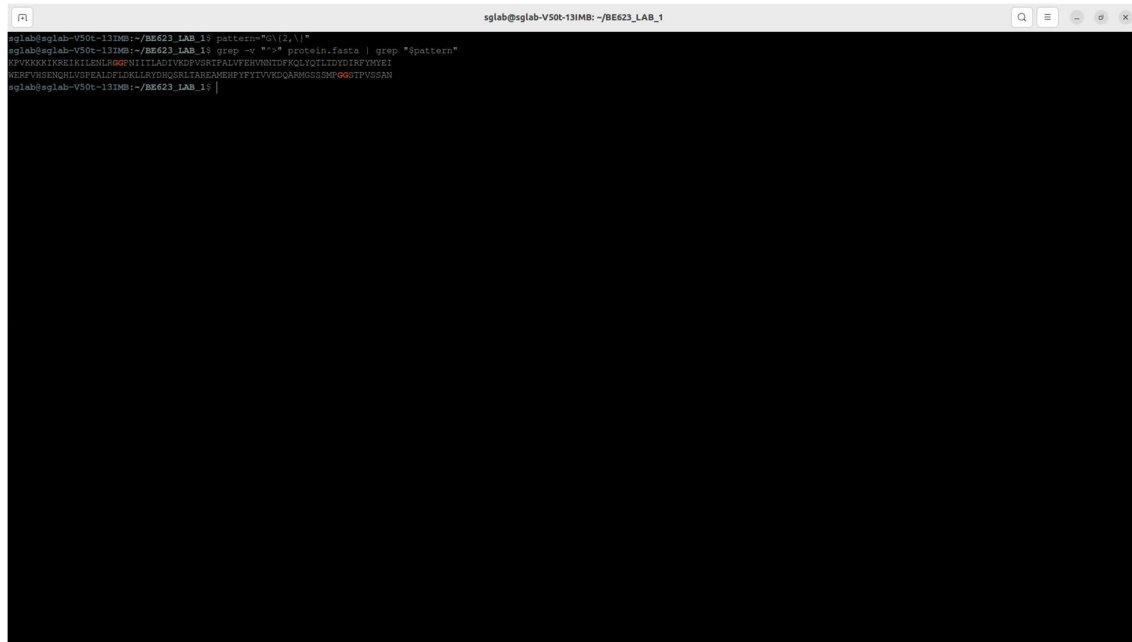
8. Print only the sequence lines (ignore headers) from sequence5.fasta that contain the letter P.

[illegible]

9. Store the filename sequence5.fasta in a variable called seq and print the number of sequences in it (headers count as sequences).

```
sglab@sglab-V50t-131MB: ~/BE623_LAB_1
sglab@sglab-V50t-131MB:~/BE623_LAB_1$ seq="sequence5.fasta"
sglab@sglab-V50t-131MB:~/BE623_LAB_1$ echo "Number of sequences in $seq:"
Number of sequences in sequence5.fasta:
sglab@sglab-V50t-131MB:~/BE623_LAB_1$ grep -c ">" $seq
13
sglab@sglab-V50t-131MB:~/BE623_LAB_1$
```

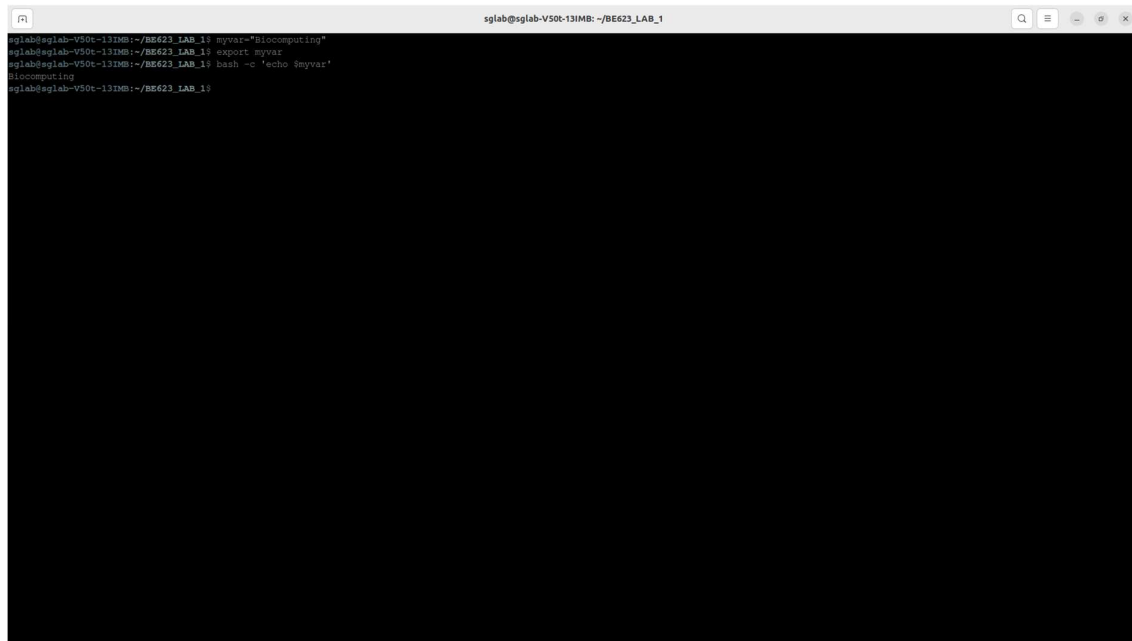
10. Store the pattern `G{2,}` in a variable and search `protein.fasta` for sequence lines (ignore headers) with 2 or more consecutive Gs.



```
sglab@sglab-V50t-131MB: ~/BE623_LAB_1
sglab@sglab-V50t-131MB:~/BE623_LAB_1$ pattern="G{2,}"
sglab@sglab-V50t-131MB:~/BE623_LAB_1$ grep -v ">" protein.fasta | grep "$pattern"
VFVRSKIRKRIHLELQGNIIITADLVKDPGRTPALAFERDRHTSPQLVITLIDYDFPHVEI
REFVVEENQHLVSEFALDFLKLILAYDHOSRLTAREAHETFTITVNDQARMSSQNGGSTIVSSN
sglab@sglab-V50t-131MB:~/BE623_LAB_1$
```

11. Store "Biocomputing" in a variable, export it, and verify that it is available inside a new shell started using:

`bash -c 'echo $VARIABLE_NAME'`



```
sglab@sglab-V50t-131MB:~/BE623_LAB_1$ myvar="Biocomputing"
sglab@sglab-V50t-131MB:~/BE623_LAB_1$ export myvar
sglab@sglab-V50t-131MB:~/BE623_LAB_1$ bash -c 'echo $myvar'
Biocomputing
sglab@sglab-V50t-131MB:~/BE623_LAB_1$
```

12. Write a shell script that checks if sequence3.fasta exists in the current folder. If yes, print the number of lines. If no, print "Missing file".



The image contains two terminal window screenshots. The top window shows a user at a prompt in a directory named ~/BE623\_LAB\_1. They run 'vi check\_file.sh' to create a script. The script's content is shown in the second window: it uses 'wc -l' to count lines if the file 'sequence3.fasta' exists, and 'echo' to print 'Missing file' otherwise. The user then runs 'bash check\_file.sh', which outputs '19 sequence3.fasta'.

```
sglab@sglab-V50t-13IMB: ~/BE623_LAB_1
sglab@sglab-V50t-13IMB:~/BE623_LAB_1$ vi check_file.sh
sglab@sglab-V50t-13IMB:~/BE623_LAB_1$ less vi.sh
vi.sh: No such file or directory
sglab@sglab-V50t-13IMB:~/BE623_LAB_1$ lesscat check_file.sh
lesscat: command not found
sglab@sglab-V50t-13IMB:~/BE623_LAB_1$ less check_file.sh
sglab@sglab-V50t-13IMB:~/BE623_LAB_1$
sglab@sglab-V50t-13IMB:~/BE623_LAB_1$ bash check_file.sh
19 sequence3.fasta
sglab@sglab-V50t-13IMB:~/BE623_LAB_1$
```

```
#!/bin/bash

if [ -f "sequence3.fasta" ]; then
    wc -l sequence3.fasta
else
    echo "Missing file"
fi

check_file.sh (END)
```

13. Using a for loop, go through all .fasta files in the current directory and print: filename, number of sequences, and file size in characters.

```
sglab@sglab-V50t-131MB: ~/BE623_LAB_1
sglab@sglab-V50t-131MB:~/BE623_LAB_1$ vi loop_file.sh
sglab@sglab-V50t-131MB:~/BE623_LAB_1$ bash loop_file.sh
sequence1.fasta 1 sequences 467 characters
sequence2.fasta 1 sequences 974 characters
sequence3.fasta 4 sequences 1710 characters
sequence4.fasta 2 sequences 1000 characters
sequence5.fasta 4 sequences 2374 characters
sequence6.fasta 13 sequences 4229 characters
sequence7.fasta 1 sequences 7951 characters
sglab@sglab-V50t-131MB:~/BE623_LAB_1$
```

```
#!/bin/bash
for file in *.fasta; do
    count=$(grep -c '^>' "$file")
    size=$(wc -c < "$file")
    echo "$file $count sequences $size characters"
done

loop_file.sh* 6L, 143B
```



14. Modify the above loop so that it only prints files with more than 3 sequences.

```
sglab@sglab-V50t-131MB: ~/BE623_LAB_1
sglab@sglab-V50t-131MB:~/BE623_LAB_1$ vi loop_modified.sh
sglab@sglab-V50t-131MB:~/BE623_LAB_1$ bash loop_modified
bash: loop_modified: No such file or directory
sglab@sglab-V50t-131MB:~/BE623_LAB_1$ bash loop_modified.sh
sequence1.fasta: 4 sequences 1710 characters
sequence1.fasta: 4 sequences 2374 characters
sequence3.fasta: 13 sequences 4229 characters
sglab@sglab-V50t-131MB:~/BE623_LAB_1$
```

```
sglab@sglab-V50t-131MB: ~/BE623_LAB_1
#!/bin/bash
for file in *.fasta; do
    count=$(grep -c '^>' "$file")
    if [ "$count" -gt 3 ]; then
        size=$(wc -c < "$file")
        echo "$file $count sequences $size characters"
    fi
done

*loop_modified.sh* 9L 101B
6/2-26 311
```

15. From `sequence5.fasta`, extract only the sequence lines (no headers) that contain 3 or more cysteines (C). Save the output to a file named `cys_rich.txt`. Ensure the output file contains no empty lines.

[illegible]