

# Text Processing

Session 5 (And last)

# Table of contents

**01**

Introduction to Text  
Processing

**02**

Basic Text Manipulation  
commands

**03**

Searching with “grep”  
command

**04**

Regex and Pattern  
Matching



01

# **Introduction to Text Processing**

# What is Text processing?

- Text processing is the automated manipulation of text data using tools or scripts, usually to extract, transform, or analyze information.
- It makes managing big files like logs or configs a lot easier, saves time, and avoids mistakes

# Why Text Processing is Crucial ?

## Efficiency

Efficient text processing helps in extracting valuable information quickly.

## Data Analysis

Processing logs and configuration files aids in monitoring troubleshooting, and performance tuning

## Manipulation

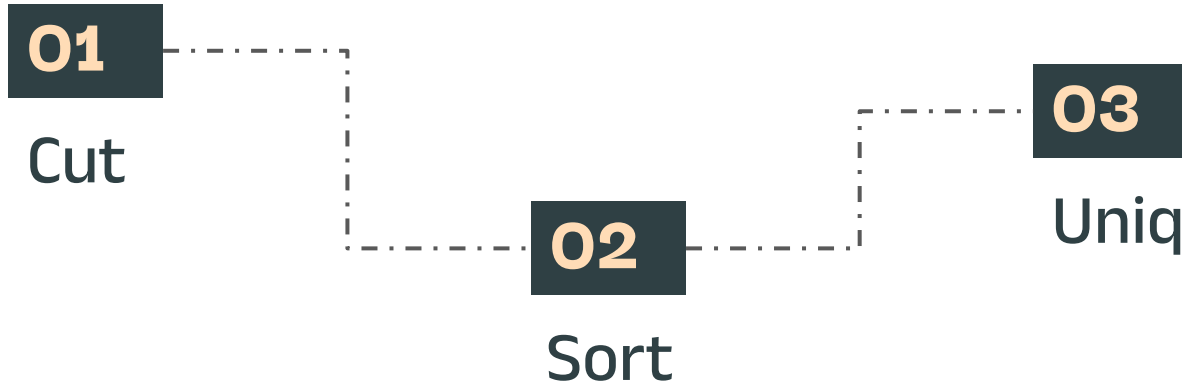
Modifying, rearranging, or cleaning text to prepare it for further use or analysis



02

# **Basic text manipulation commands**

**There are many commands that help in text manipulation**  
**But we'll only mention some of them**





**“cut”  
command**



# “cut” command

- A powerful text processing utility that extracts specific sections from each line of files.
- Perfect for working with delimited text files, log files, and formatted output from other commands.
- Basic syntax :

```
cut [option] [file] (options are necessary)
```

# Options

<b>-f</b>	To extract fields/columns
<b>-d</b>	To specify a specific delimiter other than the default (Tab)

```
sheikhwalter ~/OSC/Text-processing
>> cut -f 1,3 grades_tab.txt
Name    Physics
Alice   92
Bob     84
Charlie 91
David   95
Eve     55
```

```
sheikhwalter ~/OSC/Text-processing
>> cut -d ',' -f 1,2,4 grades.csv
Name,Math,Chemistry
Alice,88,86
Bob,81,76
Charlie,61,75
Diana,67,100
Ethan,75,87
```

# Options

<code>-f <i>n</i>, <i>k</i></code>	Extract the <i>n</i> th and the <i>k</i> th fields
<code>-f -<i>n</i></code>	Extract from the start to the <i>n</i> th field
<code>-f <i>n</i>-</code>	Extract from the <i>n</i> th field to the end
<code>-f <i>n</i>-<i>k</i></code>	Extract from the <i>n</i> th field to the <i>k</i> th field

Works with -b and -c in the same way

# Options

<b>-b</b>	To extract specific bytes
<b>-c</b>	To extract specific characters

```
sheikhwalter ~/OSC/Text-processing  
>> cut -b 1-4 baby  
waaa  
gogo  
gaga  
ab32  
ab5d
```

```
sheikhwalter ~/OSC/Text-processing  
>> cut -c 1-4 baby  
waaa  
gogo  
gaga  
ab32  
ab5d
```

# Options

-b	To extract specific bytes
-c	To extract specific characters

```
while ((optc = getopt_long (argc, argv, "b:c:d",
    != -1)
{
    switch (optc)
    {
        case 'b':
        case 'c':
            /* Build the byte list. */
            byte_mode = true;
            FALLTHROUGH;
        case 'f':
            /* Build the field list. */
            if (spec_list_string)
```

Actually, there's no  
difference between them  
(for now at least)

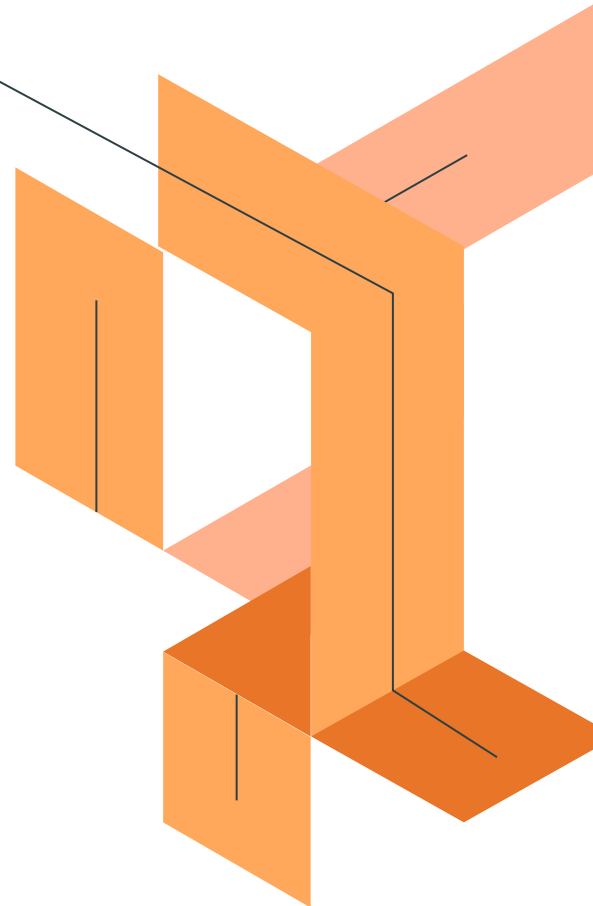


**“sort”  
command**

# “sort” command

- The **sort** command is a powerful text processing utility that arranges lines of text files in alphabetical or numerical order.
- It supports various options for customizing sort behavior, handling different data types, and processing complex datasets.
- **Basic syntax :**

```
sort [option] [file]
```



By default, the **sort** command:

- Sorts lines alphabetically
- Is case-sensitive (uppercase before lowercase)
- Uses the entire line for sorting

```
$ sort fruits.txt
```

```
Banana
```

```
Grape
```

```
apple
```

```
kivi
```

```
orange
```



# Options

<b>-n</b>	Numeric sort (treats multi-digit numbers as numbers not strings)
<b>-h</b>	Human-readable numeric sort (2K, 1M)
<b>-r</b>	Reverse order
<b>-f</b>	Case-insensitive sort
<b>-u</b>	Remove duplicates
<b>-k <i>n</i></b>	Sort by the <i>n</i> th column
<b>-t</b>	Specify field separator (default is Tab)
<b>-c</b>	Check if file is sorted

# Examples

## Case-Insensitive Sort

```
$ sort fruits.txt
```

Banana

Grape

apple

kivi

orange

```
$ sort -f fruits.txt
```

apple

Banana

Grape

kivi

orange

# Examples

Sort numeric values

numbers.txt

```
$ sort -n numbers.txt
```

1

2

5

10

100

10

2

100

5

1

# Examples

Sort by a specific column

```
$ sort -k2 -n employees.txt  
Alice 28 Designer  
Eva 31 Engineer  
John 35 Developer  
Mike 39 Analyst  
Bob 42 Manager
```

```
John 35 Developer  
Alice 28 Designer  
Eva 31 Engineer
```

employees.txt



# Examples

Output sorted values without duplicates

```
$ sort -u colors.txt  
blue  
green  
orange  
red  
yellow
```

colors.txt

```
red  
blue  
green  
red  
yellow  
blue  
orange
```

# Examples

## Human-Readable Sort

```
$ sort -h sizes.txt  
2K  
15K  
1M  
100M  
3G
```

sizes.txt

```
2K  
1M  
15K  
100M  
3G
```

# Examples

## Sort with custom delimiter

```
>> sort -t ',' -k 2 grades.csv  
Charlie,61,50,75,95,88  
Diana,67,60,100,58,50  
Ethan,75,74,87,62,57  
Bob,81,68,76,88,98  
Alice,88,56,86,87,72
```

## Check if the file is sorted

```
sort -c file.txt  
# Outputs an error message if not sorted  
# Nothing is shown if sorted
```



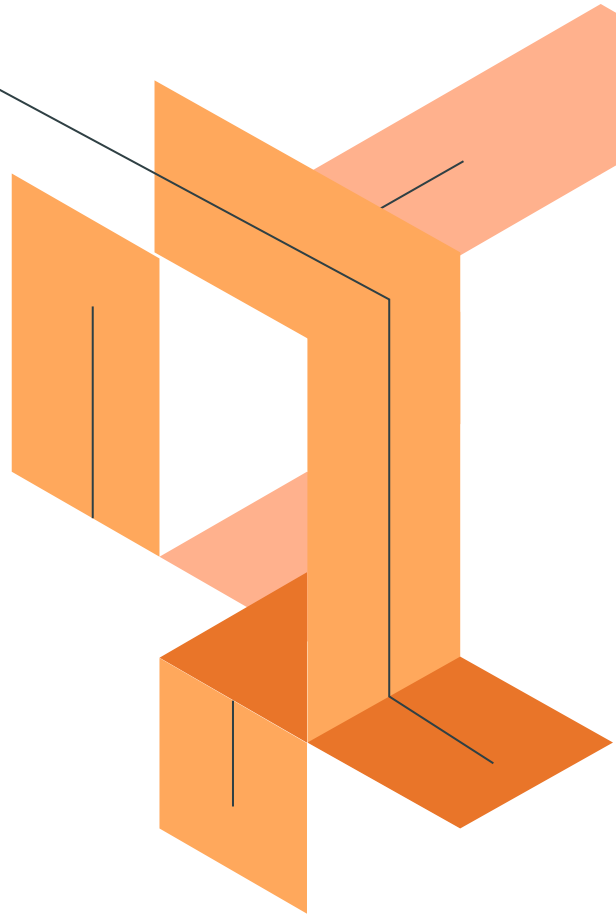
# **“uniq” command**



# “uniq” command

- It removes any duplicate lines and sends the results to standard output.
- It is often used in conjunction with sort to clean the output of duplicates.
- **Basic syntax:**  

```
uniq [option] [file]
```
- By default , output will be all lines without duplication



# Options

<b>-c</b>	Print each output line with the number of occurrences
<b>-d</b>	Display only duplicate lines
<b>-u</b>	Display only unique lines
<b>-i</b>	Ignore case sensitive

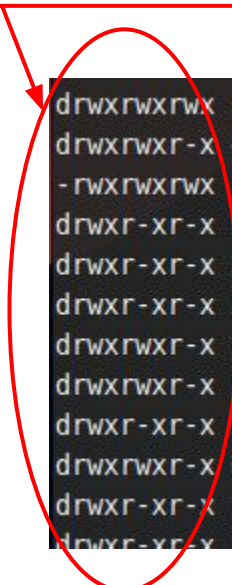


# **Hands on**

# **1**

1. Go to your home directory
2. Display all the unique file permissions only

Hint: these are the file permissions



```
drwxrwxrwx 1 root      root      16384 Aug 10 23:50 'Big Files'
drwxrwxr-x 4 sheikwalter sheikwalter 4096 Mar 22 22:59 Coding
-rwxrwxrwx 1 sheikwalter sheikwalter 70 Jul 15 22:24 Deltarune.sh
drwxr-xr-x 2 sheikwalter sheikwalter 4096 Jul 25 04:34 Desktop
drwxr-xr-x 3 sheikwalter sheikwalter 4096 Jul 25 04:34 Documents
drwxr-xr-x 5 sheikwalter sheikwalter 4096 Aug 13 23:18 Downloads
drwxrwxr-x 2 sheikwalter sheikwalter 4096 Jul 25 03:17 fake-libasound2
drwxrwxr-x 5 sheikwalter sheikwalter 4096 Apr 23 10:29 interview
drwxr-xr-x 2 sheikwalter sheikwalter 4096 Mar 18 22:07 Music
drwxrwxr-x 6 sheikwalter sheikwalter 4096 Aug 7 22:46 OSC
drwxr-xr-x 2 sheikwalter sheikwalter 4096 Aug 13 23:56 Pictures
drwxr-xr-x 2 sheikwalter sheikwalter 4096 Feb 10 2025 Public
```

## Solution:

1. `cd ~`

2. `ls -l | cut -c -10 | sort | uniq`

Or

`ls -l | cut -b -10 | sort | uniq`

Or

`ls -l | cut -f 1 -d ' ' | sort | uniq`

```
>> ls -l | cut -c -10 | sort | uniq
drwxrwxrwx
drwxrwxr-x
drwxr-xr-x
-rw-rw-r--
-rwxrwxrwx
total 9843
```



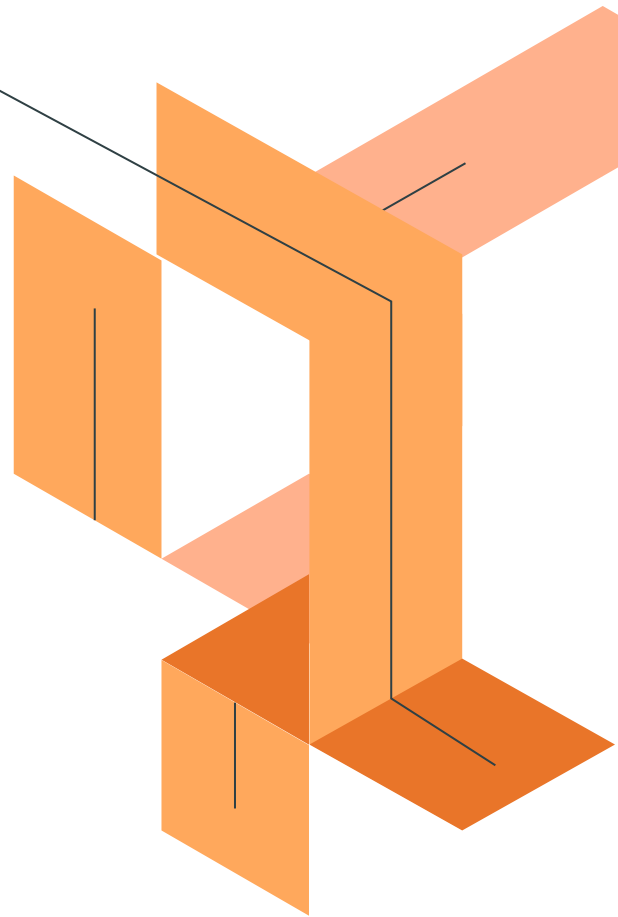
03

---

## Searching with “grep” command

# “grep” command

- “Grep” stands for ‘Global Regular expression print’.
- Used to search files for the occurrence of a string of characters that matches a specified pattern
- Command syntax:  
`grep [option] [text/pattern] [file]`



# Options

<b>-n</b>	Print the number of each line that contain the matching pattern
<b>-c</b>	Count matching lines only
<b>-o</b>	Print <b>only</b> the matched part (not the whole line)
<b>-i</b>	Ignore case sensitive
<b>-v</b>	Print all lines that doesn't match text
<b>-l</b>	Print name of each file that contain match text
<b>-L</b>	Print name of each file that <b>doesn't</b> contain match text
<b>-r</b>	Recursive search in files inside a directory (prints matches with file paths)



# Options for context control

-A <i>N</i>	Print <i>N</i> lines <b>After</b> the match.
-B <i>N</i>	Print <i>N</i> lines <b>Before</b> the match.
-C <i>N</i>	Print <i>N</i> lines <b>Before and After</b> the match.



# **Hands on**

# **2**

Find all lines in `/etc/passwd` that mention the word `bash` (case-sensitive) and display the number of each line.

## Solution:

```
grep -n bash /etc/passwd
```

```
>> grep -n bash /etc/passwd  
1:root:x:0:0:root:/root:/bin/bash  
47:sheikhwalter:x:1000:1000:SheikhWalter,,,:/home/sheikhwalter:/bin/bash
```

# Let's take a Break!



# Text can be written in 3 ways:

## 1. Text only without quotes

- Good for simple patterns without special characters or whitespaces
- Lets the shell interpret the special characters then send it to grep
- Words separated by spaces become separate arguments
- Can't safely reference variables in scripts if they contain spaces or special characters

```
grep text filename
```

# Text can be written in 3 ways:

## 2. Text between double quotes (" ")

- Same as without quotes but keeps the string as one argument, even with spaces
- Allows referencing variables in scripts
- Some backslashes and characters (like `\n`, `\\`) are interpreted by the shell
- Can cause issues when using regex

```
grep "text" filename
```

# Text can be written in 3 ways:

## 3. Text between single quotes ('')

- Handles the whole text as a literal string with no shell interpretation and sends it to grep to do its tricks
- Safest when using regex, backslashes, or special characters
- Can't use variables in scripts

```
grep 'text'  
filename
```

The slide features decorative orange geometric shapes in the corners. On the left, there is a vertical orange bar and a diagonal orange bar meeting at a right angle. On the right, there is a similar arrangement of orange bars forming a corner. The text is centered in a dark blue font.

**But wait.**

**What does Regex even mean?**





04

# Regex and Pattern Matching

# Regular Expressions

- Regular Expressions (often called **Regex** or **Regexp**) are sequences of characters used to search, find, and manipulate text using patterns.
- They are powerful tools used in Linux programs like **grep**, **bash**, **sed**, and many programming languages.

**Regex comes in a few versions,  
main ones are:**

1. **BRE**: Basic Regular Expressions
2. **ERE**: Extended Regular Expressions
3. **PCRE**: Perl-Compatible Regular Expressions

# Regular Expressions

The default version in tools like `grep` and `sed` is **BRE**, which is stricter as many characters need to be escaped.

However, `grep` gives you the option to use **ERE** which is easier to write or **PCRE** for more advanced features like lookaheads.

<code>grep -E</code>	For Extended Regular Expressions ( <b>ERE</b> )
<code>grep -P</code>	For Perl-Compatible Regular Expressions ( <b>PCRE</b> )

# Types of Special characters

Metacharacters that are special by default:	Characters that become special when escaped by <b>backslash \</b> :
<ul style="list-style-type: none"><li>• .</li><li>• *</li><li>• ^</li><li>• \$</li><li>• [ ]</li></ul>	<ul style="list-style-type: none"><li>• \ </li><li>• \+</li><li>• \?</li><li>• \{ \}</li><li>• \( \)</li></ul>

# Important note:

All characters that are special by default can be made literal by escaping them with **backslash** `\`, even the **backslash** itself.

```
sheikhwalter ~/OSC/session
>> grep 'th\\s' message
Don't screw th\s up

sheikhwalter ~/OSC/session
>> grep '2\\*2' message
you know that 2*2 is 4
```

# Usages of Special characters

## 1. Anchors:

It matches according to the position of the pattern in the line, not the characters only

^	Matches lines starting with the pattern
\$	Matches lines ending with the pattern

```
sheikhwalter ~/OSC/session
>> grep '^Khayat' message
Khayat went to the college today, Khayat is my friend!

sheikhwalter ~/OSC/session
>> grep 'end$' message
He's my friend to the end of time, I don't want our friendship to end
```

## 2. Repetitions

.	matches any single character
?	matches zero or one of the previous item
*	matches zero or more of the previous item
+	matches one or more of the previous item

```
sheikhwalter ~/OSC/session  
>> grep 'H.ts' words  
Hats Hits Huts
```

```
sheikhwalter ~/OSC/session  
>> grep 'jars\?' words  
jar jars jarssss
```

```
sheikhwalter ~/OSC/session  
>> grep 'jars*' words  
jar jars jarsssss
```

```
sheikhwalter ~/OSC/session  
>> grep 'jars\+' words  
jar jars jarsssss
```

## 2. Repetitions

<code>\{n\}</code>	matches exactly <b>n</b> times of the previous item
<code>\{n,\}</code>	matches at least <b>n</b> times of the previous item
<code>\{,m\}</code>	matches at most <b>m</b> times of the previous item
<code>\{n,m\}</code>	matches from <b>n</b> to <b>m</b> times of the previous item

Shows that the pattern was found 4 times not 3 as it seems in the first command

```
sheikhwalter ~/OSC/session
```

```
>> grep 'b\{2\}' words
```

```
bbbb
```

```
bbb
```

```
bb
```

```
sheikhwalter ~/OSC/session
```

```
>> grep -o 'b\{2\}' words
```

```
bb
```

```
bb
```

```
bb
```

```
bb
```



## 2. Repetitions

At least 2 "b"s

```
sheikhwalter ~/OSC/session  
>> grep -o 'b\{2,\}' words  
bbbb  
bbb  
bb
```

At most 2 "b"s

```
sheikhwalter ~/OSC/session  
>> grep -o 'b\{,2\}' words  
b  
bb  
bb  
bb  
b  
bb
```

More than 2 and less  
than 3 "b"s

```
sheikhwalter ~/OSC/session  
>> grep -o 'b\{2,3\}' words  
bbb  
bbb  
bb
```

### 3. Alternation \

Alternation lets you match **one pattern OR another** using `|` character, which is useful for full patterns

```
sheikhwalter ~/OSC/session
>> grep 'Apple|Banana' words
Apple
Banana
```

---

But what if we only want to match **one character out of several**, like one of the 10 first alphabet letters?

```
grep 'a|b|c|d|e|f|g|h|i|j'
filename
```

???

That would be really tedious 😞

That's where **character classes** come in

## 4. Character classes

Character classes let you match any single character from a set without writing multiple full patterns

We write the set we want between **square brackets** `[ ]`

Type	Example	Syntax
Any of a set	Match <code>a</code> , <code>b</code> , or <code>c</code>	<code>[abc]</code>
A range	Match numbers between <code>1</code> and <code>8</code>	<code>[1-8]</code>
Negation	Match anything except <code>a</code> , <code>b</code> , or <code>c</code>	<code>[^abc]</code>

## 4. Character classes

Match any line  
ending with **s** or **e**

```
sheikhwalter ~/OSC/session  
>> grep '[se]$\n' words  
jar jars jarsssss  
Apple
```

Match any line not  
ending with **s** or **e**

```
sheikhwalter ~/OSC/session  
>> grep '[^se]$\n' words  
  
bbb  
bb  
b  
watch  
Banana  
1 2 3 4 5
```

Match all numbers  
from **2** to **4**

```
sheikhwalter ~/OSC/session  
>> grep '[2-4]\n' words  
1 2 3 4 5
```

## 4. Character classes

**POSIX Character Classes:** These are special sets that work inside brackets

<code>[[:digit:]]</code>	Digits (numbers)
<code>[[:lower:]]</code>	Lowercase letters
<code>[[:upper:]]</code>	Uppercase letters
<code>[[:alpha:]]</code>	All alphabetic letters
<code>[[:alnum:]]</code>	Alphanumeric characters (letters and numbers)
<code>[[:punct:]]</code>	Punctuation characters
<code>[[:print:]]</code>	All printable Characters including <code>alnum</code> , <code>punct</code> and <code>whitespace</code> (doesn't include characters like <code>\n</code> , <code>\t</code> , <code>\b</code> , etc...)
<code>[[:space:]]</code>	All space characters including <code>whitespace</code> , <code>\n</code> , <code>\t</code> , <code>\b</code> , etc...

## 4. Character classes

### Examples:

Numbers only

```
sheikhwalter ~/OSC/session  
>> grep '[:digit:]' grocery  
Bread 20  
Oranges 10  
Mac & cheese 3
```

Alphabetic characters only

```
sheikhwalter ~/OSC/session  
>> grep '[:alpha:]' grocery  
Bread 20  
Oranges 10  
Mac & cheese 3
```

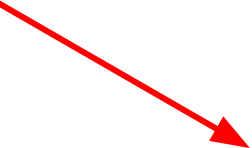
Punctuation characters only

```
sheikhwalter ~/OSC/session  
>> grep '[:punct:]' grocery  
Mac & cheese 3
```

## 4. Character classes

### Examples:

It didn't include the new line character (`\n`)



```
sheikhwalter ~/OSC/session  
>> cat grocery  
Bread 20
```

```
Oranges 10
```

```
Mac & cheese 3
```

```
sheikhwalter ~/OSC/session  
>> grep '[:,print:]' grocery  
Bread 20  
Oranges 10  
Mac & cheese 3
```



# **Hands on**

# **3**



# Search in “random.txt” for dates

Note: dates can be **dd/mm/yyyy** or **dd-mm-yyyy**

## Solution:

```
grep '[0-9]\{2\}[/-][0-9]\{2\}[/-][0-9]\{4\}' random.txt
```

```
>> grep '[0-9]\{2\}[-/][0-9]\{2\}[-/][0-9]\{4\}' random.txt  
12-03-2019  
16/08/2025  
31-12-1999  
21/05/2018  
01/12/2020
```

## 5. Grouping and Capturing

Grouping lets you treat part of a pattern as a single unit which is very useful in many scenarios.

It works by writing the grouped pattern between **escaped braces** `\( \)`

```
grep '\(pattern\)' filename
```

## 5. Grouping and Capturing

Grouping allows you to:

1. Apply repetition to a group

Each group is treated as a single entity which makes applying repetition a lot easier

```
sheikhwalter ~/OSC/session  
>> grep '\(waaa\)\{3\}' baby  
waaawaaawaaa
```

```
sheikhwalter ~/OSC/session  
>> grep '\(go\)\{3\}' baby  
gogogo
```

```
sheikhwalter ~/OSC/session  
>> grep '\(ga\)\{3\}' baby  
gagaga
```

## 5. Grouping and Capturing

Grouping allows you to:

2. Use alternation with sub-patterns

```
sheikhwalter ~/OSC/session  
>> grep '\(Cat\|Dog\)' animals  
CatDog  
DogLionCat
```

Matches either **Cat** or **Dog**.

## 5. Grouping and Capturing

Grouping allows you to:

### 3. Capture part of a match for extraction or replacement

This is useful in grep if the captured group is used in another place in the pattern.

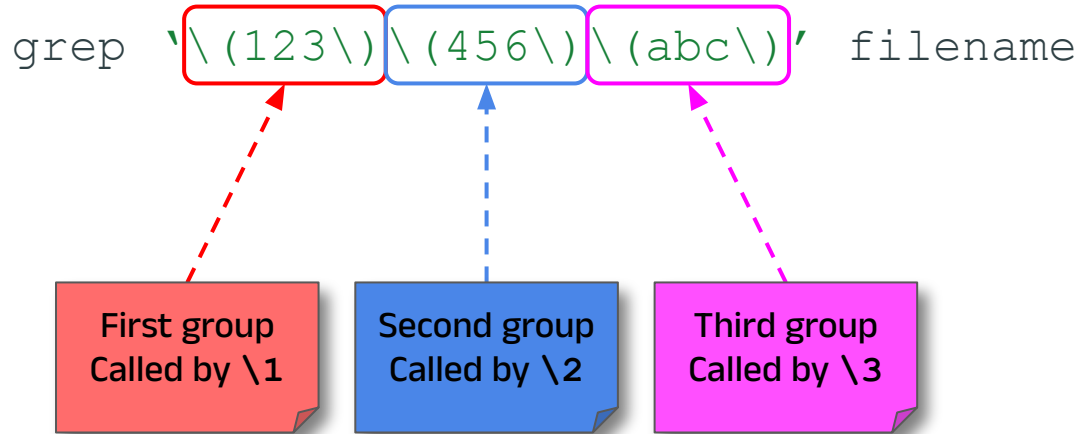
Suppose we have a pattern like this: `ab.....abab`

We can capture the group from the beginning and reuse it later on like this:

```
sheikhwalter ~/OSC/session
>> grep '\(ab\).....\1\{2\}' baby
ab321f3abab
ab5d8v3abab
```

## 5. Grouping and Capturing

How capturing works?





# Hands on

# 4

# Search in “random.txt” for emails

Example emails: john@example.com  
corp@dummy.co.uk  
200017001@fcis.asu.edu.eg

## Solution:

```
grep '^[^@ ]\+@[^@. ]\+\(\.[[:alpha:]]\{2,\}\)\+' random.txt
```

```
>> grep '^[^@ ]\+@[^@. ]\+\(\.[[:alpha:]]\{2,\}\)\+' random.txt  
john.doe@example.com  
jane_smith99@mail.org  
support@company.co.uk  
contact@website.com  
2029170321@fcis.asu.edu.eg
```



## 5. Grouping and Capturing

### What is `sed` command and why it needs capturing?

Capturing is useful in `grep`, but it's also important in commands like `sed` which is mostly used to replace or delete some parts in text.

Without digging into details, `sed` command uses Regex to match patterns then replace or delete them.

It works like this:

```
sed -option 's/pattern/replacement/flags' filename
```

<b>s:</b>	Substitute command (there are other commands but we won't dive into them now)
<b>pattern:</b>	The desired pattern we want to replace
<b>replacement:</b>	What will the pattern be replaced by
<b>flags:</b>	Modify how the substitution behaves
<b>option:</b>	There are different options for <code>sed</code> like <code>-i</code> which replaces the text directly inside the file

## 5. Grouping and Capturing

What is `sed` command and why it needs capturing?

In some cases we just need to change the order of some patterns or change the format

**For example:** You have names listed in the format “`Lastname, Firstname`” and you want to swap them to “`Firstname Lastname`”

With capturing, we can group the `Firstname` and `Lastname` sections then reuse them in the replacement in one single command!

Capturing in `sed` works the same as in `grep` by grouping with escaped braces `\ ( \)` and calling each group by it's number `\1 \2 \3`

We'll use the command below for demonstration:

```
sed 's/^\([[:alpha:]]+\), \([[:alpha:]]+\)$/\2 \1/w namesFormatted' names
```

The captured pattern

The replacement

Write the output in this file

## 5. Grouping and Capturing

What is `sed` command and why it needs capturing?

Here's the command in action:

```
~/OSC/Text-processing
>> cat names
Johnson, Alice
Smith, Bob
Brown, Charlie
Elkhayat, Abdulrahman
~/OSC/Text-processing
>> sed 's/^\([[:alpha:]]\+\), \([[:alpha:]]\+\)$/\2 \1/w namesFormatted' names
Alice Johnson
Bob Smith
Charlie Brown
Abdulrahman Elkhayat
~/OSC/Text-processing
>> cat namesFormatted
Alice Johnson
Bob Smith
Charlie Brown
Abdulrahman Elkhayat
```

Recalling previously captured groups

\1

\2



**Thanks!**