



subprocess

Programming (for biologists)
BIOL 7800

Previously...



You've been writing code to accomplish tasks

New

Classes

Methods

Attributes

Functions

```
code_example.py — /Users/bcf/Dropbox/Classes/BIOL7800/temp/lecture20

code_example.py
1  #!/usr/bin/env python
2  # -*- coding: utf-8 -*-
3
4  """
5  (c) 2016 Brant Faircloth || http://faircloth-lab.org/
6  All rights reserved.
7  This code is distributed under a 3-clause BSD license. Please see
8  LICENSE.txt for more information.
9  Created on 17 March 2016 16:01 CDT (-0500)
10 """
11
12 import re
13 import argparse
14 from collections import Counter
15
16 # import pdb
17
18
19 def get_args():
20     """Get arguments from CLI"""
21     parser = argparse.ArgumentParser(
22         description="Reads a file in; write some stuff out"
23     )
24     parser.add_argument(
25         "--infile",
26         required=True,
27         type=str,
28         help="The path to the input file"
29     )
30     parser.add_argument(
31         "--output",
32         required=True,
33         type=str,
34         help="The path to the output file"
35     )
36     return parser.parse_args()
37
38
39 def get_and_clean_file(f):
40     with open(f, 'r') as infile:
41         text = infile.read()
42     return re.findall("\w+[-]*\w+", text.lower())
43
44
45 def pretty_print_counts(cnt, top=20):
46     for word, cnt in cnt.most_common(top):
47         print("{: <20}{:}".format(word, cnt))
48
49
50 def pretty_write_counts(o, cnt):
51     with open(o, 'w') as outfile:
52         # don't need to sort here - most_common gives us sorted output.
53         for word, cnt in cnt.most_common():
54             outfile.write("{}\t{}\n".format(word, cnt))
55
56
57 def main():
58     args = get_args()
59     clean_text = get_and_clean_file(args.infile)
60     cnt = Counter(clean_text)
61     pretty_print_counts(cnt)
62     pretty_write_counts(args.output, cnt)
63
64
65 if __name__ == '__main__':
66     main()
```

File 0 Project 0 ✓ No Issues code_example.py 66:11 LF UTF-8 Python

But often...

Programs **already exist** that do something you want to do

Shell commands

tail
head
wc
cat

Sequence assembly

velvet
trinity
abYss

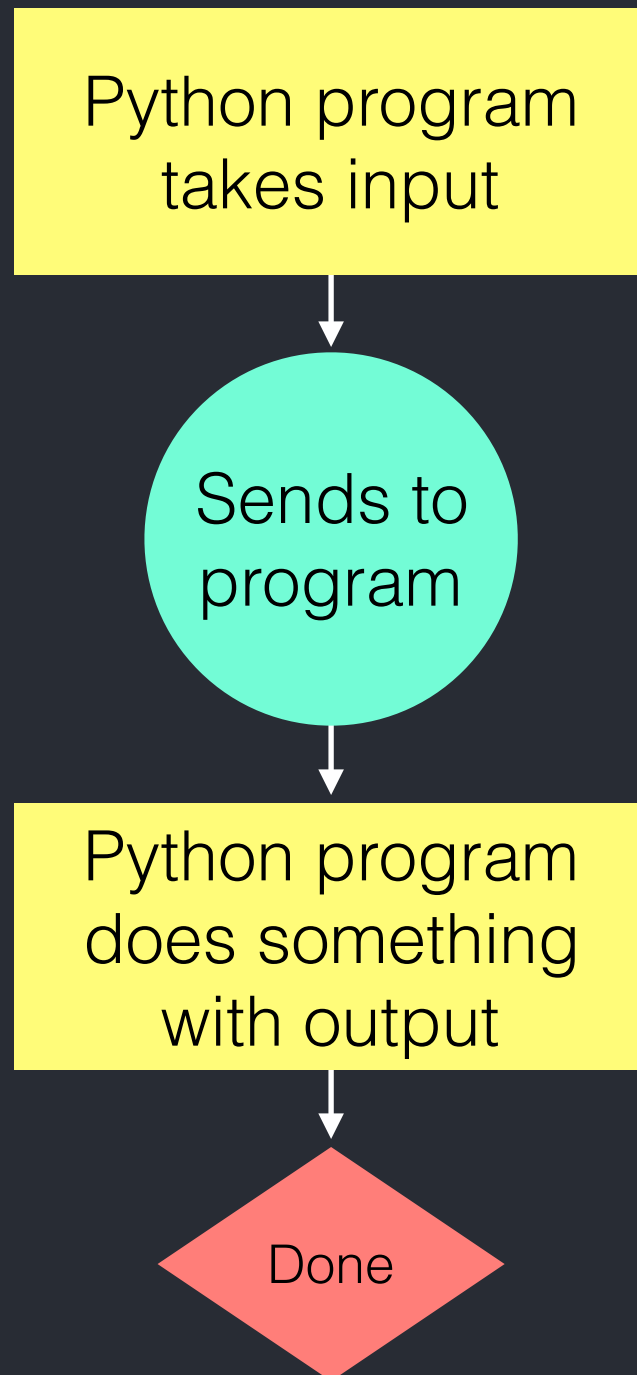
Alignment

bwa
bowtie
lastz
blast
mafft
muscle

Or, any other external program you want to run
(from within your Python code)

AKA Pipelining

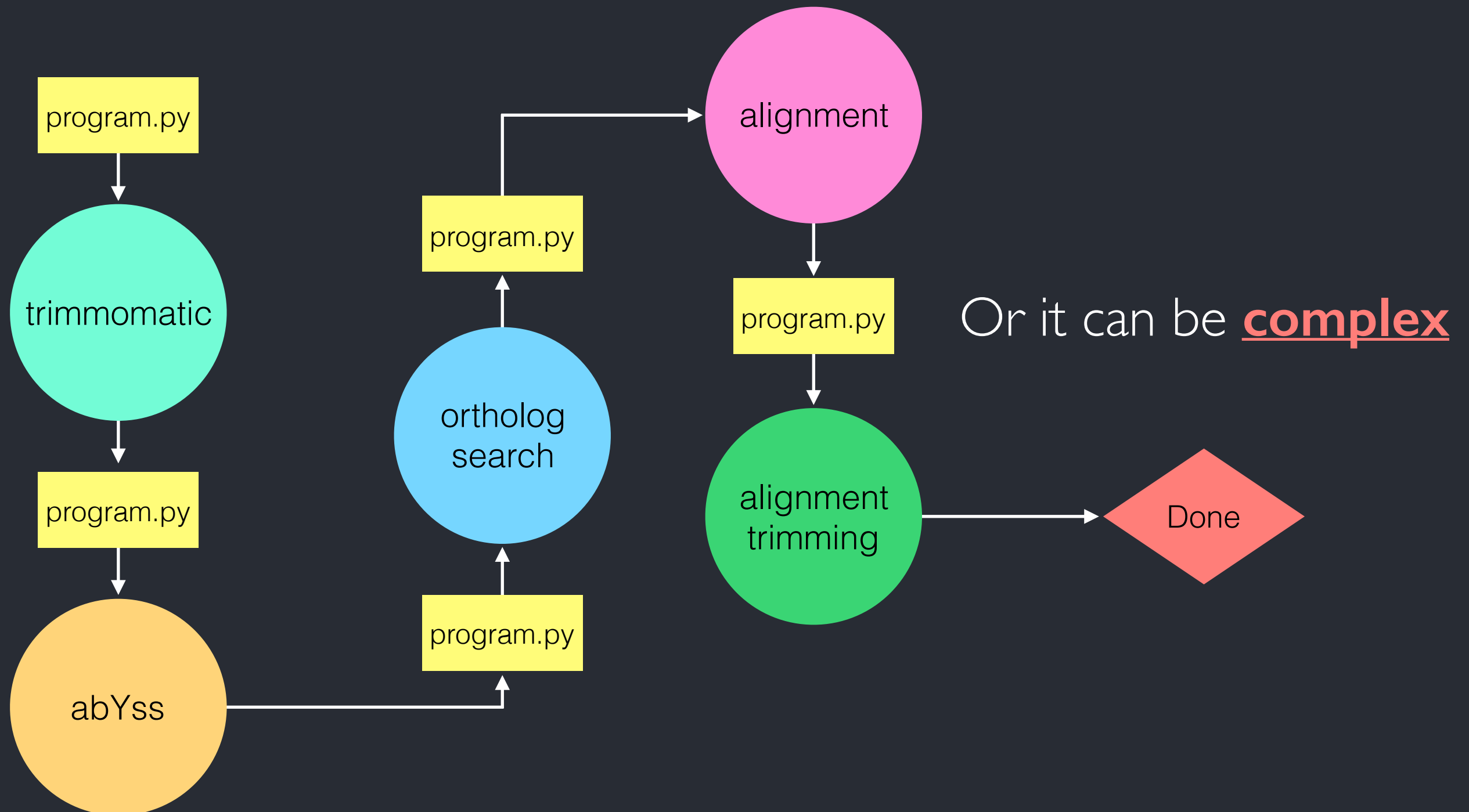
(or building a pipeline)



This can be simple

AKA Pipelining

(or building a pipeline)



But, how do you do all of this?

subprocess

A python (standard) module that allows you to

"spawn new processes, connect to their input/output/error pipes and obtain their return codes"

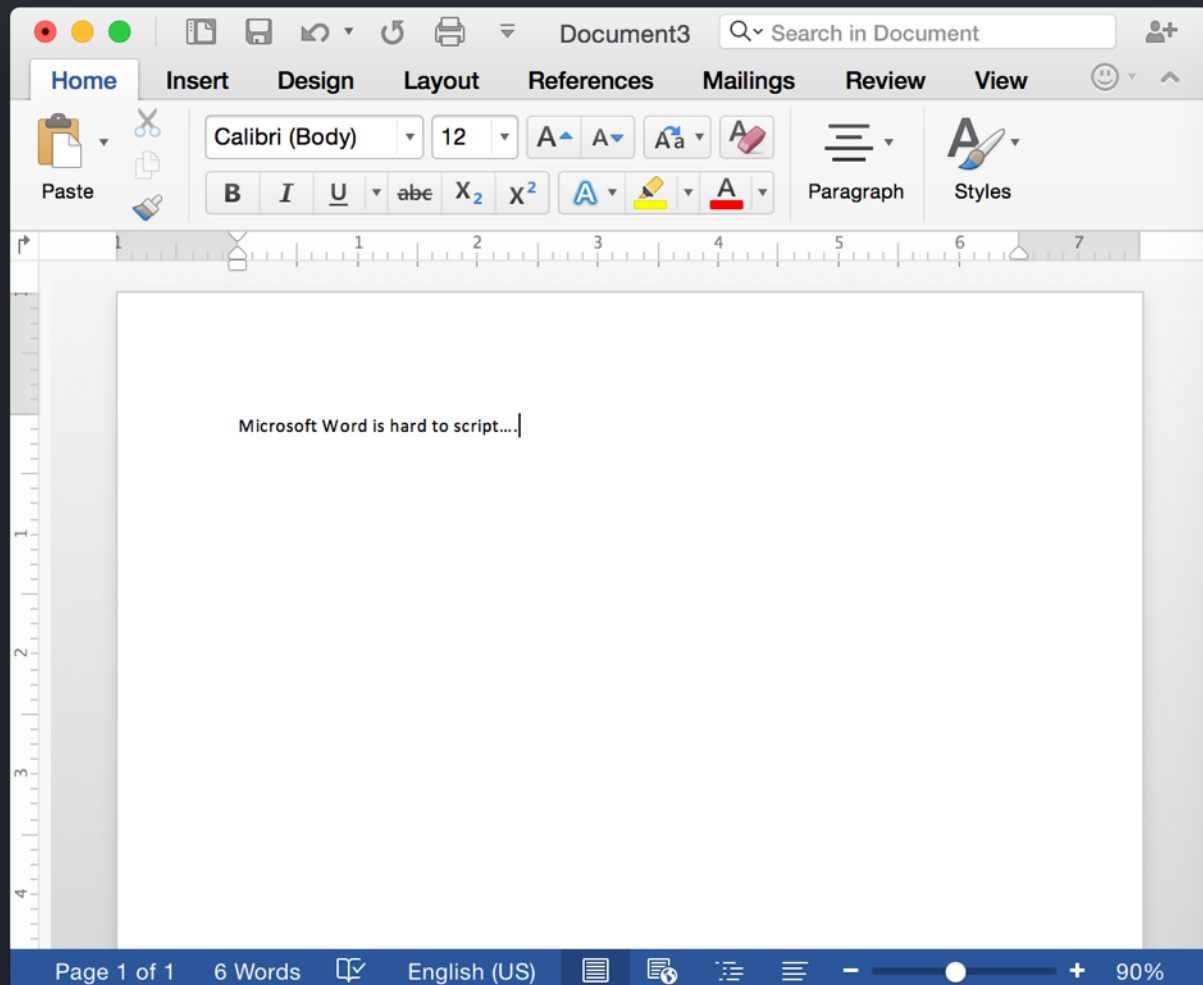
- python 3.5 doc

subprocess

What does **subprocess** require to make this magic happen?

Basically, that you have a "scriptable" program
you want to run from python

Not so much GUI



More like command-line

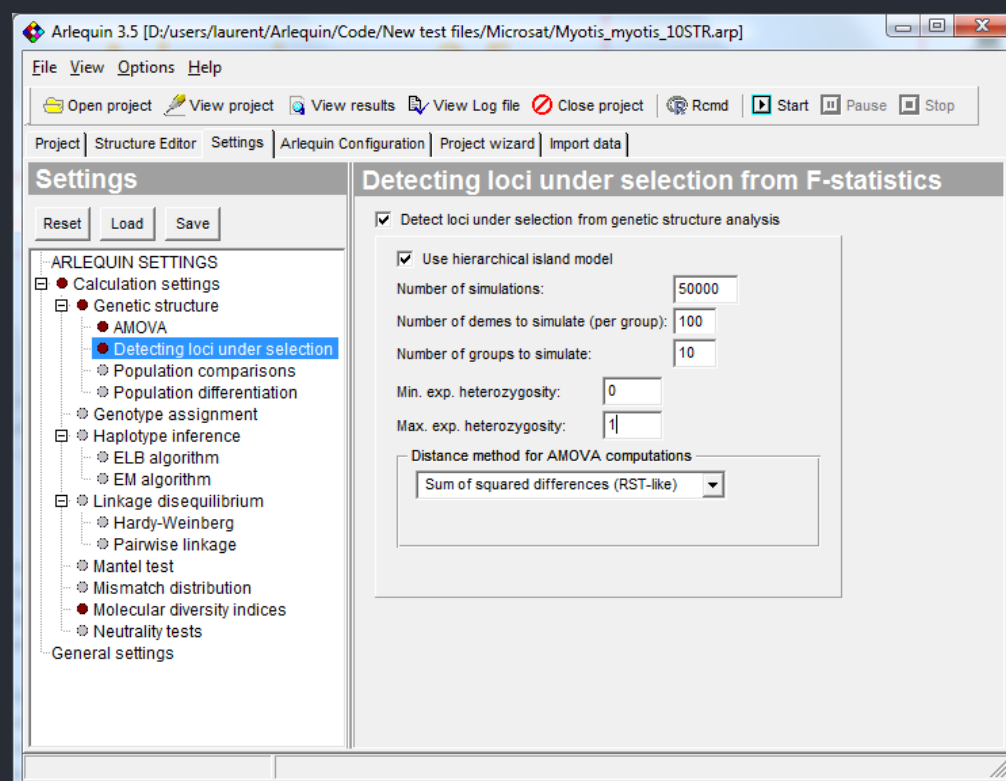
A screenshot of a terminal window titled '1. exit (zsh)'. It shows a command-line session where a script named 'configureBclToFastq.pl' is being executed. The command is: `~/bin/bcl2fastq-1.8.4/bin/configureBclToFastq.pl \`. The script takes several arguments: `--input-dir /nfs/data1/illumina-raw/junco-sparrow/Data/Intensities/L001/ \`, `--output junco-data \`, `--sample-sheet SampleSheet.csv \`, and `--use-bases-mask Y*,I*,Y80`. The prompt indicates the user is 'bcf' at 'brant-4' in the '~' directory.

subprocess

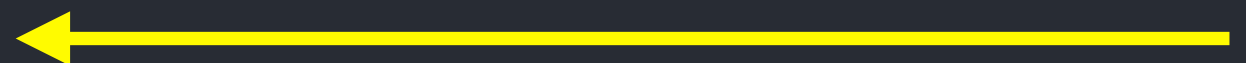
What does **subprocess** require to make this magic happen?

Basically, that you have a "scriptable" program
you want to run from python

Generally, that means a command-line program
but not always



A good example is Arlequin



Has both GUI and command-line interaction

anatomy of subprocess

In: `import subprocess`

← Import `subprocess`

`run()` method is responsible
for running external program

program (and arguments)
goes in a `list`

In: `subprocess.run(['ls'])`

← Use `run` to run prog

Out: `CompletedProcess(args=['ls'], returncode=0)`

↑
`subprocess` returns a
`CompletedProcess` object

↖
The `returncode` is programming lingo for
`success = 0`
`failure = 1`

anatomy of subprocess

```
In: import subprocess
```

```
In: my_process = subprocess.run(['ls'])
```

```
In: dir(my_process)
```

```
Out:
```

```
'args',  
'check_returncode',  
'returncode',  
'stderr',  
'stdout'
```

```
In: my_process.args
```

```
Out: ['ls']
```

```
In: my_process.check_returncode()
```

```
Out:
```

```
In: print(my_process.check_returncode())
```

```
Out: None
```

We can store
CompletedProcess
object of subprocess in a
variable



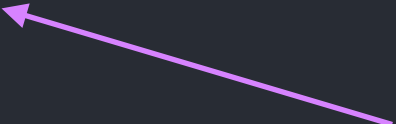

These are the
CompletedProcess
attributes & methods



We can access the **.args**
attribute



We can use the
check_returncode()
method to make sure
program ran



check_returncode() returns
None if program ran successfully

anatomy of subprocess

In: `import subprocess`

These arguments are **wrong**
(they don't work)

In: `my_process = subprocess.run(['ls', '-ZZZ'])`

In: `my_process.args`

Out: `['ls', '-ZZZ']`

← Use **run** to run prog
(with bad arguments)

In: `my_process.returncode`

Out: `1`

In: `my_process.check_returncode()`

Out:

← Let's look at actual
returncode

← Now, let's **check** it...

```
CalledProcessError                                Traceback (most recent call last)
<ipython-input-21-c11267dce8c6> in <module>()
----> 1 my_process.check_returncode()

/Users/bcf/anaconda/envs/py35/lib/python3.5/subprocess.py in check_returncode(self)
    660         if self.returncode:
    661             raise CalledProcessError(self.returncode, self.args, self.stdout,
--> 662                                     self.stderr)
    663
    664
```

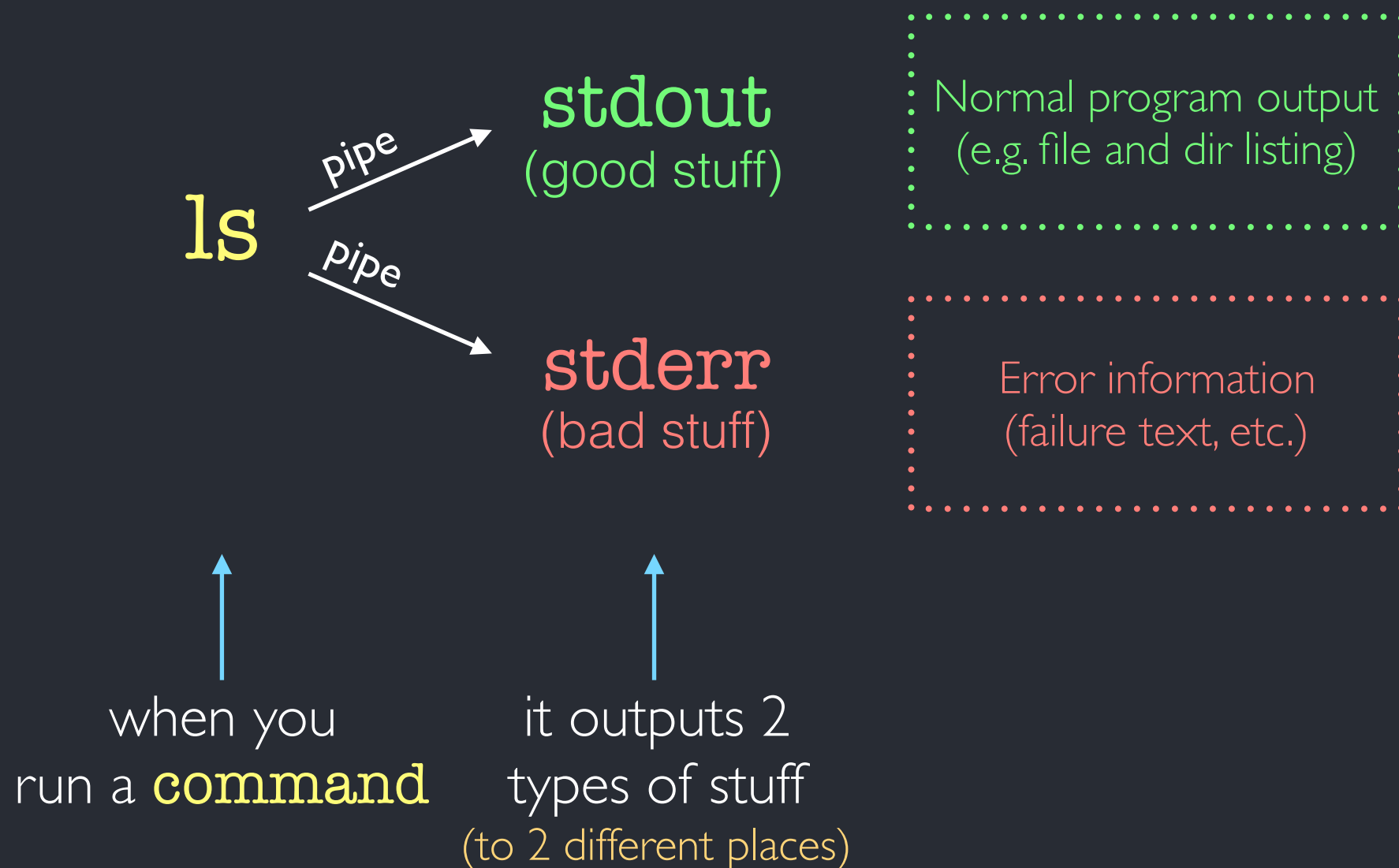
CalledProcessError: Command '['ls', '-ZZZ']' returned non-zero exit status 1

← **Error !**

anatomy of subprocess

okay, that's great, but where's our output?

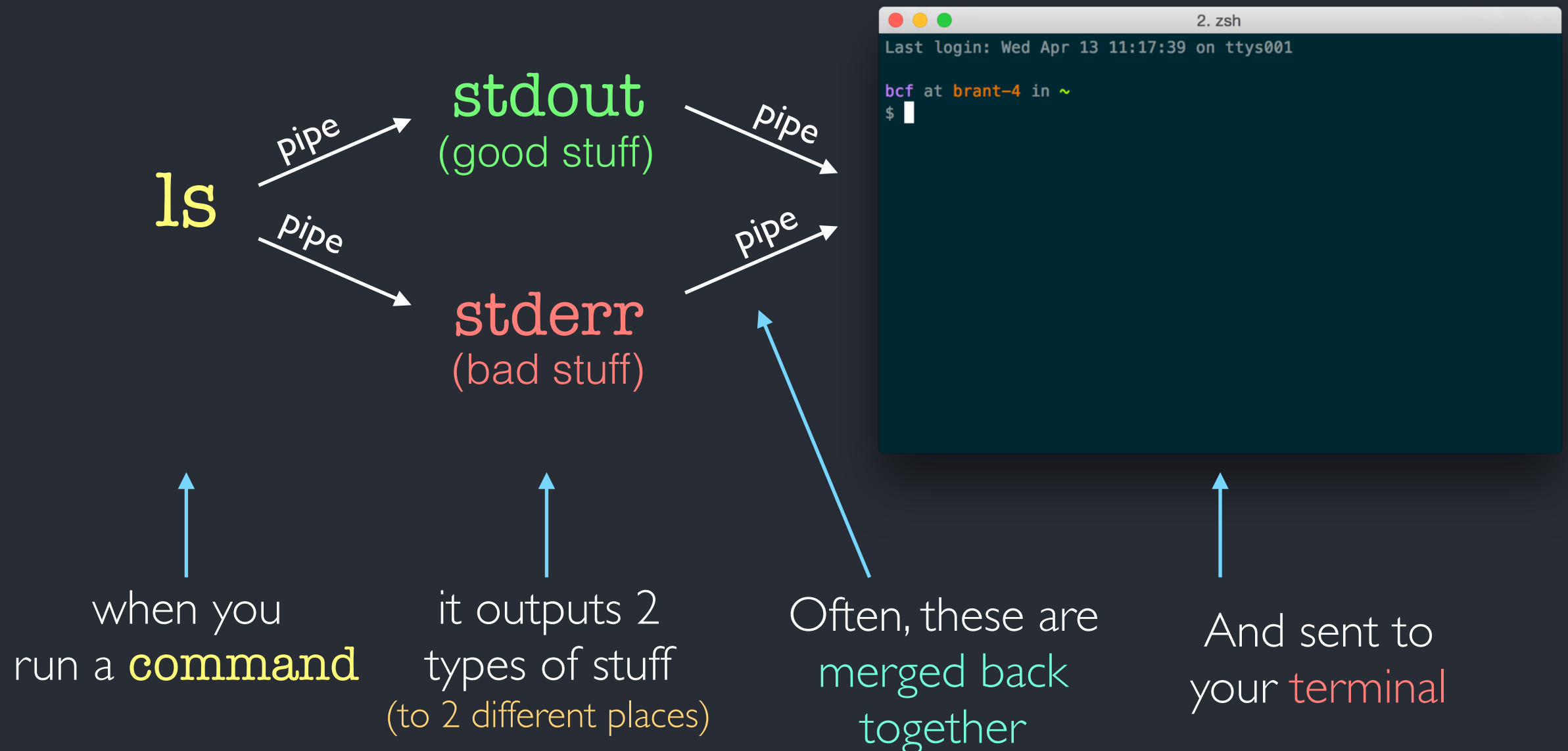
(a quick diversion into the **land of pipes**)



anatomy of subprocess

okay, that's great, but where's our output?

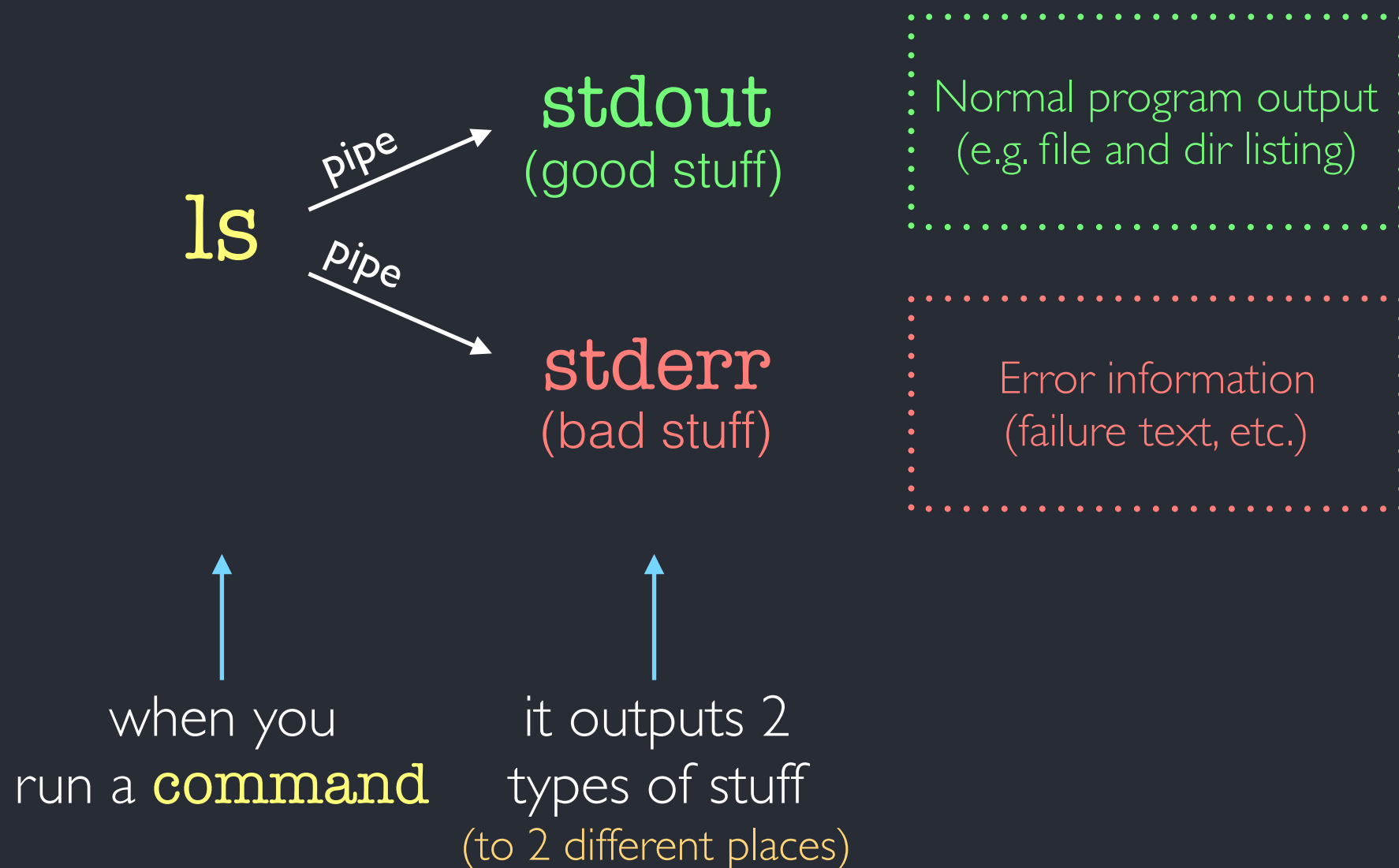
When you run a command from your shell



anatomy of subprocess

When you run a command from subprocess

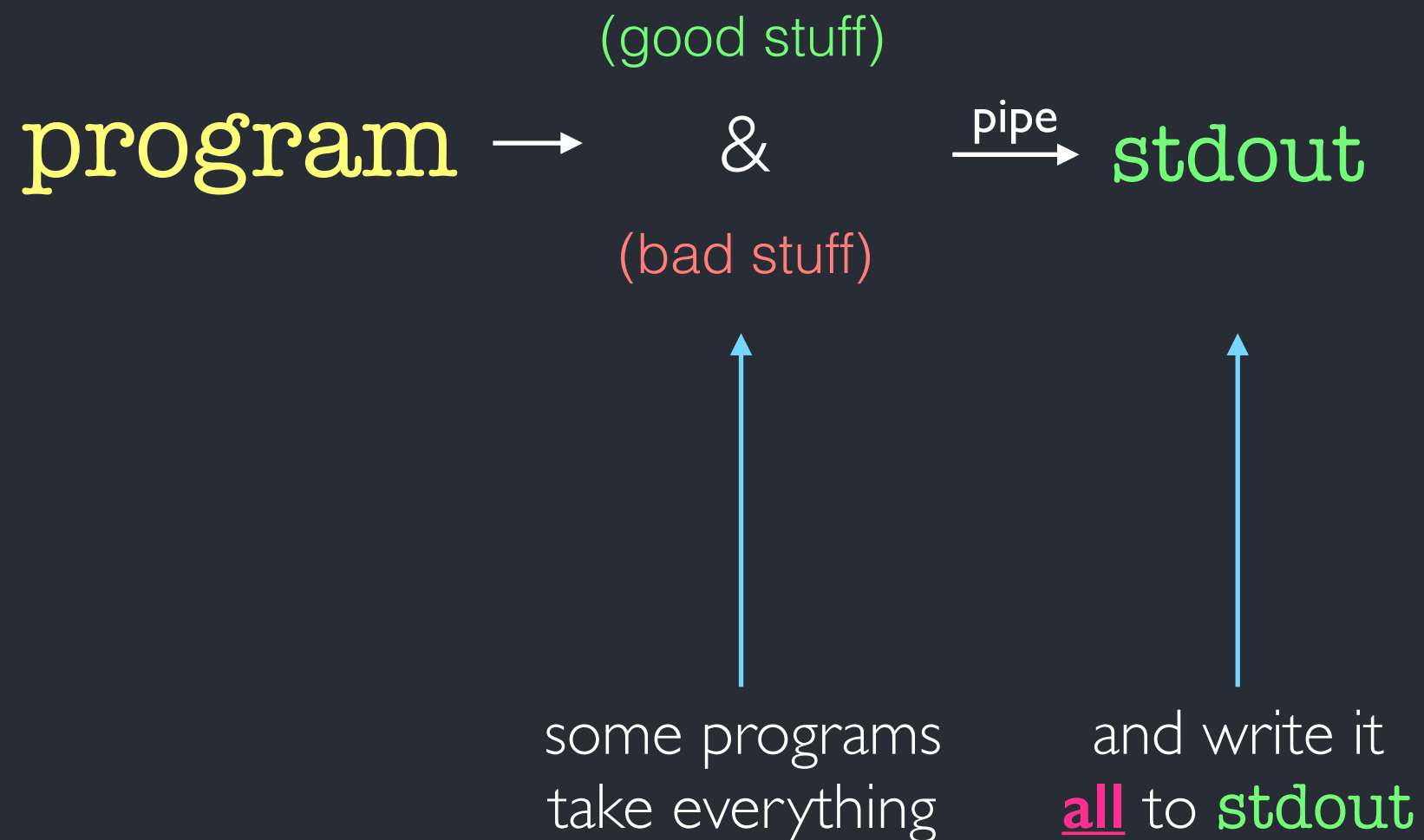
The ideal situation



anatomy of subprocess

When you run a command from subprocess

But it's not always ideal...



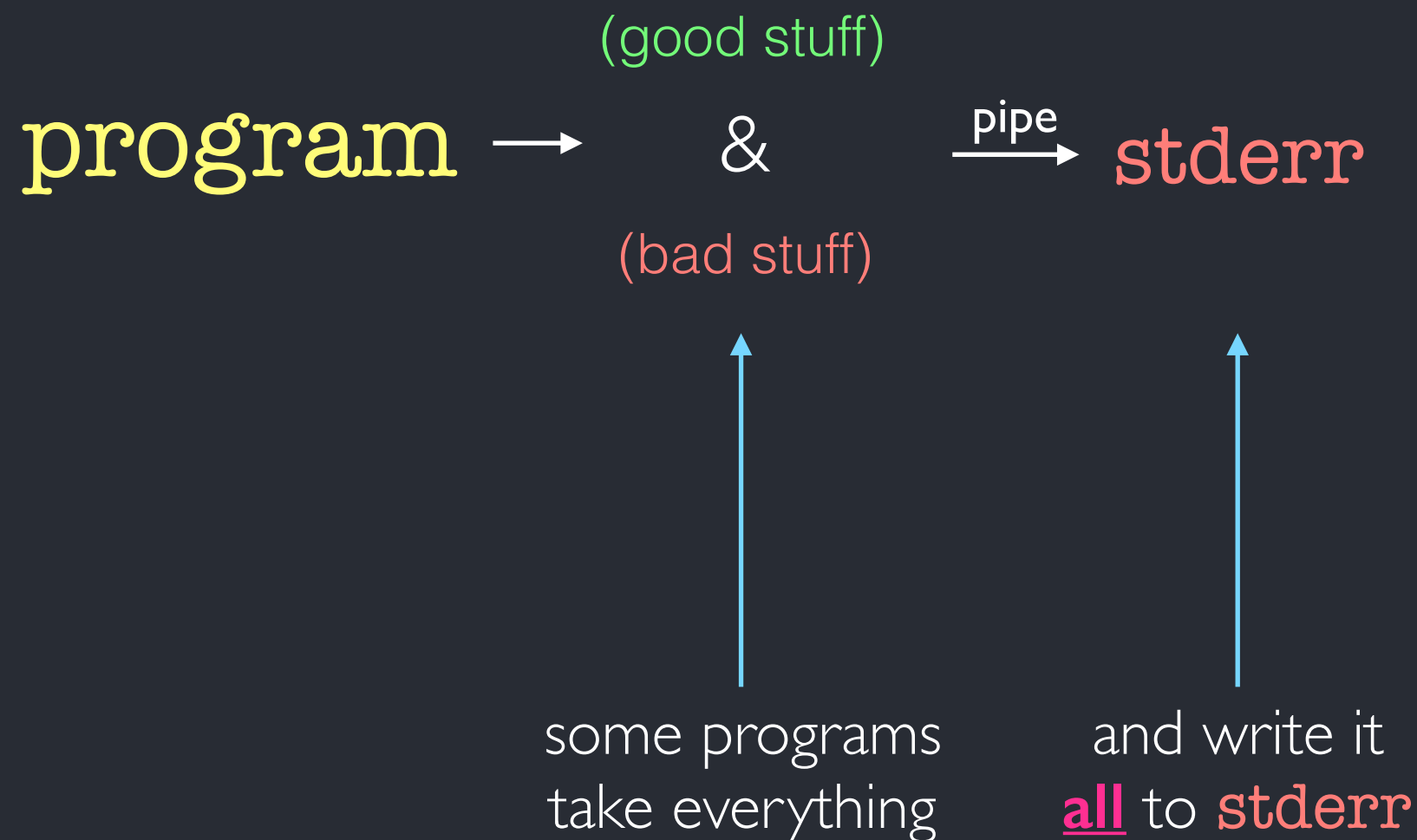
Normal program output
(e.g. file and dir listing)

Error information
(failure text, etc.)

anatomy of subprocess

When you run a command from subprocess

But it's not always ideal...

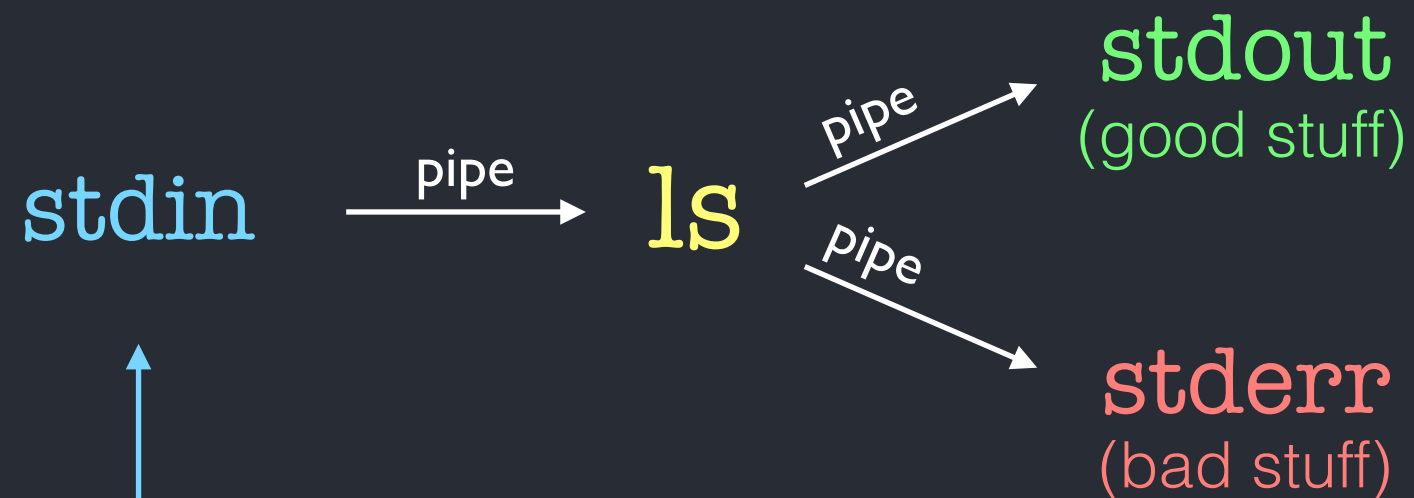


Normal program output
(e.g. file and dir listing)

Error information
(failure text, etc.)

anatomy of subprocess

One more little wrinkle...



Normal program output
(e.g. file and dir listing)

Error information
(failure text, etc.)

you can also **send** information to
some programs on **stdin**

anatomy of subprocess

So, *generally*, we have 3 types of things to worry about



And, the *stdout* is sometimes different from what the *program* can write to a file (although we can *also* store *stdout* and *stdin* in a file)

anatomy of subprocess

In: `import subprocess`


program (and arguments)
goes in a `list`



set keyword stdout
to `subprocess.PIPE`



We're only
capturing `stdout`



In: `my_proc = subprocess.run(['ls'], stdout=subprocess.PIPE)`

In: `my_proc.args`

Out: `['ls']`

In: `my_proc.returncode`

Out: `0`

← See if program ran
`successfully`

In: `my_proc.stdout`

← stdout is in the `stdout`
attribute

Out: `b'Applications\nCreative Cloud Files\nDesktop\nDocuments\nDownloads\nDropbox\nDropbox (faircloth-lab)\nLibrary\nMesquite_Support_Files\nMovies\nMusic\nPictures\nPublic\nVirtual Machines\nanaconda\nbin\ngit\ngithub-app\nnotebooks\nsrc\ntmp\n'`

← stdout is returned as a
`bytes` object, which is not
the same as a string

(more in a minute)

anatomy of subprocess

```
In: my_proc = subprocess.run(['ls'],  
                             stdout=subprocess.PIPE,  
                             stderr=subprocess.PIPE)
```

← Capture **stdout** and **stderr**

```
In: my_proc.returncode
```

```
Out: 0
```

← See if program ran **successfully**
(it did)

```
In: my_proc.stdout
```

```
Out: b'Applications\nCreative (...truncated...)'
```

← stdout is in the **stdout**
attribute

```
In: my_proc.stderr
```

```
Out: b''
```

← There is nothing in **stderr**
attribute
(program ran successfully)

anatomy of subprocess

These arguments are **wrong**
(they don't work)



```
In: my_proc = subprocess.run(['ls', '-ZZZ'],  
                             stdout=subprocess.PIPE,  
                             stderr=subprocess.PIPE)
```



Capture **stdout** and **stderr**

```
In: my_proc.returncode
```

```
Out: 1
```



See if program ran **successfully**
(it didn't)

```
In: my_proc.stdout
```

```
Out: b''
```



There is nothing in **stdout**
(didn't run)

```
In: my_proc.stderr
```

```
Out: b'ls: illegal option -- Z\nusage: ls [-  
ABCFGHL0PRSTUwabcdefghiklmnopqrstuw1] [file ...]\n'
```



But there is something in **stderr**
(because the program failed)

anatomy of subprocess

Can we do something with the `stdout`?

```
In: my_proc = subprocess.run(['ls'],  
                             stdout=subprocess.PIPE,  
                             stderr=subprocess.PIPE)
```

← Capture `stdout` and `stderr`

```
In: my_proc.stdout
```

```
Out: b'Applications\nCreative Cloud Files\nDesktop\nDocuments\nDownloads\nDropbox\nDropbox (faircloth-lab)\nLibrary\nMesquite_Support_Files\nMovies\nMusic\nPictures\nPublic\nVirtual Machines\nanaconda\nbin\ngit\ngithub-app\nnotebooks\nsrc\ntmp\n'
```

← `stdout` is returned as a `bytes` object, which is not the same as a string

```
In: type(my_proc.stdout)
```

```
Out: bytes
```

← `stdout` is returned as a `bytes` object, which is not the same as a string

`bytes` objects are a `primitive` way to encode a string

anatomy of subprocess

Can we do something with the `stdout`?

```
In: my_proc = subprocess.run(['ls'],  
                             stdout=subprocess.PIPE,  
                             stderr=subprocess.PIPE)
```

← Capture `stdout` and `stderr`

```
In: my_proc.stdout
```

```
Out: b'Applications\nCreative Cloud Files\nDesktop\nDocuments\nDownloads\nDropbox\nDropbox (faircloth-lab)\nLibrary\nMesquite_Support_Files\nMovies\nMusic\nPictures\nPublic\nVirtual Machines\nanaconda\nbin\ngit\ngithub-app\nnotebooks\nsrc\ntmp\n'
```

← `stdout` is returned as a `bytes` object, which is not the same as a `string`

```
In: my_proc.stdout.split("\n")
```

```
Out: TypeError: a bytes-like object is required, not 'str'
```

And `bytes` objects need to be *encoded* to `string` objects

← This fails, because `my_proc.stdout` is a `bytes` object, not a string

anatomy of subprocess

Two ways around the problem

Method I

```
In: my_proc = subprocess.run(['ls'],  
                             stdout=subprocess.PIPE,  
                             stderr=subprocess.PIPE)
```

```
In: my_proc.stdout.decode("utf-8").split("\n")
```

← Use the **bytes** object's **decode()** method to convert it to "utf-8" (or other) encoding

```
Out:  
['Applications',  
 'Creative Cloud Files',  
 'Desktop',  
 'Documents',  
 'Downloads',  
 'Dropbox',  
 'Dropbox (faircloth-lab)',  
 'Library',  
 'Mesquite_Support_Files',  
 'Movies',  
 (...truncated...)  
 '']
```


anatomy of subprocess

Two ways around the problem

Method 2

```
In: my_proc = subprocess.run(['ls'],  
                             stdout=subprocess.PIPE,  
                             stderr=subprocess.PIPE,  
                             universal_newlines=True)
```

Pass the `universal_newlines`
argument to
← `subprocess.run()`

```
In: type(my_proc.stdout)  
Out: str
```

← `stdout` is now a `string`

```
In: my_proc.stdout.split("\n")  
Out:
```

← And, we can easily
split a `string`

```
['Applications',  
 'Creative Cloud Files',  
 'Desktop',  
 'Documents',  
 'Downloads',  
 'Dropbox',  
 (...truncated...)  
 '']
```

anatomy of subprocess

Pipes are nice, but we may want to write `stdout/stderr` to a file (particularly when there is a lot of it - our pipes can get full)

```
In: with open('stdout.txt', 'w') as stdout_file:
    with open('stderr.txt', 'w') as stderr_file:
        my_proc = subprocess.run(['ls'],
                                stdout=stdout_file,
                                stderr=stderr_file,
                                universal_newlines=True)
```

↖ ↗ Open two files, one each to capture `stdout` and `stderr`

```
In: print(open('stdout.txt', 'r').read())
```

Out:

```
Applications
Creative Cloud Files
Desktop
Documents
Downloads
Dropbox
(...truncated...)
```

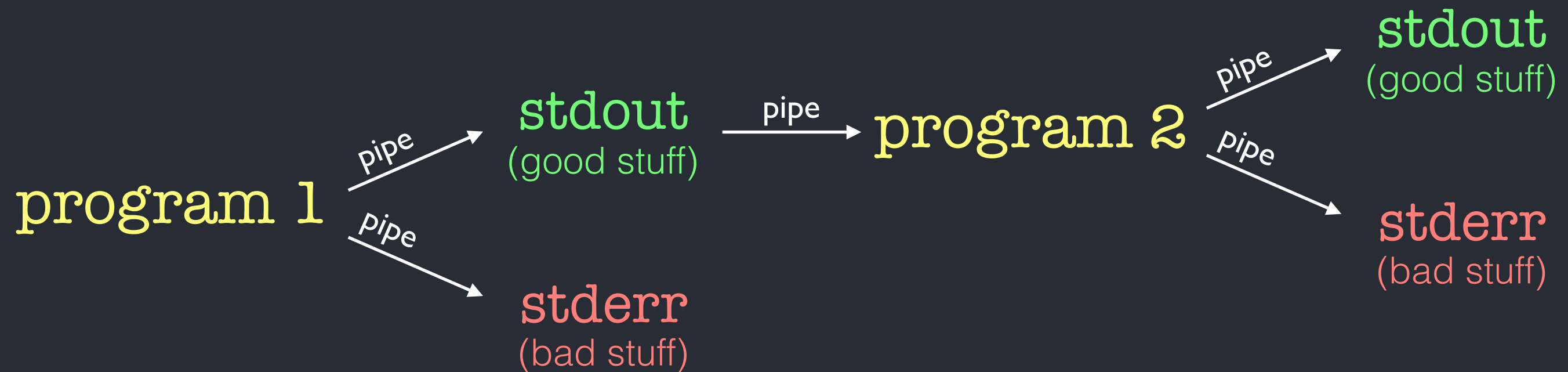
← The lines that were going to our `stdout` pipe are now being written to a file

Same for `stderr`
(but there's nothing there)

We can now do `anything we want` with these file (parse them, etc.)

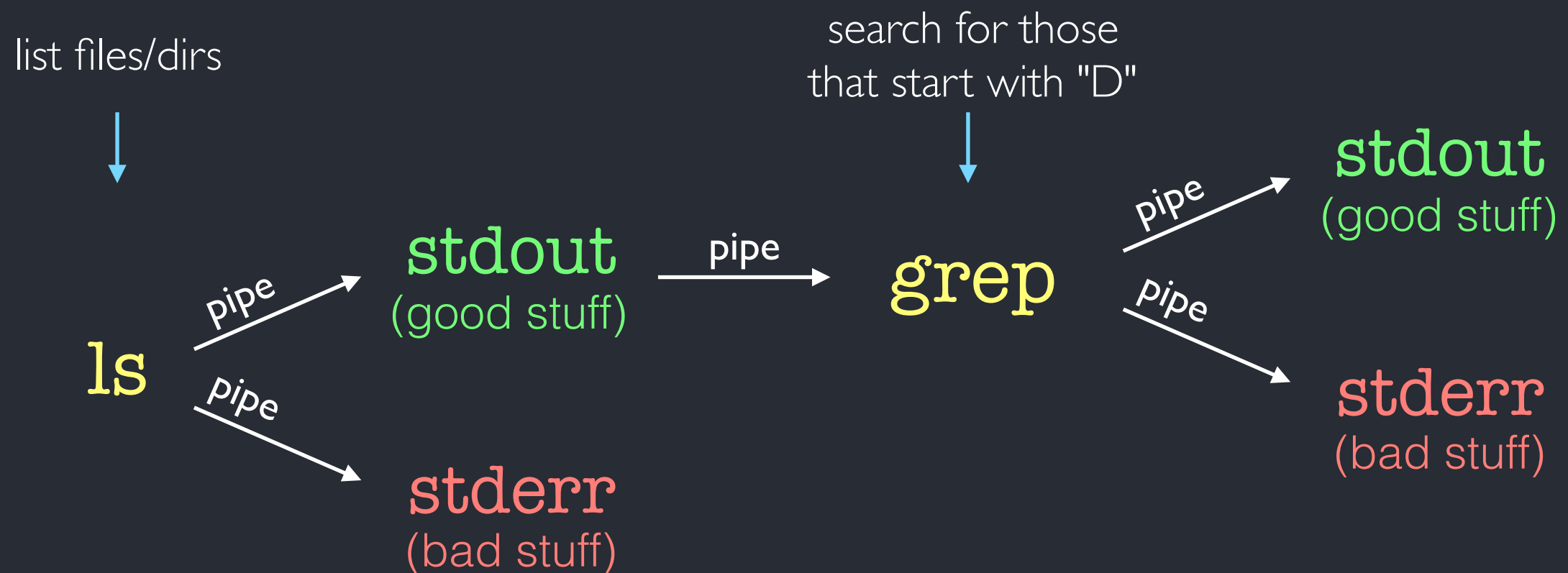
anatomy of subprocess

What if we want to chain several programs together ?



anatomy of subprocess

What if we want to chain several programs together ?



or, as we would type it in the shell

```
ls | grep "^D"
```

anatomy of subprocess

Chaining programs together by pipes

```
In: my_proc_1 = subprocess.run(['ls'],  
                                stdout=subprocess.PIPE,  
                                stderr=subprocess.PIPE)
```

← Run **ls**, send the output to **stdout**

```
my_proc_2 = subprocess.run(['grep', '^D'],  
                            input=my_proc_1.stdout,  
                            stdout=subprocess.PIPE,  
                            stderr=subprocess.PIPE)
```

Take the **my_proc_1.stdout** as **input**, and run **grep** against it

```
print(my_proc_2.stdout)
```

```
Out: b'Desktop\nDocuments\nDownloads\nDropbox\nDropbox (faircloth-lab)\n'
```

stdout is returned as a **bytes** object

anatomy of subprocess


Chaining programs together by pipes (with **string** output)

```
In: my_proc_1 = subprocess.run(['ls'],  
                                stdout=subprocess.PIPE,  
                                stderr=subprocess.PIPE,  
                                universal_newlines=True)
```

```
my_proc_2 = subprocess.run(['grep', '^D'],  
                            input=my_proc_1.stdout,  
                            stdout=subprocess.PIPE,  
                            stderr=subprocess.PIPE,  
                            universal_newlines=True)
```

```
print(my_proc_2.stdout)
```

We need to specify **universal_newlines** in both processes for **string** output



Out:

```
Desktop  
Documents  
Downloads  
Dropbox  
Dropbox (faircloth-lab)
```

stdout is returned as a **string**



anatomy of subprocess

What about programs with **lots** of **parameters**?
e.g.

```
java -jar /usr/bin/gatk -T UnifiedGenotyper -nt 12 -R genome.fasta -I genome.bam -gt_mode DISCOVERY -glm INDEL
```

```
In: command = ['java', '-jar', '/usr/bin/gatk', '-T',  
'UnifiedGenotyper', '-nt', '12', '-R', 'genome.fasta', '-I',  
'genome.bam', '-gt_mode', 'DISCOVERY', '-glm', 'INDEL']
```

Because this is just a **list**, you can make the list and place it in a variable...

```
In: my_proc_1 = subprocess.run(  
    command,  
    stdout=subprocess.PIPE,  
    stderr=subprocess.PIPE,  
    universal_newlines=True)
```

Then, pass the variable to **subprocess.run**

anatomy of subprocess

What about programs with **lots** of **parameters**?
e.g.

```
java -jar /usr/bin/gatk -T UnifiedGenotyper -nt 12 -R genome.fasta -I genome.bam -gt_mode DISCOVERY -glm INDEL
```

In: `import shlex`

In: `shlex.split('java -jar /usr/bin/gatk -T UnifiedGenotyper
-nt 12 -R genome.fasta -I genome.bam -gt_mode DISCOVERY -glm
INDEL')`

Out:

```
['java',  
'-jar',  
'/usr/bin/gatk',  
'-T',  
'UnifiedGenotyper',  
'-nt',  
'12',  
'-R',  
'genome.fasta',  
'-I',  
'genome.bam',  
'-gt_mode',  
'DISCOVERY',  
'-glm',  
'INDEL']
```



You can also use
the **shlex** module
to help you split
complex command
line statements



Then, use the list that
shlex makes as your
input