# Preparatory Course Informatics for Life Scientists

An Introduction to Python 5:

Filesystem, Files, and Command Line Interface

Philipp Thiel

September 13, 2022

# Modules

- Now we look into some interesting packages from the Python standard library

- There are many more but the chosen ones are really basic but important

- Package **os**
  $\rightarrow$ working with the file system

- Packages **sys** or **argparse**
  $\rightarrow$ working with input parameters to a Python program

- Working with files, which is built-in

# Package: os

- Package **os** has functionality to interact with the operating system
- It is part of the Python standard library

```python
import os

# Return the list of entries in the current working directory
print( os.listdir() )

# Return the list of entries in the directory 'folder2'
print( os.listdir('folder2') )

# Return the absolute path to the current working directory
print( os.getcwd() )

# Change the current working directory to 'folder2'
os.chdir('folder2')
```

# Package: os

```
myscript.py
1   import os
2
3   # Create a directory with name 'folderX'
4   # The directory to be created must not exist
5   os.mkdir('folderX')
6
7   # Rename directory 'folderX' (src) to 'folderY' (dest)
8   # src needs to exist, dest must not exist
9   # Works also for files
10  os.rename('folderX', 'folderY')
11
12  # Remove directory with name 'folderX'
13  # The directory to be deleted
14  #  - needs to exist
15  #  - needs to empty
16  os.rmdir('folderX')
17
```

# Package: os.path

- Module **os.path** provides extremely helpful path utilities
- Highly recommended to achieve cross-platform compatibility

myscript.py

```python
from os import path

p = '/home/charly/diary.txt'

# Returns the lowest (last) element in a path p
print( path.basename(p) )  # 'diary.txt'

# Returns the path without the lowest (last) element
print( path.dirname(p) )  # '/home/charly/'

# Returns normalized absolutized version of given path
print( path.abspath(p) )  # '/home/charly/diary.txt'

# Joins path components intelligently,
# especially matching the underlying OS
print( path.join('/home/charly', 'diary.txt') )
```

# Package: os.path

```python
myscript.py
from os import path

# Assume the following file (and path) exists
# /home/charly/diary.txt

# Some very helpful check functions

path.isdir( '/home/charly/diary.txt' )
# False

path.isdir( '/home/charly' )
# True

path.isfile( '/home/charly/diary.txt' )
# True

path.isfile( '/home/charly' )
# False
```

# Package: shutil

- In addition to package os the package **shutil** can be recommended
- High-level operations on files and collections of files

- Functionality to work with files is built-in
- If file is compressed or archived you will need additional modules

myscript.py

```python
1   # Assume the file 'diary.txt' exists
2
3   # Open the specified file
4   # The return value is a file object to work with
5   # The second parameter specifies either
6   #  - 'r': open for reading (default)
7   #  - 'w': open for writing
8   #  - 'a': open for appending
9
10  f = open('diary.txt', 'r')
11
12  # do some work
13  pass
14
15  # Close the file
16  f.close()
```

# Reading and Writing Files

- **All files you open must be closed**
- Otherwise you risk data loss or undesired program behaviour

- Modern Python takes that burden from you: **with** statement
- As soon as the with block is left Python takes care of closing the file
- Even in case of error or unexpected behaviour

```
myscript.py
1  # Assume the file 'diary.txt' exists
2
3  with open("diary.txt", 'r') as f:
4      # do some work
5      pass
```

# Reading and Writing Files

- **All files you open must be closed**
- Otherwise you risk data loss or undesired program behaviour

- Modern Python takes that burden from you: **with** statement
- As soon as the with block is left Python takes care of closing the file
- Even in case of error or unexpected behaviour

myscript.py
```python
1  # Assume the file 'diary.txt' exists
2
3  with open("diary.txt", 'r') as f:
4      # do some work
5      pass
```

Side note:

- The keyword **pass** is a very important one
- What is it doing? Exactly nothing!
- With this cool feature it can serve as a useful placeholder
- Think about in which situations it can be helpful

# Reading and Writing Files

- In Bioinformatics we mainly work with **plain-text** files
- Such files contain human readable characters already
- Human readable: files are no encoded, compressed or binary
- Human readable: not necessarily readable text

# Reading and Writing Files

- In Bioinformatics we mainly work with **plain-text** files
- Such files contain human readable characters already
- Human readable: files are no encoded, compressed or binary
- Human readable: not necessarily readable text

```
1
2      Marvin  01211112152D
3
4     3  2  0  0  0  0              999 V2000
5     -0.4125     0.7145     0.0000 H    0  0  0  0  0  0  0  0  0  0  0  0
6      0.0000     0.0000     0.0000 O    0  0  0  0  0  0  0  0  0  0  0  0
7     -0.4125    -0.7145     0.0000 H    0  0  0  0  0  0  0  0  0  0  0  0
8     2  1  1  0  0  0  0
9     2  3  1  0  0  0  0
10   M  END
```

water.mol

Example: MOL file to store topology of chemical molecules

# Reading and Writing Files

```
1   >sp|P01308|INS_HUMAN Insulin OS=Homo sapiens OX=9606 GN=INS PE=1 SV=1
2   MALWMRLLPLLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAED
3   LQVGQVELGGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN
```

Example: FASTA file to store sequence data

# Reading and Writing Files

```
1   >sp|P01308|INS_HUMAN Insulin OS=Homo sapiens OX=9606 GN=INS PE=1 SV=1
2   MALWMRLLPLLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAED
3   LQVGQVELGGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN
```

Example: FASTA file to store sequence data

- Plain-text files usually consist of one or more lines of 'data'
- Invisible characters indicate line breaks (OS dependent)
- They belong to a set of special and so-called **whitespace** characters
- They can be exploited by programs to read files line by line

# Reading and Writing Files

P01308.fa

```
1    >sp|P01308|INS_HUMAN Insulin OS=Homo sapiens OX=9606 GN=INS PE=1 SV=1
2    MALWMRLLPLLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAED
3    LQVGQVELGGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN
```

myscript.py

```python
1   with open("P01308.fa", 'r') as f:
2       # Function read() consumes the entire file content
3       data = f.read()
4       print(data)
```

```
1   $ python myscript.py
```

# Reading and Writing Files

**P01308.fa**

```
1  >sp|P01308|INS_HUMAN Insulin OS=Homo sapiens OX=9606 GN=INS PE=1 SV=1
2  MALWMRLLPLLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAED
3  LQVGQVELGGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN
```

**myscript.py**

```python
1  with open("P01308.fa", 'r') as f:
2      # Function read() consumes the entire file content
3      data = f.read()
4      print(data)
```

```
1  $ python myscript.py
2  >sp|P01308|INS_HUMAN Insulin OS=Homo sapiens OX=9606 GN=INS PE=1 SV=1
3  MALWMRLLPLLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAED
4  LQVGQVELGGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN
5
6  $
```

# Reading and Writing Files

P01308.fa

```
1    >sp|P01308|INS_HUMAN Insulin OS=Homo sapiens OX=9606 GN=INS PE=1 SV=1
2    MALWMRLLPLLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAED
3    LQVGQVELGGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN
```

myscript.py

```python
with open("P01308.fa", 'r') as f:
    # Function readline() consumes a single line
    data = f.readline()
    print(data)
```

```
$ python myscript.py
```

# Reading and Writing Files

P01308.fa

```
1  >sp|P01308|INS_HUMAN Insulin OS=Homo sapiens OX=9606 GN=INS PE=1 SV=1
2  MALWMRLLPLLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAED
3  LQVGQVELGGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN
```

myscript.py

```python
1  with open("P01308.fa", 'r') as f:
2      # Function readline() consumes a single line
3      data = f.readline()
4      print(data)
```

```
1  $ python myscript.py
2  >sp|P01308|INS_HUMAN Insulin OS=Homo sapiens OX=9606 GN=INS PE=1 SV=1
3
4  $
```

# Reading and Writing Files

P01308.fa

```
1   >sp|P01308|INS_HUMAN Insulin OS=Homo sapiens OX=9606 GN=INS PE=1 SV=1
2   MALWMRLLPLLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAED
3   LQVGQVELGGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN
```

myscript.py

```python
1   with open("P01308.fa", 'r') as f:
2       line = f.readline()
3       while line != '':
4           print(line)
5           line = f.readline()
```

```
1   $ python myscript.py
```

# Reading and Writing Files

```
1  >sp|P01308|INS_HUMAN Insulin OS=Homo sapiens OX=9606 GN=INS PE=1 SV=1
2  MALWMRLLPLLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAED
3  LQVGQVELGGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN
```

myscript.py

```python
1  with open("P01308.fa", 'r') as f:
2      line = f.readline()
3      while line != '':
4          print(line)
5          line = f.readline()
```

```
1  $ python myscript.py
2  >sp|P01308|INS_HUMAN Insulin OS=Homo sapiens OX=9606 GN=INS PE=1 SV=1
3
4  MALWMRLLPLLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAED
5
6  LQVGQVELGGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN
7
8  $
```

# Reading and Writing Files

P01308.fa

```
1  >sp|P01308|INS_HUMAN Insulin OS=Homo sapiens OX=9606 GN=INS PE=1 SV=1
2  MALWMRLLPLLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAED
3  LQVGQVELGGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN
```

myscript.py

```python
1  with open("P01308.fa", 'r') as f:
2      # strip('chars') removes leading and trailing characters
3      # Without parameter it removes whitespace characters
4      line = f.readline().strip()
5      while line != '':
6          print(line)
7          line = f.readline().strip()
```

```
1  $ python myscript.py
```

# Reading and Writing Files

```
1   >sp|P01308|INS_HUMAN Insulin OS=Homo sapiens OX=9606 GN=INS PE=1 SV=1
2   MALWMRLLPLLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAED
3   LQVGQVELGGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN
```

myscript.py

```python
1   with open("P01308.fa", 'r') as f:
2       # strip('chars') removes leading and trailing characters
3       # Without parameter it removes whitespace characters
4       line = f.readline().strip()
5       while line != '':
6           print(line)
7           line = f.readline().strip()
```

```
1   $ python myscript.py
2   >sp|P01308|INS_HUMAN Insulin OS=Homo sapiens OX=9606 GN=INS PE=1 SV=1
3   MALWMRLLPLLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAED
4   LQVGQVELGGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN
5   $
```

15

# Reading and Writing Files

```
myscript.py
1  lines = []
2  with open("P01308.fa", 'r') as f:
3      line = f.readline().strip()
4      while line != '':
5          lines.append(line)
6          line = f.readline().strip()
7
8  with open("P01308-copy.fa", 'w') as f:
9      for line in lines:
10         # The write() method writes a given string to the output file
11         f.write(line)
```

```
$ python myscript.py
$ cat P01308-copy.fa
```

# Reading and Writing Files

myscript.py

```python
lines = []
with open("P01308.fa", 'r') as f:
    line = f.readline().strip()
    while line != '':
        lines.append(line)
        line = f.readline().strip()

with open("P01308-copy.fa", 'w') as f:
    for line in lines:
        # The write() method writes a given string to the output file
        f.write(line)
```

```
$ python myscript.py
$ cat P01308-copy.fa
>sp|P01308|INS_HUMAN Insulin OS=Homo sapiens OX=9606 GN=INS PE=1 SV=1MALWMRLLPLLALLALWGPDPAAAFVN
```

# Reading and Writing Files

```python
myscript.py
1   lines = []
2   with open("P01308.fa", 'r') as f:
3       line = f.readline().strip()
4       while line != '':
5           lines.append(line)
6           line = f.readline().strip()
7
8   with open("P01308-copy.fa", 'w') as f:
9       for line in lines:
10          # write() does not insert a line break by default!
11          f.write(line + '\n')
```

```
1   $ python myscript.py
2   $ cat P01308-copy.fa
```

# Reading and Writing Files

```
myscript.py
1  lines = []
2  with open("P01308.fa", 'r') as f:
3      line = f.readline().strip()
4      while line != '':
5          lines.append(line)
6          line = f.readline().strip()
7
8  with open("P01308-copy.fa", 'w') as f:
9      for line in lines:
10         # write() does not insert a line break by default!
11         f.write(line + '\n')
```

```
1  $ python myscript.py
2  $ cat P01308-copy.fa
3  >sp|P01308|INS_HUMAN Insulin OS=Homo sapiens OX=9606 GN=INS PE=1 SV=1
4  MALWMRLLPLLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAED
5  LQVGQVELGGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN
```

# Reading and Writing Files

```
myscript.py
1  with open("P01308.fa", 'r') as i, open("P01308-copy.fa", 'w') as o:
2      line = i.readline().strip()
3      while line != '':
4          o.write(line + '\n')
5          line = i.readline().strip()
```

A very compact version for text file copying by read-write : )

- I need to apologize … the examples above are really bad style!
- Guess why?

# Command Line Interface

- I need to apologize … the examples above are really bad style!
- Guess why?

```
myscript.py
1  with open("P01308.fa", 'r') as f:
2      data = f.read()
3      print(data)
```

Bad …

# Command Line Interface

- I need to apologize ... the examples above are really bad style!
- Guess why?

myscript.py

```python
1  with open("/home/charly/data/P01308.fa", 'r') as f:
2      data = f.read()
3      print(data)
```

Worse ...

# Command Line Interface

- I need to apologize … the examples above are really bad style!
- Guess why?

myscript.py

```
1   with open("/home/charly/data/P01308.fa", 'r') as f:
2       data = f.read()
3       print(data)
```

Worse …

myscript.py

```
1   with open("../data/P01308.fa", 'r') as f:
2       data = f.read()
3       print(data)
```

Even worse … or at least as worse …

# Command Line Interface

- I need to apologize … the examples above are really bad style!
- Guess why?

```
myscript.py
1  with open("/home/charly/data/P01308.fa", 'r') as f:
2      data = f.read()
3      print(data)
```

Worse …

- Paths written literally in the source code are bad
- We call that practice **hard coding**
- There are several reasons for badness:
  - → We cannot easily use another FASTA input file
  - → Cross-platform incompatible (example code breaks on Win)
  - → Created output files will exist the second run: overwrite

- Good solution: pass file names as arguments to the script
- The standard library package **sys** can help us here

# Command Line Interface: module sys

- Good solution: pass file names as arguments to the script
- The standard library package **sys** can help us here

myscript.py

```python
1  import sys
2
3  input_file = sys.argv[1]
4
5  with open(input_file, 'r') as f:
6      data = f.read()
7      print(data)
```

```
1  $ python myscript.py P01308.fa
```

# Command Line Interface: module sys

- Good solution: pass file names as arguments to the script
- The standard library package **sys** can help us here

myscript.py

```python
import sys

input_file = sys.argv[1]

with open(input_file, 'r') as f:
    data = f.read()
    print(data)
```

```
$ python myscript.py P01308.fa
>sp|P01308|INS_HUMAN Insulin OS=Homo sapiens OX=9606 GN=INS PE=1 SV=1
MALWMRLLPLLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAED
LQVGQVELGGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN
$
```

# Command Line Interface: module sys

```
1  import sys
2
3  print( sys.argv )
```

```
1  $ python myscript.py a b c -d1 -e 2 3
```

# Command Line Interface: module sys

myscript.py

```python
1   import sys
2
3   print( sys.argv )
```

```
$ python myscript.py a b c -d1 -e 2 3
['myscript.py', 'a', 'b', 'c', '-d1', '-e', '2', '3']
$
```

# Command Line Interface: module sys

```
myscript.py
1  import sys
2
3  print( sys.argv )
```

```
1  $ python myscript.py a b c -d1 -e 2 3
2  ['myscript.py', 'a', 'b', 'c', '-d1', '-e', '2', '3']
3  $
```

- This solution is well suited for assignments or simple use cases

- Parsing, interpretation and error handling needs to be done manually

- Gets unhandy with more an especially optional arguments

# Command Line Interface: module argparse

- A better solution for building **command line interfaces** (CLI) exists
- A CLI defines how a program can be called from the command line
- It defines optional and mandatory arguments, flags, files, …

- The standard library module **argparse** is highly recommended
- It provides functionality to define a command line interface
- All kinds of arguments can be predefined and evaluated
- It parses the command line call automatically
- It checks for errors and prints appropriate messages
- It automatically creates and on demand prints a help section
- It allows to access the arguments in a structured way

```
myscript.py
1  import argparse as ap
2
3  parser = ap.ArgumentParser(description='Read and print FASTA files.')
4  parser.add_argument('infile', type=str, help='FASTA input file')
5  parser.add_argument('-s', '--strip',
6                      action='store_true',
7                      help='Strip read lines')
8
9  args = parser.parse_args()
```

```
1  $ python myscript.py
```

# Command Line Interface: module argparse

```
1  import argparse as ap
2
3  parser = ap.ArgumentParser(description='Read and print FASTA files.')
4  parser.add_argument('infile', type=str, help='FASTA input file')
5  parser.add_argument('-s', '--strip',
6                      action='store_true',
7                      help='Strip read lines')
8
9  args = parser.parse_args()
```

```
1  $ python myscript.py
2  usage: myscript.py [-h] [-s] infile
3  myscript.py: error: the following arguments are required: infile
4  $
```

# Command Line Interface: module argparse

```
myscript.py
1  import argparse as ap
2
3  parser = ap.ArgumentParser(description='Read and print FASTA files.')
4  parser.add_argument('infile', type=str, help='FASTA input file')
5  parser.add_argument('-s', '--strip',
6                      action='store_true',
7                      help='Strip read lines')
8
9  args = parser.parse_args()
```

```
1  $ python myscript.py -h
```

# Command Line Interface: module argparse

`myscript.py`

```python
import argparse as ap

parser = ap.ArgumentParser(description='Read and print FASTA files.')
parser.add_argument('infile', type=str, help='FASTA input file')
parser.add_argument('-s', '--strip',
                    action='store_true',
                    help='Strip read lines')

args = parser.parse_args()
```

```
$ python myscript.py -h
usage: myscript.py [-h] [-s] infile

Read and print a FASTA file.

positional arguments:
  infile        FASTA input file

optional arguments:
  -h, --help   show this help message and exit
  -s, --strip  Strip read lines
$
```

# Command Line Interface: module argparse

```
myscript.py
1  import argparse as ap
2
3  parser = ap.ArgumentParser(description='Read and print FASTA files.')
4  parser.add_argument('infile', type=str, help='FASTA input file')
5  parser.add_argument('-s', '--strip',
6                      action='store_true',
7                      help='Strip read lines')
8
9  args = parser.parse_args()
10
11 print( args.infile )
12 print( args.strip )
```

```
1  $ python myscript.py P01308.fa
```

# Command Line Interface: module argparse

myscript.py

```
1  import argparse as ap
2
3  parser = ap.ArgumentParser(description='Read and print FASTA files.')
4  parser.add_argument('infile', type=str, help='FASTA input file')
5  parser.add_argument('-s', '--strip',
6                      action='store_true',
7                      help='Strip read lines')
8
9  args = parser.parse_args()
10
11 print( args.infile )
12 print( args.strip )
```

```
1  $ python myscript.py P01308.fa
2  P01308.fa
3  False
4  $
```

# Command Line Interface: module argparse

myscript.py

```
1  import argparse as ap
2
3  parser = ap.ArgumentParser(description='Read and print FASTA files.')
4  parser.add_argument('infile', type=str, help='FASTA input file')
5  parser.add_argument('-s', '--strip',
6                      action='store_true',
7                      help='Strip read lines')
8
9  args = parser.parse_args()
10
11 print( args.infile )
12 print( args.strip )
```

```
1  $ python myscript.py P01308.fa -s
```

# Command Line Interface: module argparse

```
1  import argparse as ap
2
3  parser = ap.ArgumentParser(description='Read and print FASTA files.')
4  parser.add_argument('infile', type=str, help='FASTA input file')
5  parser.add_argument('-s', '--strip',
6                      action='store_true',
7                      help='Strip read lines')
8
9  args = parser.parse_args()
10
11 print( args.infile )
12 print( args.strip )
```

```
1  $ python myscript.py P01308.fa -s
2  P01308.fa
3  True
4  $
```

# Diving Deeper …

Why closing files is important

→ https://realpython.com/why-close-file-python/

Reading and parsing plain text files in Python

→ https://realpython.com/python-csv/

→ https://realpython.com/python-json/

→ https://docs.python.org/3/library/xml.etree.elementtree.html

More about CLI in Python

→ https://realpython.com/command-line-interfaces-python-argparse/

→ https://docs.python.org/3/library/argparse.html

p001: A third method can be used to read text file content: readlines()
Please find out what it is doing.

p002: There are methods closely relate top strip(): lstrip() and rstrip()
Please familiarize yourself with them and how to use them.

p003: Extract the header lines from a fasta file without symbol '>' using one of the strip methods.
Write the extracted and stripped lines into a new file.

p004: Another very important and helpful string method is split()
Familiarize yourself with it. Apply it on the FASTA header lines to learn how to use it.

p005: Look at the snippet where the newline character '\n' was used to write the FASTA file.
Please try to modify this snippet to work without adding the newline character.

p006: A second method can be used to write content to a text file: writelines()
Please find out what it is doing.

p007: Write a program with the following properties:
- Has 3 functions: getFolderContent(…), analyseFolders(…), writeFolderStats(…)
- Takes 2 input parameters that are parsed with argparse:
  (1) the absolute path to a folder
  (2) the name of an output file to write content
- getFolderContent(…) returns a dict with all items in the specified folder
- analyseFolders(…) annotates in the dict if item is a folder or a directory
- writeFolderStats(…) writes annotated dict to the specified output file sorted by item name

# License and Contributors

Attribution 4.0 International (CC BY 4.0)

*Main Author*
Philipp Thiel

*Additional Contributors*
n.a.