

# Intermediate Linux Course: Commandline and Basic Scripting

---

## Table of Content

---

Part I: More Commandline Tools .....	3
Command-line Tools.....	3
GZIP.....	3
TAR.....	3
GREP.....	4
SED .....	5
AWK .....	6
Hints .....	7
Quoting.....	7
Expanding and Escaping.....	8
PART II: Basic Shell Scripting.....	9
What is a Script? .....	9
Script Naming and Organization.....	9
Running a Script.....	9
Basic Structure of a Shellscript.....	10
Readability and Documentation .....	10
Command Grouping and Sequences .....	11
Control Structures .....	12
Conditional Statements.....	12
Loops .....	15
Making Scripts Flexible .....	16
Configurable Scripts .....	16
Defining your own Commandline Options and Arguments.....	17
Ensuring a Sensible Exit Status.....	17
Why is the exit status important after all? .....	18
Tips and Tricks.....	18
Script Debugging.....	18
Command Substitution .....	18
Create Temporary Files.....	19
Cleanup Temporary Files.....	19
About Bio-IT .....	20
Links and Further Information .....	20
Live-CDs .....	20
Acknowledgements .....	21
Index.....	22

# Part I: More Commandline Tools

## Command-line Tools

### GZIP

**gzip** is a compression/decompression tool.

When used on a file (without any parameters) it will compress it and replace the file by a compressed version with the extension '.gz' attached:

```
# ls textfile*
textfile
# gzip textfile
# ls textfile*
textfile.gz
#
```

To revert this / to uncompress, use the parameter -d:

```
# ls textfile*
textfile.gz
# gzip -d textfile
# ls textfile*
textfile
#
```



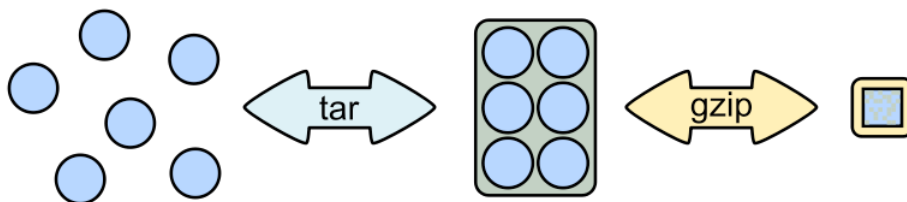
As a convenience, on most Linux systems, a shellscript named "gunzip" exists which simply calls "gzip -d"

### TAR

**tar** (tape archive) is a tool to handle archives. Initially it was created to combine multiple files/directories to be written onto tape, it is now the standard tool to collect files for distribution or archiving.

**tar** stores the permissions of the files within an archive and also copies special files (such as symlinks etc.), which makes it an ideal tool for archiving...

Usually **tar** is used in conjunction with a compression tool such as **gzip** to create a compressed archive:



source: [Th0msn80](#) (Wikipedia)

The most common command-line switches are:

- c create an archive
- t test an archive
- x extract an archive
- z use gzip compression
- f *filename* filename of the archive



Don't forget to specify the target filename. It needs to follow the -f parameter. Although you can combine options like such: "tar -czf archive.tar" the order matters, so "tar -cfz archive.tar" will not do what you want...

Creating an archive containing two files:

```
# tar -cf archive.tar textfile1 textfile2
#
```

Listing the contents of an archive:

```
# tar -tf archive.tar
textfile1
textfile2
#
```

Extracting an archive:

```
# tar -xf archive.tar
#
```

Creating and extracting a compressed archive containing two files:

```
# tar -czf archive.tar.gz textfile1 textfile2
# tar -xzf archive.tar.gz
#
```

## GREP

Find lines matching a pattern in textfiles

Usage: `grep [options] pattern file(s)`

```
# grep -i ensembl P04637.txt
DR   Ensembl; ENST00000269305; ENSP00000269305; ENSG00000141510.
DR   Ensembl; ENST00000359597; ENSP00000352610; ENSG00000141510.
DR   Ensembl; ENST00000419024; ENSP00000402130; ENSG00000141510.
DR   Ensembl; ENST00000420246; ENSP00000391127; ENSG00000141510.
DR   Ensembl; ENST00000445888; ENSP00000391478; ENSG00000141510.
DR   Ensembl; ENST00000455263; ENSP00000398846; ENSG00000141510.
#
```

Useful options:

- v: Print lines that do *not* match
- i: Search case-insensitive
- l: List files with matching lines, not the lines itself
- L: List files without matches
- c: Print count of matching lines for each file

Count the number of fasta sequences (they start with a ">") in a file:

```
# grep -c ">" twofiles.fasta
2
#
```

List all files containing the term "Ensembl":

```
# grep -l Ensembl *.txt
P04062.txt
P12931.txt
#
```

## SED

**sed** is a Stream E<sup>D</sup>itor, it modifies text (text can be a file or a pipe) on the fly.

Usage: 'sed command file',

The most common usecases are:

<b>Substitute</b> TEXT by REPLACEMENT:	's/TEXT/REPLACEMENT/'
<b>Transliterate</b> the characters x□a, and y□b:	'y/xy/ab/'
<b>Print</b> lines containing PATTERN:	'/PATTERN/p'
<b>Delete</b> lines containing PATTERN:	'/PATTERN/d'

```
# echo "This is text." | sed 's/text/replaced stuff/'
This is replaced stuff.
#
```

By default, text substitution are performed only once per line. You need to add a trailing '**g**' option, to make the substitution 'global' ('s/TEXT/REPLACEMENT/**g**'), meaning **all** occurrences in a line are substituted (not just the first in each line). Note the difference:

```
# echo "ACCAAGCATTGGAGGAATATCGTAGGTAAA" | sed 's/A/_/'
_CCAAGCATTGGAGGAATATCGTAGGTAAA
#
# echo "ACCAAGCATTGGAGGAATATCGTAGGTAAA" | sed 's/A/_/g'
_CC__GC__TTGG__GG__T_TCGT__GGT__
#
```

When used on a file, **sed** prints the file to standard output, replacing text as it goes along:

```
# echo "This is text" > textfile
# echo "This is even more text" >> textfile
# sed 's/text/stuff/' textfile
This is stuff
This is even more stuff
#
```

**sed** can also be used to print certain lines (not replacing text) that match a pattern. For this you leave out the leading '**s**' and just provide a pattern: '/PATTERN/p'. The trailing letter determines, what **sed** should do with the text that matches the pattern ('p': print, 'd': delete)

```
# sed '/more/p' textfile
This is text
This is even more text
This is even more text
#
```

As **sed** by default prints each line, you see the line that matched the pattern, printed twice. Use option '**-n**' to suppress default printing of lines.

```
# sed -n '/more/p' textfile
This is even more text
#
```

Delete lines matching the pattern:

```
# sed '/more/d' textfile
This is text
#
```

Multiple **sed** statements can be applied to the same input stream by prepending each by option **'-e'** (edit):

```
# sed -e 's/text/good stuff/' -e 's/This/That/' textfile
That is good stuff
That is even more good stuff
#
```

Normally, **sed** prints the text from a file to standard output. But you can also edit files in place. Be careful - this will change the file! The **'-i'** (in-place editing) won't print the output. As a safety measure, this option will ask for an extension that will be used to rename the original file to. For instance, the following option **'-i.bak'** will edit the file and rename the original file to **textfile.bak**:

```
# sed -i.bak 's/text/stuff/' textfile
# cat textfile
This is stuff
This is even more stuff
# cat textfile.bak
This is text
This is even more text
#
```

## AWK

**awk** is more than just a command, it is a complete text processing language (the name is an abbreviation of the author's names).

Each line of the input (file or pipe) is treated as a record and is broken into fields.

Generally, **awk** commands are of the form: **'condition { action }'**, where:

- condition is typically an expression
- action is a series of commands

If no condition is given, the action is applied to each line, otherwise just to the lines that match the condition.

```
# awk '{print}' textfile
This is text
This is even more text
# awk '/more/ {print}' textfile
This is even more text
#
```

**awk** reads each line of input and automatically splits the line into columns. These columns can be addressed via **\$1**, **\$2** and so on (**\$0** represents the whole line).

So an easy way to print or rearrange columns of text is:

```
# echo "Bob likes Sue" | awk '{print $3, $2, $1}'
Sue likes Bob
# echo "Master Obi-Wan has lost a planet" | awk '{print $4,$5,$6,$1,$2,$3}'
lost a planet Master Obi-Wan has
#
```



**awk** splits text by default on whitespace (spaces or tabs), which might not be ideal in all situations. To change the field separator (FS), use option '-F' (remember to quote the field separator):

```
# echo "field1,field2,field2" | awk -F',' '{print $2, $1}'
field2 field1
#
```

Note two things here: First, the field separator is not printed, and second, if you want to have space between the output fields, you actually need to separate them by a comma or they will be catenated together...

```
echo "field1,field2,field2" | awk -F',' '{print $1 $2 $3}'
field1field2field3
#
```

You can also combine the pattern matching and the column selection techniques:

```
# awk '/more/ {print $3}' textfile
even
#
```

**awk** really is powerful in filtering out columns, you can for instance print only certain columns of certain lines. Here we print the third column of those lines where the fourth column is 'more':

```
# awk '$4=="more" {print $3}' textfile
even
#
```

Note the double equal signs "==" to check for equality and note the quotes around "more".

If you want to match a field, but not exactly, you can use '~' instead of '==':

```
# awk '$4~"ore" {print $3}' textfile
even
#
```

## Hints

### Quoting

In Programming it is often necessary to "glue together" certain words. Usually, a program or the shell splits sentences by whitespace (space or tabulators) and treats each word individually. In order to tell the computer that certain words belong together, you need to "quote" them, using either single (') or double (") quotes. The difference between these two is generally that within double quotes, variables will be expanded, while everything within single quotes is treated as string literal.

When setting a variable, it doesn't matter which quotes you use:

```
# MYVAR=This is set
-bash: is: command not found
# MYVAR='This is set'
# echo $MYVAR
This is set
# MYVAR="This is set"
# echo $MYVAR
This is set
#
```

However, it does matter, when using (expanding) the variable:

Double quotes:

```
# export MYVAR=123
# echo "the variable is $MYVAR"
the variable is 123
# echo "the variable is set" | sed "s/set/$MYVAR/"
the variable is 123
#
```

Single quotes:

```
# export MYVAR=123
# echo 'the variable is $MYVAR'
the variable is $MYVAR
# echo "the variable is set" | sed 's/set/$MYVAR/'
the variable is $MYVAR
#
```

Weird things can happen when parsing data/text that contains quote characters:

```
# MYVAR='Don't worry'; echo $MYVAR
> # you need to press Ctrl-C to abort
# MYVAR="Don't worry"; echo $MYVAR
Don't worry
#
```

## Expanding and Escaping

You already learned how to expand a variable such that its value is used instead of its name:

```
# export MYVAR=123
# echo "the variable is $MYVAR"
the variable is 123
```

“Escaping” a variable is the opposite, ensuring that the literal variable name is used instead of its value:

```
# export MYVAR=123
# echo "the \$MYVAR variable is $MYVAR"
the $MYVAR variable is 123
```



The “escape character” is usually the backslash “\”.

# PART II: Basic Shell Scripting

## What is a Script?

A script is nothing else than a number of shell command place together in a file. The simplest script is maybe just a complex oneliner that you don't want to type each time again. More complex scripts are seasoned with control elements (conditions and loops) which allow for a sophisticated command flow. scripts might allow for configuration and customization, thus allowing one script to be flexibly used in several different environments.



Whatever you do in a script, you can also do on the commandline. This is also the first way to test your scripts step by step!

## Script Naming and Organization

It is good practice – though not technically required – to give your scripts an extension which specifies their type. I.e. “.sh” for Bourne Shell and Bourne Again Shell scripts, “.csh” for C-Shell scripts. Sometimes “.bash” for Bourne Again Shell scripts is used.

We recommend to either store all scripts in one location (e.g. ~/bin) and add this location to your `$PATH` variable or to store the scripts together with the files that are processed by the script.



If you use scripts to process data, then the scripts should probably be archived together with the data files

## Running a Script

There are basically three ways to run a script:

- a) the location to your script is not in your `$PATH` variable, then you have to specify the full path to the script:

```
# /here/is/my/script.sh
[...]
```

- b) the location to the script is in the `$PATH` variable, then you can simply type its name:

```
# script.sh
[...]
```

In both situations, the script will need to have execute permissions to be run. If for some reason you can only read but not execute the script, then it can still be run by

- c) specifying the interpreter. The full path (relative or absolute) to script has to be provided in this case, no matter whether the script location is already contained in `$PATH` or not:

```
# /bin/sh /here/is/my/script.sh
[...]
```

## Basic Structure of a Shellscript

Shellscripts have the following general structure:

- A line starting with "#!" which defines the interpreter (i.e. the program used to run the script). This line is called the "shebang line" and must be the first line in a script
- A section where the configuration takes place, e.g. paths, options and commands are defined and it is made sure, that all prerequisites are met
- A section where the actual processing is done. This includes error handling
- A controlled exit sequence, which includes cleaning up all temporary files and returning a sensible exit status

This is merely a recommendation to keep your scripts well structured. None of these sections are mandatory.

## Readability and Documentation

Make your script easily readable. Use comments and whitespace and avoid super compact but hardly understandable commandlines. Always take into account, that not only the shell, also human beings will probably have to read and understand your script.

Even if your script is very simple – document it! This helps others understand what you did, but – most important – it helps you remember what you did, when you have to reuse the script in the future.

Documentation is done either by writing comments into the script or by creating a special documentation file (README.txt or similar). Documenting in the script can be done in several ways:

- A preamble in the script, outlining the purpose, parameters and variables of the script as well as some information about authorship and and perhaps changes
- Within the script as blocks of text or "End of line" comments

To write a comments use the hash sign ("#"). Everything after a "#" is ignored when executing a script.

```
#!/bin/sh
```

Shebang line

```
#
# myscript.sh
#
# General purpose script for extracting Glycine
# occurrences in a datafile.
#
# Usage: myscript.sh datafile
#
# Exit values: 1: No datafile given or file
#               doesn't exist
#               2: No Glycine found
#
# Author: Me, myself and I
# Date:   Heidelberg, December 12., 2012
#
```

Preamble with a short description, usage information, authorship etc. etc.

```
# --- Configuration ---
GREPCMD=/bin/grep
DATAFILE=$1
```

Configuration

```
# --- Check prerequisites ---
# first check for $1
if [ -z $DATAFILE ]
then
    echo "No datafile given" 1>&2 # print on STDERR
    echo "USAGE: $0 datafile"
    exit 1
fi

# then check if the file exists
if [ ! -f $DATAFILE ]
then
    echo "Datafile $DATAFILE does not exist!" 1>&2
    exit 1
fi
```

Checking prerequisites and sane environment

```
# --- Now processing---
$GREPCMD -q Glycine $DATAFILE # Where is Glycine?
```

This is what you actually wanted to do

```
# --- Exit ---
if [ $? -eq 0 ]
then
    exit 0
else
    exit 2
fi
```

Ensure a valid and meaningful exit status

## Command Grouping and Sequences

Commands can be concatenated to be executed one after the other unconditionally or based on the success of the respective previous command:

**cmd1; cmd2** – execute commands in sequence

Create a directory and change into it

```
# mkdir a; cd a
```

**cmd1 && cmd2** – execute cmd2 only if cmd1 was successful

Confirm that /etc exists

```
# cd /etc && echo "/etc exists"
```

**cmd1 || cmd2** – execute cmd2 only if cmd1 was **not** successful

Warn if /etc doesn't exist

```
# cd /etc || echo "/etc is missing!"
```

**( cmds )** – groups commands to create one single output stream. The commands are run in a subshell (i.e. a new shell is opened to run them)

Change into /etc and list content. You are still in the same directory as you were before

```
# pwd
/home/fthommen
# (cd /etc; ls)
# pwd
/home/fthommen
#
```

**{ cmds; }** – groups commands to create one single output stream. The commands are run in the current (!) shell. The opening “{” must be followed by a blank and the last command must be succeeded by a “;”

Change into `/etc` and list content. You are still in `/etc` after the bracketed expression (compare to the example above)

```
# pwd
/home/fthommen
# { cd /etc; ls; }
[...]
# pwd
/etc
#
```

## Control Structures

The following syntax elements will be described for sh/bash and for csh/tcsh. However since this course is mainly about sh/bash, examples will only be given for sh/bash. Some notes about csh/tcsh specialities might be given in the text.

This is only a selection of the most useful or most common elements. There are much more in the manpages. All shells offer myriads of possibilities which cannot possibly be demonstrated in this course.

Some of the described features might be specific to bash and not be available in a classical Bourne Shell on other systems.

## Conditional Statements

### *if – then – else*

This is the most basic conditional statement: Do something depending on certain conditions. The basic syntax is

#### *sh/bash*

```
if condition1
then
    statements
elif condition2
    more statements
[...]
else
    even more statements
fi
```

#### *csh/tcsh*

```
if (condition) then
    statements
else if (condition2) then
    more statements
[...]
else
    even more statements
endif
```

Conditions can be a) the exit status of a command or b) the evaluation of a logical or arithmetic expression:

- a) **Evaluating the exit status of a command:** Simply use the command as condition

Example

```
if grep -q root /etc/passwd
then
    echo root user found
else
    echo No root user found
fi
```



To evaluate the exit status of a command in `csch/tcsh`, it must be placed within curly brackets with blanks separating the brackets from the command: `if ( { grep -q root /etc/passwd } ) then [...]`



Redirect the output of the command to be evaluated to `/dev/null` if you are only interested in the exit status and if the command doesn't have a "quiet" option.

Note: Redirection of commands in conditions does not work for `csch/tcsh`

- b) **Evaluating of conditions or comparisons:** Conditions and comparisons are evaluated using a special command `test` which is usually written as `"["` (no joke!). As `"["` is a command, it **must** be followed by a blank. As a speciality the `"["` command **must** be ended with `" ]"` (note the preceding blank here)



In `csch/tcsh` the `test/[` command is not needed. Conditions and comparisons are directly placed within the round braces.

*sh/bash*

*csch/tcsh*

#### File conditions

<code>-e file</code>	<i>file</i> exists	<code>-e file</code>
<code>-f file</code>	<i>file</i> exists and is a regular file	<code>-f file</code>
<code>-d file</code>	<i>file</i> exists and is a directory	<code>-d file</code>
<code>-r file</code>	<i>file</i> exists and is readable	<code>-r file</code>
<code>-w file</code>	<i>file</i> exists and is writable	<code>-w file</code>
<code>-x file</code>	<i>file</i> exists and is executable	<code>-x file</code>
<code>-s file</code>	<i>file</i> exists and has a size > 0	
	<i>file</i> exists and has zero size	<code>-z file</code>

#### String Comparisons

<code>-n s1</code>	String <i>s1</i> has non-zero length	
<code>-z s1</code>	String <i>s1</i> has zero length	
<code>s1 = s2</code>	Strings <i>s1</i> and <i>s2</i> are identical	<code>s1 == s2</code>
<code>s1 != s2</code>	Strings <i>s1</i> and <i>s2</i> differ	<code>s1 != s2</code>
<code>string</code>	String <i>string</i> is not null	

#### Integer Comparisons

<code>n1 -eq n2</code>	<i>n1</i> equals <i>n2</i>	<code>n1 == n2</code>
<code>n1 -ge n2</code>	<i>n1</i> is greater than or equal to <i>n2</i>	<code>n1 &gt;= n2</code>
<code>n1 -gt n2</code>	<i>n1</i> is greater than <i>n2</i>	<code>n1 &gt; n2</code>
<code>n1 -le n2</code>	<i>n1</i> is less than or equal to <i>n2</i>	<code>n1 &lt;= n2</code>
<code>n1 -lt n2</code>	<i>n1</i> is less than <i>n2</i>	<code>n1 &lt; n2</code>
<code>n1 -ne n2</code>	<i>n1</i> is not equal to <i>n2</i>	<code>n1 != n2</code>

#### Combination of conditions

<code>! cond</code>	True if condition <i>cond</i> is not true	<code>! cond</code>
<code>cond1 -a cond2</code>	True if conditions <i>cond1</i> and <i>cond2</i> are both true	<code>cond1 &amp;&amp; cond2</code>
<code>cond1 -o cond2</code>	True if conditions <i>cond1</i> or <i>cond2</i> is true	<code>cond1    cond2</code>

Examples: Test for the existence of `/etc/passwd`

```
if [ -e /etc/passwd ]
then
    echo /etc/passwd exists
else
    echo /etc/passwd does NOT exist
fi
```

OR

```
if test -e /etc/passwd
then
    echo /etc/passwd exists
else
    echo /etc/passwd does NOT exist
fi
```

### case

The case statement implements a more compact and better readable form of `if – elif – elif – elif` etc. Use this if your variable (and you can only check for variables with case) can have a distinct number of valid values. A typical usage of case will follow later.

The basic syntax is

*sh/bash*

```
case variable in
    pattern1)
        statements
        ;;
    pattern2)
        statements
        ;;
    [...]
    *)
        statements
        ;;
esac
```

*csh/tcsh*

```
switch (variable)
    case pattern1:
        statements
        breaksw
    case pattern2:
        statements
        breaksw
    default:
        statements
endsw
```



"\*", "?" and "[...]" can be used for the patterns



The `*)` (sh/bash) and `default:` (csh/tcsh) patterns are "catch-all" patterns which match everything not matched above. It is often used to detect invalid values of variable.



Multiple patterns can be handled by separating them with "|" in sh/bash or by successive `case` statements in csh/tcsh.

Examples: Check if /opt/ or /usr/ paths are contained in \$PATH

```
case $PATH in
  */opt/* | */usr/* )
    echo /opt/ or /usr/ paths found in $PATH
    ;;
  *)
    echo '/opt and /usr are not contained in $PATH'
    ;;
esac
```

## Loops

### for / foreach

The for and foreach statements respectively will loop through a list of given values and run the given statements for each run:

*sh/bash*

```
for variable in list
do
  statements
done
```

*csh/tcsh*

```
foreach variable (list)
  commands
end
```

*list* is a list of strings, separated by whitespaces

Examples: List all files in /tmp in a bulleted list

```
for FILE in /tmp/*
do
  echo " * $FILE"
done
```

OR

```
for FILE in `ls /tmp`
do
  echo " * $FILE"
done
```

### while / until

The while and until loops execute your commands while (or until respectively) a certain condition is met

*sh/bash*

```
while condition
do
  statements
done
```

```
until condition
do
  statements
done
```

*csh/tcsh*

```
while (condition)
  commands
end
```

**N.A.**

The conditions are constructed the same way as those used in if statements.

### "Manual" loop control

Instead of (or additionally to) the built-in loop control in for/foreach, while and until loops, you can control exiting and continuing them with "break" and "continue":

**break** “breaks out” of the innermost loop (loops can be nested!) and continues after the end of the loop.

**continue** skips the rest of the current (innermost) loop and starts the next iteration

## Making Scripts Flexible

Scripts are most useful, if they can be reused. Copying scripts and changing them to fit the new situation is time-consuming and error-prone. Additionally if you add an improvement to the current script, then all previous versions will stay without it. Having one script with the possibility to configure it, is usually the better way. Customization of scripts can be achieved by either using variables or by adding the possibility to use your own commandline options and arguments.

### Configurable Scripts

#### Using Variables

Any value – be it paths, commands or options – that are specific to individual applications or your script, should not be “hardcoded” (i.e. used literally within the script) but assigned to variables:

Bad example: You have to change two instances of the path each time you want to list an other directory:

```
#!/bin/sh

echo "The directory /etc contains the following files:"
ls /etc
```

Good example: The path is now in a variable and only one instance has to be changed each time (less work, less errors)

```
#!/bin/sh

$MYDIR=/etc

echo "The directory $MYDIR contains the following files:"
ls $MYDIR
```

Of course, you’ll still have to modify the script each time you want to list the content of an other directory. A more flexible way of customization would be to use a settings file.

#### Using a Settings File

Instead of having your configurable section within the script, it can be “outsourced” in its own file. This file is basically a shellscript which is run within the primary script. To run commands from a file within the *current* environment, the commands **source** (bash, csh/tcsh) or **.** (dot) (sh/bash) are used:

The settings file, e.g. **settings.ini**:

```
$MYDIR=/etc
```

The script:

```
#!/bin/sh

. ./settings.ini

echo "The directory $MYDIR contains the following files:"
ls $MYDIR
```

## Defining your own Commandline Options and Arguments

The best way to configure a script is to allow for your own commandline options and arguments. Commandline arguments are available the script as so-called positional parameters \$1, \$2, \$3 etc. \$0 contains the name of the script.

If you run the script

```
#!/bin/sh

echo The script is $0
echo The first commandline option is $1
echo The second commandline option is $2
```

With two arguments, you'll get the following output:

```
# ./script.sh ABC DEF
The script is ./script.sh
The first commandline option is ABC
The second commandline option is DEF
#
```

In many cases you'll not know how many parameters are given on the commandline. In these cases you can use `shift` to loop through them. `shift` removes \$1 and moves all other positional parameters one position to the right: \$2 becomes \$1, \$3 becomes \$2 etc.

Some more variables are important when dealing with commandline parameters:

- \$#: Current number of positional parameters
- \*: All positional parameters
- @: All positional parameters. If used within double quotes ("@"), then it will expand to the list of all positional parameters, where each parameter is individually quoted

With the help of \$#, `shift`, `case` and the positional parameters we can now check all the commandline parameters:

```
while [ "$#" -gt 0 ]
do
  case $1 in
    -h) echo "Sorry, no help available!" # not very helpful, is it?
        exit 1                         # exit with error
        ;;

    -v) VERBOSE=1                       # we may use $VERBOSE later
        ;;

    -f) shift
        FILE=$1                        # Aha, -f requires an
                                       # additional argument
        ;;

    *)  echo "Wrong parameter!"
        exit 1                         # exit with error
  esac
  shift
done
```

## Ensuring a Sensible Exit Status

If you don't provide your own exit status, then the script will return the exit status of the last executed command. In many cases this might be what you want, but very often it isn't.

Consider the following script which is a real example from real life and happened to me personally:

```
#!/bin/sh

[... do something that fails ...]

echo "End of the script"
```

This script will *always* succeed, as the `echo` command hardly ever fails. You will – from the exit status of the script – never be able to detect, that something went wrong. Instead in such cases we should handle the exit codes of the commands we run within the script. The most important variable in this context is

`$?`: The exit status of the last run command

With its help we can keep track of the exit status of all our important processing steps and finally return a sensible value:

```
#!/bin/sh
mystatus=0;

[... do something that might fail ...]
if [ $? -ne 0 ]
then
    mystatus=1
fi

[... do something else that might fail, too ...]
[ $? -ne 0 ] && mystatus=1      # same as above.  Do you understand
                                # this?

echo "End of the script"
exit $mystatus
```

### Why is the exit status important after all?

First when you use your script within other scripts, you’ll probably need to be able to check, if it has succeeded. There might be other ways (e.g. checking outputfiles for certain strings, checking directly the textual output of the script etc.), but these ways are usually cumbersome and require lots of coding. Exit values are easy to check.

Second: Other tools and systems might also use the exit status of your script. E.g. the cluster system uses your job’s exit status to assess, if it has run successfully or not. Returning success even in case of failure will result in lots of complications in case a problem occurs. It took us several days to find the bug above.

## Tips and Tricks

### Script Debugging

`sh/bash` and `csh/tcsh` have both an option “`-x`” which helps debugging a script by echoing each command before executing it.

### Command Substitution

You can use the output of a command and assign it to a variable or use it right away as text string, by using the command substitution operators “```” (backticks, backquotes) or “`$(...)`”. The backtick operator works in all shells, while `$(...)` only works in bash.

Three variants for the same (print out who you are in English text):

```
# ME=`whoami`  
# echo I am $ME  
I am fthommen  
#
```

```
# ME=$(whoami)  
# echo I am $ME  
I am fthommen  
#
```

```
# echo I am `whoami`  
I am fthommen  
#
```

## Create Temporary Files

You can create temporary files with `mktemp`. By default it will create a new file in `/tmp` and print its name:

```
# mktemp  
/tmp/tmp.Yaafh19370  
#
```

## Cleanup Temporary Files

It is considerate, good practice and sometimes even important, to clean up temporary data before ending a script. A simple way – which will not cover all cases, though – could be to store all created temporary files in a variable and remove them all before exiting the script:

```
#!/bin/sh  
ALL_TEMPFILES=""      # store a list of all temporary files here  
  
TEMPFILE1=`mktemp`  
ALL_TEMPFILES="$ALL_TEMPFILES $TEMPFILE1"  
  
TEMPFILE2=`mktemp`  
ALL_TEMPFILES="$ALL_TEMPFILES $TEMPFILE2"  
  
[... process, process, process ...]  
  
rm -f $ALL_TEMPFILES  
exit
```

## About Bio-IT

---

Bio-IT is a community project aiming to develop and strengthen the bioinformatics user community at EMBL Heidelberg. It is made up of members across the different EMBL Heidelberg units and core facilities. The project works to achieve these aims, firstly, by providing a forum for discussing and sharing information and ideas on computational biology and bioinformatics, focused on the Bio-IT portal <http://bio-it.embl.de>. Secondly, we organise and participate in a range of different networking and social activities aiming to strengthen ties across the community.

## Links and Further Information

---

- A full 500 page book about the Linux commandline for free(!): LinuxCommand.org (<http://linuxcommand.org/>)
- Another nice introduction: “A beginner's guide to UNIX/Linux” (<http://www.mn.uio.no/astro/english/services/it/help/basic-services/linux/guide.html>)
- The “commandline starter” chapter of an O'Reilly book: Learning Debian GNU/Linux – Issuing Linux Commands ([http://oreilly.com/openbook/debian/book/ch04\\_01.html](http://oreilly.com/openbook/debian/book/ch04_01.html))
- A nice introduction to Linux/UNIX file permissions: “chmod Tutorial” (<http://catcode.com/teachmod/>)
- Linux Cheatsheets (<http://www.cheat-sheets.org/#Linux>)
- For the technically interested:  
Linux Filesystem Hierarchy Standard (<http://www.pathname.com/fhs/>) and Linux Standard Base (<http://www.linuxfoundation.org/collaborate/workgroups/lsb>)
- Unix commands applied to bioinformatics ([http://rous.mit.edu/index.php/Unix\\_commands\\_applied\\_to\\_bioinformatics](http://rous.mit.edu/index.php/Unix_commands_applied_to_bioinformatics))
- BioPieces (<http://code.google.com/p/biopieces/>)

Real printed paper books:

- Dietz, M., , Praxiskurs Unix-Shell, O'Reilly (highly recommended!)
- Herold, H., awk & sed, Addison-Wesley
- Robbins, A., sed & awk Pocket Reference, O'Reilly
- Robbins, A. and Beebe, N., Classic Shell Scripting, O'Reilly
- Siever, E. et al., Linux in a Nutshell, O'Reilly

## Live-CDs

---

A Live-CD is a complete bootable computer operating system which runs in the computer's memory, rather than loading from the hard disk drive. It allows users to experience and evaluate an operating system without installing it or making any changes to the existing operating system on the computer.

Just download an ISO-Image, burn it onto a CD/DVD and insert it into your DVD-Drive to boot your computer with Linux!

**Fedora Live CD**<http://fedoraproject.org/wiki/FedoraLiveCD>

This Live CD contains everything the Fedora Linux operating system has to offer and it's everything you need to try out Fedora — you don't have to erase anything on your current system to try it out, and it won't put your files at risk. Take Fedora for a test drive, and if you like it, you can install Fedora directly to your hard drive straight from the Live Media desktop.

**Knoppix**<http://knopper.net/knoppix>

Knoppix is an operating system based on Debian designed to be run directly from a CD / DVD or a USB flash drive, one of the first of its kind for any operating system. When starting a program, it is loaded from the removable medium and decompressed into a RAM drive. The decompression is transparent and on-the-fly. More than 1000 software packages are included on the CD edition and more than 2600 are included on the DVD edition. Up to 9 gigabytes can be stored on the DVD in compressed form.

**BioKnoppix**<http://bioknoppix.hpcf.upr.edu/>

Bioknoppix is a customized distribution of Knoppix Linux Live CD. With this distribution you just boot from the CD and you have a fully functional Linux OS with open source applications targeted for the molecular biologist. Beside using RAM, Bioknoppix doesn't touch the host computer, being ideal for demonstrations, molecular biology students, workshops, etc.

**Vigyaan**<http://www.vigyaancd.org>

Vigyaan is an electronic workbench for bioinformatics, computational biology and computational chemistry. It has been designed to meet the needs of both beginners and experts.

**BioSlax**<http://www.bioslax.com/>

BioSLAX is a live CD/DVD suite of bioinformatics tools that has been released by the resource team of the Bioinformatics Center (BIC), National University of Singapore (NUS).

**Acknowledgements**

---

EMBL Logo © EMBL Heidelberg

Bio-IT Logo © Bio-IT Project. EMBL Heidelberg

All other graphics as declared

## Index

.	16	gzip	3
<b>\$</b>		<b>I</b>	
\$?	18	if	12
\$@	17	Integer Comparisons	13
\$*	17	<b>L</b>	
\$#	17	Loops	15
<b>A</b>		for	15
awk	6	foreach	15
<b>B</b>		until	15
backtick operator	18	while	15
<b>break</b>	16	<b>M</b>	
<b>C</b>		mktemp	19
case	14	<b>P</b>	
Combination of conditions	13	positional parameters	17
Command Substitution	18	<b>Q</b>	
conditional statements		Quoting	7
case	14	Double quotes	8
if	12	Single quotes	8
<b>continue</b>	16	<b>S</b>	
Control Structures	12	Script Debugging	18
<b>D</b>		<b>sed</b>	5
Debugging	18	source	16
dot command	16	String Comparisons	13
<b>E</b>		<b>T</b>	
Escape	8	tar	3
Escaping	8	Temporary Files	19
Exit Status	17	Cleanup	19
Expand	8	Creation	19
<b>F</b>		<b>U</b>	
File conditions	13	until	15
for	15	<b>W</b>	
foreach	15	while	15
<b>G</b>			
grep	4		

# Input-Output (IO) Redirections

Three IO “channels” are available by default:

- **Standard input (STDIN, Number: 0):** The input for your program, normally your keyboard but can be an other program (when using pipes or IO redirection)
- **Standard output (STDOUT, Number: 1):** Where your program writes its regular output to. Normally your terminal
- **Standard error (STDERR, Number: 2):** Where your programs normally write their error message to. Normally your terminal

Input, output and error messages can be redirected from their default “targets” go others. If using the file descriptor numbers (0, 1, 2) in redirections, then there must be no whitespace between the numbers and the redirection operators.

Redirect to `/dev/null` to discard the output.

<code>cmd &gt; afile</code>	Write the output of <code>cmd</code> into <code>afile</code> . This will overwrite <code>afile</code> .
<code>cmd &gt;&gt; afile</code>	Write the output of <code>cmd</code> into <code>afile</code> . This will add to <code>afile</code>
<code>cmd &gt; /dev/null</code>	Discard the output of <code>cmd</code>
<code>cmd &gt; afile 2&gt;&amp;1</code>	Write the output of <code>cmd</code> into <code>afile</code> (overwriting the file!) and write STDERR to the same place
<code>cmd &gt;&gt; afile 2&gt;&amp;1</code>	Add the output and error messages of <code>cmd</code> into <code>afile</code>
<code>cmd &gt; afile 2&gt; afile</code>	Same as above
<code>cmd &gt;&gt; afile 2&gt;/dev/null</code>	Add the output of <code>cmd</code> to <code>afile</code> and discard error messages
<code>cmd &gt; /dev/null 2&gt;&amp;1</code> <code>cmd &gt; /dev/null 2&gt;/dev/null</code> <code>cmd &gt;&amp; /dev/null</code>	Three time the same: Discard output and error messages completely
<code>cmd1 &lt; cmd2</code>	Use output of <code>cmd2</code> as standard input for <code>cmd1</code>

See also <http://www.catonmat.net/blog/bash-one-liners-explained-part-three/>,  
<http://www.catonmat.net/blog/bash-redirections-cheat-sheet/> and  
[http://wiki.bash-hackers.org/howto/redirection\\_tutorial](http://wiki.bash-hackers.org/howto/redirection_tutorial).

# Variables

The shell knows two types of variables: “Local” shell variables and “global” exported environment variables. By convention, environment variables are written in uppercase letters.

Shell variables are only available to the current not inherited when you start another shell or script from the commandline. Consequently, these variables will not be available for your shells/scripts.

Environment variables are inherited in shells and scripts started from your current.

## Setting, Exporting and Removing Variables

Variables are set (created) by assigning them a value

```
# MYVAR=something
```

There must be no whitespace around the equal sign. To create an environment variable, **export** is used. You can either export while assigning a value or in a separate step. Both of the following procedures are equivalent:

```
# export MYGLOBALVAR="something else"
```

```
# MYGLOBALVAR="something else"  
# export MYGLOBALVAR      # No "$" in front of the variable!
```

Variables are removed with **unset**:

```
# unset MYVAR
```

Assigning a variable an empty value (**MYVAR=**) will *not* remove it!

## Listing Variables

You can list all your current environment variables with **env** and all shell variables with **set**. The list of shell variables will also contain all environment variables

## Variable Inheritance

Only environment variables will be available in shells and scripts started from your current shell. However in shell commands run in subshells (i.e. commands run within round brackets) also local (shell) variables of your current shell are available.

## Examples

---

Consider the following small shellscript `vartest.sh`:

```
#!/bin/sh
echo $MYLOCALVAR
echo $MYGLOBALVAR
echo -----
```

We will use it in the following examples to illustrate the various variable inheritances:

```
# export MYGLOBALVAR="I am global"
# MYLOCALVAR="I am local"
#
# ./vartest.sh

I am global
-----
#
# . ./vartest.sh
I am local
I am global
-----
#
# (echo $MYGLOBALVAR; echo $MYLOCALVAR)
I am global
I am local
#
```

Set the variables

Run the script normally, i.e. in a new shell

“source” the script, i.e. run it within your current shell

Access the variables in a subshell

# Exercises

## Commandline tools

### I. TAR & GZIP

1. Use gzip to compress the file P12931.txt
2. Decompress the resulting file P12931.txt.gz (revert previous command)
3. Use tar to create an archive containing all fasta files in the current directory into an archive called "fastafiles.tar"
4. Use gzip to compress the archive "fastafiles.tar".
5. How can you achieve the two previous steps "using tar to create archive" and "gzip the archive" in one command?
6. Test (list the contents of) the compressed archive "fastafiles.tar.gz"
7. Download the compressed PDB file for entry 1Y57 from rcsb.org (eg. "wget <http://www.rcsb.org/pdb/files/1Y57.pdb.gz>") and decompress it.

### II. GREP

1. Which of the DNA files ENST.\* contains 'TATATCTAA' as part of the sequence?
2. Which of the DNA files ENST.\* contains 'CAACAAA' as part of the sequence?
3. Considering the previous example, would you consider grep a suitable tool to perform motif searches? Why not? Try to find the pattern 'CAACAAA' by manual inspection of the first two lines of each sequence.
4. Count the number of ATOMs (lines starting with 'ATOM') in the file 1Y57.pdb. Does this number agree with the annotated number of atoms (Search the REMARKs for 'protein atoms')

### III. SED

1. Use sed to print only those lines that contain "version" in the files P12931.txt and P04062.txt
2. Use sed to change the text 'sequence version 3' to 'sequence version 4' in the files P12931.txt and P04062.txt (without actually changing the files, just printing)
3. Use sed to update the text 'entry version 3' to 'entry version 4' in the files P12931.txt and P05480.txt (this time, make the changes directly in the files)
4. Replace all occurrences of 'r' by 'l' and 'l' by 'r' (at the same time) in the file PROTEINS.txt (so that 'structural' becomes 'stluctular')

### IV. AWK

1. Use awk to print only those lines that contain "version" in the files P12931.txt and P05480.txt and compare the procedure to sed.
2. For all FASTA files print only the second item of the header (split on "|")  
eg. for ">sp|P12931|SRC\_HUMAN Proto-oncogene", print only "P12931"
3. The file 'P12931.csv' contains phosphorylation sites in the protein P12931. (If the file 'P12931.csv' does not exist, use 'wget' to download it from "<http://phospho.elm.eu.org/byAccession/P12931.csv>").
  - a. Column three of this file lists the amino acid position of the phosphorylation site. You are only interested in position 17 of the protein. Try to use 'grep' to filter out all these lines containing '17'.
  - b. Now use awk to show all lines containing '17'.

- c. Next try show only those lines where column three equals 17 (Hint: The file is semicolon-separated...).
- d. Finally print the PMIDs (column 6) of all lines that contain '17' in column 3.

## Quoting and Escaping

1. Familiarize yourself with quoting and escaping.
  - a. Run the following commands to see the difference between single and double quotes when expanding variables:

```
# echo "$HOSTNAME"
# echo '$HOSTNAME'
```

- b. Next, use ssh to login to a different machine to run the same command there, again using both quoting methods:

```
# ssh teach01@pc-atcteach01 'echo $HOSTNAME'
```

```
# ssh teach01@pc-atcteach01 "echo $HOSTNAME"
```

Closely inspect the results; is that what you were expecting? Discuss this with your neighbour.

## Basic scripting

1. **gzip** can use nine different levels of compression. Read 'gzip -help' or 'man gzip' to find out how to compress a file while keeping the original file unchanged and how to use these different compression levels. Then write a **for**-loop to iterate over all possible levels, compressing the file 1Y57.pdb and compare the sizes of the resulting nine compressed files.
2. Write a script that combines the two files "P00523.fasta" and "P12931.fasta" into one file called "twofiles.fasta" and run /g/software/bin/clustalw2 on these files. (OPTIONAL) Finally, view the resulting alignment "twofiles.aln" using /g/software/bin/clustalx.
3.
  - a. Write a script that tests the existence of a directory and if it exists, prints out the content list as a nice bulleted list (using "\*" as bullets). If the directory does not exist, then print out the error message of "ls <directory>".
  - b. Modify the above script, so that the directory can be given as commandline argument. Print out an error message, if no commandline argument is given or if the given argument is not a directory.
  - c. Modify the script, so that the directory must be given as an option (e.g. myscript.sh -d <directory>)
  - d. If you used file testing expressions, then modify the script to only use "ls" and its exit status to test for the directory existence. "ls" should only be run once (!) in the course of the whole script (Hint: Use temporary files)
4. (OPTIONAL) Write a for loop to get all (unique) PDB ids (column 13) from file P12931.csv and retrieve the corresponding PDB file using this URL scheme: "http://www.rcsb.org/pdb/files/ID.pdb"
 

eg. "wget <http://www.rcsb.org/pdb/files/1Y57.pdb>"

# Solutions to the Exercises

## Commandline tools

### I. TAR & GZIP

1. Use gzip to compress the file P12931.txt

```
# gzip P12931.txt
```

2. Decompress the resulting file P12931.txt.gz (revert previous command)

```
# gunzip P12931.txt.gz # or gzip -d P12931.txt.gz
```

3. Use tar to create an archive containing all fasta files in the current directory into an archive called "fastafiles.tar"

```
# tar -c -f fastafiles.tar *.fasta
```

4. Use gzip to compress the archive "fastafiles.tar".

```
# gzip fastafiles.tar
```

5. How can you achieve the two previous steps "using tar to create archive" and "gzip the archive" in one command?

```
# tar -c -z -f fastafiles.tar.gz *.fasta # note the -z
```

6. Test (list the contents of) the compressed archive "fastafiles.tar.gz"

```
# tar -tf fastafiles.tar.gz
```

7. Download the compressed PDB file for entry 1Y57 from rcsb.org (eg. "wget <http://www.rcsb.org/pdb/files/1Y57.pdb.gz>") and decompress it.

```
# wget "http://www.rcsb.org/pdb/files/1Y57.pdb.gz"
# gunzip 1Y57.pdb.gz
```

### II. GREP

1. Which of the DNA files ENST.\* contains 'TATATCTAA' as part of the sequence?

```
# grep TATATCTAA ENST.*
ENST00000380152.fasta:ACGGAAGAATGTGAGAAAAATAAGCAGGACACAATTA
CAACTAAAAAATATATCTAA
ENST00000544455.fasta:ACGGAAGAATGTGAGAAAAATAAGCAGGACACAATTA
CAACTAAAAAATATATCTAA
```

2. List only the names of the DNA files ENST.\* that contain 'CAACAAA' as part of the sequence.

```
# grep -l CAACAAA ENST.*
ENST00000380152.fasta
ENST00000544455.fasta
```

3. Considering the previous example, would you consider grep a suitable tool to perform motif searches? Why not? Try to find the pattern 'CAACAAA' by manual inspection of the first two lines of each sequence.

Answer: When using grep as a motif searching tool, you need to keep in mind that grep (like sed and awk) is line-oriented, meaning that by default it only searches for a

given motif in a single line. In the given example, upon manual inspection you will find the given motif also in the file ENST00000530893.fasta, which grep missed. You would need to think about how to do multi-line searches (eg. Removing line-breaks etc.)

- Count the number of ATOMs (lines starting with 'ATOM') in the file 1Y57.pdb.
- Does this number agree with the annotated number of atoms (Search the REMARKs for 'protein atoms')

```
# grep -c "ATOM" 1Y57.pdb
3632
# grep -i "protein atoms" 1Y57.pdb
REMARK      3      PROTEIN ATOMS                : 3600
```

This means there are 3600 atoms annotated in this PDB file, however we counted 3632. This is because grep also counted any occurrence of "ATOM" within REMARKS. We can avoid this by either filtering out the remarks:

```
# grep -v "REMARK" 1Y57.pdb | grep -c "ATOM"
3600
```

...or by telling grep to only count those lines that start with "ATOM":

```
# grep "^ATOM" 1Y57.pdb
3600
```

### III. SED

- Use sed to print only those lines that contain "version" in the files P12931.txt and P04062.txt

```
# sed -n '/version/p' P12931.txt P04062.txt
```

- Use sed to change the text 'sequence version 3' to 'sequence version 4' in the files P12931.txt and P04062.txt (without actually changing the files, just printing)

```
# sed 's/sequence version 3/sequence version 4/'
P12931.txt P04062.txt
```

- Use sed to update the text 'entry version 3' to 'entry version 4' in the files P12931.txt and P05480.txt (this time, make the changes directly in the files)

```
# sed -i.bak 's/sequence version 3/sequence version 4/'
P12931.txt P04062.txt
```

- Replace all occurrences of 'r' by 'l' and 'l' by 'r' (at the same time) in the file PROTEINS.txt (so that 'structural' becomes 'stluctular')

```
# sed 'y/rRlL/lLrR/' PROTEINS.txt
```

### IV. AWK

- Use awk to print only those lines that contain "version" in the files P12931.txt and P05480.txt and think about how this procedure is different to sed.

```
# awk '/version/ {print}' P12931.txt P05480.txt
```

This is very similar to sed, you also have to use the slashes '/' to define the search pattern. However the sed notation is slight a little more concise...

- For all FASTA files print only the second item of the header (split on "|") eg. for ">sp|P12931|SRC\_HUMAN Proto-oncogene", print only "P12931"

```
# awk -F"|" " '/>/ {print $2}' *.fasta
```

3. The file 'P12931.csv' contains phosphorylation sites in the protein P12931. (If the file 'P12931.csv' does not exist, use 'wget' to download it from "<http://phospho.elm.eu.org/byAccession/P12931.csv>").

- a. Column three of this file lists the amino acid position of the phosphorylation site. You are only interested in position 17 of the protein. Try to use 'grep' to filter out all these lines containing '17'.

```
# grep 17 P12931.csv
```

- b. Now use awk to show all lines containing '17'.

```
# awk '/17/ {print}' P12931.csv
```

- c. Next try show only those lines where column three equals 17 (Hint: The file is semicolon-separated...).

```
# awk -F";" " '$3==17 {print}' P12931.csv
```

- d. Finally print the PMIDs (column 6) of all lines that contain '17' in column 3.

```
# awk -F";" " '$3==17 {print $6}' P12931.csv
```

## Quoting and Escaping

1. Familiarize yourself with quoting and escaping.
- a. Run the following commands to see the difference between single and double quotes when expanding variables:

```
# echo "$HOSTNAME"
# echo '$HOSTNAME'
```

- b. Next, use ssh to login to a different machine to run the same command there, again using both quoting methods:

```
# ssh pc-atcteach01 'echo $HOSTNAME'
```

```
# ssh pc-atcteach01 "echo $HOSTNAME"
```

Closely inspect the results; is that what you were expecting? Discuss this with your neighbour.

## Basic scripting

1. **gzip** can use nine different levels of compression. Read 'gzip -help' or 'man gzip' to find out how to compress a file while keeping the original file unchanged and how to use these different compression levels. Then write a **for**-loop to iterate over all possible levels, compressing the file 1Y57.pdb and compare the sizes of the resulting nine compressed files.

```
# --- Configuration ---
GZIP=gzip
DATAFILE=$1

# --- Check prerequisites ---
# first check for $1
if [ -z $1 ]
then
    echo ?No datafile given? 1>&2 # print on STDERR
    echo ?USAGE: $0 datafile?
    exit 1
fi

# then check if the file exists
if [ ! -f $DATAFILE ]
then
    echo ?Datafile $DATAFILE does not exist!? 1>&2
    exit 1
fi

# --- Now processing---
for level in 1 2 3 4 5 6 7 8 9
do
    echo "compressing $DATAFILE at level $level"
    $GZIP -c -$level $DATAFILE > $DATAFILE.$level.gz
done
```

- ~~2. Write a script that combines the two files "P00523.fasta" and "P12931.fasta" into one file called "twofiles.fasta" and run /g/software/bin/clustalw2 on these files. (OPTIONAL) Finally, view the resulting alignment "twofiles.aln" using /g/software/bin/clustalx.~~

3.

- a. Write a script that tests the existence of a directory and if it exists, prints out the content list as a nice bulleted list (using "\*" as bullets). If the directory does not exist, then print out the error message of "ls <directory>".

```
# --- Configuration ---
DIRECTORY=/tmp

# --- Now processing---
# check if the file exists
if [ ! -d $DIRECTORY ]
then
# DIRECTORY does not exist!
    ls $DIRECTORY
    exit 2
else
    for item in $(ls $DIRECTORY)
    do
        echo " * " $item
    done
fi

# --- Exit ---
if [ $? -eq 0 ]
then
    exit 0
else
    exit 2
fi
```

- b. Modify the above script, so that the directory can be given as commandline argument. Print out an error message, if no commandline argument is given or if the given argument is not a directory.

```
# --- Configuration ---
DIRECTORY=$1

# --- Check prerequisites ---
# first check for $1
if [ -z $1 ]
then
    echo ?No DIRECTORY given? 1>&2 # print on STDERR
    echo ?USAGE: $0 DIRECTORY?
    exit 1
fi

# then check if the file exists
if [ ! -d $DIRECTORY ]
then
    # DIRECTORY does not exist!
    ls $DIRECTORY
    exit 2
fi

# --- Now processing---
for item in $(ls $DIRECTORY)
do
    echo " * " $item
done
```

- c. Modify the script, so that the directory must be given as an option (e.g. `myscript.sh -d <directory>`)

```
# --- Check prerequisites ---
if [ "$#" -eq 0 ]
then echo "No argument given. Please specify a directory"
    exit 1
fi

while [ "$#" -gt 0 ]
do
    case $1 in
        -d) shift # next argument should be stored as variable
              DIRECTORY=$1
              ;;

        *)  echo "Wrong parameter!"
              exit 1                                     # exit with error
    esac
    shift
done

# then check if the file exists
if [ ! -d $DIRECTORY ]
then
    # DIRECTORY does not exist!
    ls $DIRECTORY
    exit 2
fi

# --- Now processing---
for item in $(ls $DIRECTORY)
do
    echo " * " $item
done

# --- Exit ---
if [ $? -eq 0 ]
then
    exit 0
else
    exit 3
fi
```

- d. If you used file testing expressions, then modify the script to only use "ls" and its exit status to test for the directory existence. "ls" should only be run once (!) in the course of the whole script (Hint: Use temporary files)

```
# --- Check prerequisites ---
if [ "$#" -eq 0 ]
then echo "No argument given. Please specify a directory"
    exit 1
fi

while [ "$#" -gt 0 ]
do
    case $1 in
        -d) shift # next argument should be stored
            DIRECTORY=$1
            ;;

        *) echo "Wrong parameter!"
           exit 1 # exit with error
    esac
    shift
done

# check if the file exists
TMPFILE=`mktemp -t tmp`
ERRORFILE=`mktemp -t tmp`
ls $DIRECTORY > $TMPFILE 2>$ERRORFILE
if [ $? -eq 0 ]
    then # DIRECTORY does exist!
        for item in `cat $TMPFILE`
        do
            echo " * " $item
        done
    else # DIRECTORY does NOT exist!
        cat $ERRORFILE
    fi

rm -f $TMPFILE
rm -f $ERRORFILE
```

4. (OPTIONAL) Write a for loop to get all (unique) PDB ids (column 13) from file P12931.csv and retrieve the corresponding PDB files using this URL scheme:  
"http://www.rcsb.org/pdb/files//ID.pdb"  
eg. "wget <http://www.rcsb.org/pdb/files/1Y57.pdb>"

```
# This will print the 13th column (PDB) from the given file
# Use sort -u to avoid retrieving the same file again
for i in $(awk -F';' '{print $13}' P12931.csv | sort -u);
do
    if [ $i = '-' ]
        then echo "not using " $i
    elif [ $i = 'PDB' ]
        then echo "not using " $i
    else
        wget "http://www.rcsb.org/pdb/files/$i.pdb" ;
    fi;
done
```