Speed, filming, and

multiprocessing

Programming (for biologists) BIOL 7800



We've only talked about making programs functions, not making programs fast...

...and we've run only 1 process at a time





We've not talked at all about using multi-CPU resources to run our jobs.





is tricky thing... and can be achieved in a variety of ways

Write more efficient code

Use more CPU resources





Making things faster We're going to focus on the <u>following code</u>

| | example1 | .py - /Users/bcf/Dropbox/Classes/BIC | DL7800/temp/lecture21 |
|---|---------------------------------|--------------------------------------|-----------------------|
| | example1.py | example2.py | example3.py |
| 1 | #! /usr/bin/env py | thon | |
| | <pre># encoding UTF-8</pre> | | |
| | | | |
| | import random | | |
| | | | |
| | | | |
| | <pre>def sum_lots_of_rail</pre> | ndom_nums(): | |
| | sum_holder = [|] | |
| | for elem in ra | nge(100): | |
| | random_num | bers = [] | |
| | for cnt in | range(1000): | |
| | # pick | 1 random number | () |
| | random <u>.</u> | _number = random.randrang | e(0, 100000) |
| | # sum · | those | |
| | random_ | _numbers.append(random_nu | mber) |
| | # sum the . | 1000 random numbers | - \ \ |
| | sum_nolder | append(sum(random_number | 5)) |
| | # SUM THE 100 : | sums ot random numbers | |
| | s = sum(sum_no | ice (,)" format(c)) | |
| | | 15: (:,) .TOTMat(5)) | |
| | | | |
| | def main(). | | |
| | sum lots of ra | ndom nums() | |
| | 3um_cocs_or_ra | | |
| | if name == ' | main ': | |
| | | | |

Does the following 100 times Does the following 1000 times Selects 1 random numbers Adds that to a list Sums the list of 1000 numbers Adds that to a list Sums the list of 100 sums Pretty prints the output

Example

In: python example1.py Out: The sum is: 5,001,095,443

Python provides the timeit module for testing speed

You can use it several ways... the easiest is probably



if nome ___ I main I.

| | example1 | .py - /Users/bcf/Dropbox/Classes/Bl | OL7800/temp/lecture21 |
|---|-------------------------------|-------------------------------------|-----------------------|
| | example1.py | example2.py | example3.py |
| 1 | #! /usr/bin/env py | thon | |
| | <pre># encoding UTF-8</pre> | | |
| | | | |
| | import random | | |
| | | | |
| | | | |
| | <pre>def sum_lots_of_ra</pre> | ndom_nums(): | |
| | sum_holder = [| | |
| | for elem in ra | nge(100): | |
| | random_num | bers = [] | |
| | | 1 random number | |
| | random | number = random randrand | e(0, 100000) |
| | # sum | those | |
| | random | numbers.append(random nu | mber) |
| | # sum the | 1000 random numbers | |
| | sum_holder | .append(sum(random_number | ·s)) |
| | # sum the 100 | sums of random numbers | |
| | s = sum(sum_ho | lder) | |
| | print("The sum | is: {:,}".format(s)) | |
| | | | |
| | | | |
| | <pre>def main():</pre> | | |
| | sum_lots_of_ra | ndom_nums() | |
| | | | |
| | 1fname == 'i | main': | |
| | main() | | |
| | | | |
| | | | |

How can we make this faster?

(without using multiple CPUs) Increase efficiency!!

raw times: 0.601 0.6 0.598
5 loops, best of 3: 120 msec per loop

File 0 Project 0 🗸 No Issues /Users/bcf/Dropbox/Classes/BIOL7800/temp/lecture23/example1.py 1:1 LF UTF-8 Python 💕

Approach #I: Big list

| | | | | | time t | o beat | |
|--------|--------------------------------------|-------------------------------------|--------------------------|-------------|--------------------------------|----------|----------|
| | example2.py - | /Users/bcf/Dropbox/Classes/BIO | _7800/temp/lecture21 | | raw times: 0.601 0. | 6 0.598 | |
| 1 | example1.py | example2.py | example3.py | | 5 loops best of 3 | 120 mcec | ner loon |
| | #! /usr/bin/env python | | | | | | ρει τουρ |
| | # encourng UIF-8 | | | | | | |
| | import random | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | def cum lots of random | numc()• | | | | | |
| | sum bolder = [] | | | | | | |
| | for elem in range(1 | 100)• | | | | | |
| | random numbers | = [] | | | | | |
| | for cnt in ran | ne(1000): | | | محدة المناجعة والمعترية المحدد | | |
| | # pick 1 ra | andom number | | | Put all numbers in lists. | • • | |
| | random num | ber = random.randrange | (0. 100000) | | | | |
| | # sum those | | (0, 200000, | | | | |
| | random numb | _ bers.append(random_num | ber) | | | | |
| | sum holder.app | end(random_numbers) | 2017 | | | | |
| | # sum the 100 sums | of random numbers | | | | | |
| | <pre>s = sum([sum(elem)</pre> | for elem in sum holde | r]) | ←── | then compute sums | | |
| | print("The sum is: | <pre>{:,}".format(s))</pre> | | | | | |
| | | | | | | | |
| | | | | | | | |
| | <pre>def main():</pre> | | | | rout timoch 0 600 0 | 607 0 61 | |
| | <pre>sum_lots_of_random_</pre> | _nums() | | | Idw LIMES: 0.000 0. | | |
| | | | | | 5 loops, best of 3: | 121 msec | per loop |
| | if <u>name</u> == 'main_ | _': | | | 1 7 | | |
| | main() | | | | | | |
| | | | | | | | |
| | | | | | Not b | etter | |
| | | | | | | | |
| File O | Project 0 No Issues /Ilsers/hcf/Dror | nhox/Classes/BIOI 7800/temp/lecture | 23/example2 pv 1·1 LE LI | TE-8 Puthon | | | |

Approach #2: itertools

| | | | time to beat |
|----------|--|------------------------------------|---|
| | example3.py – /Users/bcf/Dropbox/Classes/BIOL7 | 800/temp/lecture21 | raw times, 0 601 0 6 0 508 |
| | example1.py example2.py | example3.py | |
| | #! /usr/bin/env python | | 5 LOOPS, DEST OT 3: 120 MSEC PER LOOP |
| | # encoding UTF-8 | | |
| | | | |
| | import random | | |
| | import itertools | | |
| | | | |
| | def sum lete of mender sums(). | | |
| | <pre>det sum_lots_ot_random_nums():</pre> | | |
| | $sum_notaer = []$ | | |
| | random numbers - [] | | |
| | for cot in range(1000). | | Put all numbers in lists |
| | # nick 1 random number | | |
| | random number = random.randrange() | 0. 100000) | |
| | # sum those | -, | |
| | random_numbers.append(random_numb | er) | |
| | sum_holder.append(random_numbers) | | |
| | # sum the 100 sums of random numbers | | |
| 19 | <pre>s = sum(list(itertools.chain(*sum_holder)</pre> |)) 🔶 | then use itertools to unpack lists, and sum |
| | <pre>print("The sum is: {:,}".format(s))</pre> | | aren ase reel cools co anpact hous, and sam |
| | | | |
| | | | raw times, 0 639 0 627 0 638 |
| | def main(): | ← | |
| | <pre>sum_lots_of_random_nums()</pre> | | 5 loops, best of 3: 125 msec per loop |
| | | | |
| | ifname == 'main': | | |
| | main() | | Not botton |
| | | | |
| | | | |
| 27 28 | <pre>main() Project 0 	Vo Issues /Users/bcf/Dropbox/Classes/BIOL7800/temp/lecture23/</pre> | /example3.pv 19: LF UTF-8 Python 📌 | <u>Not better</u> |

Approach #3: bigger random draw

| | | | time to beat |
|--|--|---|--|
| e 1 # 2 # 3 | example4.py — /Users/bcf/Dropbox/Classes/BIOL7800/temp/lecture21 xample1.py example2.py example3.py example4.py #! /usr/bin/env python # encoding UTF-8 | × | raw times: 0.601 0.6 0.598 5 loops, best of 3: 120 msec per loop |
| 4 i 5 6 | import random | | |
| 7 c 8 9 10 11 12 13 14 15 16 | <pre>def Sum_Lots_of_random_nums(): sum_holder = [] for elem in range(100): random_numbers = random.sample(list(range(0, 100000)), 1000)</pre> | | - Select 1000 numbers from a list (e.g. remove a loop) - Sum list of sum |
| 17 d 18 19 20 i | <pre>def main(): sum_lots_of_random_nums() if name == ' main ':</pre> | ł | _ raw times: 1.37 1.36 1.33 5 loops, best of 3: 265 msec per loop |
| 21 22 | main() | | <u>Much worse!</u> <u>2.2x slowdown</u> |
| File 0 Pro | oject 0 Volssues /Users/bcf/Dropbox/Classes/BIOL7800/temp/lecture23/example4.py 6:1 LF UTF-8 Python | ¥ | Why? |

Approach #4: bigger random draw alt. version

| | | | | | time to beat |
|---|--|--|--------------------------------------|----------------|---|
| <pre>example1.py 1 #! /usr/b 2 # encodin 3 4 import ra</pre> | example5.py — /Users example2.py in/env python g UTF-8 ndom | /bcf/Dropbox/Classes/Bl0 example3.py | DL7800/temp/lecture21 example4.py | example5.py | raw times: 0.601 0.6 0.598 5 loops, best of 3: 120 msec per loop |
| 5 6 7 def sum_l 8 sum_h 9 # mak | ots_of_random_nums older = [] e range once: | s(): | | | |
| 10 my_li 11 for e 12 r 13 s 14 # sum 15 s = s | <pre>st = list(range(0) lem in range(100); andom_numbers = ra um_holder.append(s) the 100 sums of r um(sum_holder)</pre> | 100000)) indom.sample(my_l sum(random_number random numbers | ist, 1000) s)) | • | Create range of number <u>only I time</u> Select 1000 numbers from a list (e.g. remove loop) |
| 16 print 17 18 | ("The sum is: {:,] | ".format(s)) | | — | - Sum list of sum |
| 19 def main(20 sum_l 21 |): ots_of_random_nums | .() | | | raw times: 0.457 0.443 0.443 5 loops, best of 3: 88.7 msec per loop |
| 22 ITname 23 main(24 | ==main':) | | | | <u>Waaaay better!</u> 1 3x speedup |
| File 0 Project 0 ✔ No | lssues /Users/bcf/Dropbox/C | asses/BIOL7800/temp/lectur | e23/example5.py 10: LF | UTF-8 Python 😚 | New best time! |

Approach #5: use list comprehension

| | the to deat |
|--|--|
| example6.py – /Users/bcf/Dropbox/Classes/BIOL7800/temp/lecture21 example1.py example2.py example3.py example4.py example5.py example6.py #! /usr/bin/env python # encoding UTF-8 3 | raw times: 0.457 0.443 0.443 5 loops, best of 3: 88.7 msec per loop |
| 4 import random 5 6 | |
| <pre>7 def sum_lots_of_random_nums(): 8 my_list = list(range(0, 100000)) 9 s = sum([sum(random.sample(my_list, 1000)) for i in range(0, 100)]) 10 print("The sum is: {:,}".format(s)) 11 12</pre> | • Create range of number <u>only I time</u> • Put all summing in a list comprehsion |
| <pre>13 def main(): 14 sum_lots_of_random_nums() 15 16 ifname == 'main': 17 main()</pre> | raw times: 0.466 0.457 0.462 5 loops, best of 3: 91.5 msec per loop |
| | <u>Not better</u> |
| | |

Approach #6: use numpy

time to beat example7.py — /Users/bcf/Dropbox/Classes/BIOL7800/temp/lecture21 raw times: 0.457 0.443 0.443 example7.py 5 loops, best of 3: 88.7 msec per loop import numpy def sum_lots_of_random_nums(): sum_holder = [] for elem in range(100): Use numpy to randomly sample integer 10 random_numbers = numpy.random.random_integers(0, 100000, 1000) sum_holder.append(sum(random_numbers)) Append sum(int) to a list s = sum(sum_holder) ← Sum list of sum print("The sum is: {:,}".format(s)) (use **regular** Python sum) def main(): sum_lots_of_random_nums() raw times: 0.044 0.0455 0.0454 if __name__ == '__main__': 5 loops, best of 3: 8.8 msec per loop main() Much better! 10x speedup New best time! File O Project 0 🗸 No Issues /Users/bcf/Dropbox/Classes/BIOL7800/temp/lecture23/example7.py 10: LF UTF-8 Python 🔧

Approach #7: use numpy (with numpy.sum)

time to beat example8.py - /Users/bcf/Dropbox/Classes/BIOL7800/temp/lecture21 raw times: 0.044 0.0455 0.0454 example8.p example6. 5 loops, best of 3: 8.8 msec per loop import numpy def sum_lots_of_random_nums(): sum_holder = [] for elem in range(100): Use numpy to randomly sample integer random_numbers = numpy.random.random_integers(0, 100000, 1000) sum_holder.append(numpy.sum(random_numbers)) \frown Append numpy.sum(int) to a list s = sum(sum_holder) ← Sum list of sum print("The sum is: {:,}".format(s)) (use **regular** Python sum) def main(): sum_lots_of_random_nums() raw times: 0.00844 0.00876 0.00837 if __name__ == '__main__': 5 loops, best of 3: 1.67 msec per loop main() Much better! 5.2x speedup (over 1st numpy version) 71.0x speedup (over 1st python version)

20

Making things faster Can we go even faster?

Approach #8: everything numpy

time to beat

| | example9.pv - /Users/bcf/Dropbox/Classes/BIOL7800/temp/lecture21 | | |
|----------|---|------------|--|
| example: | example2.p example3.p example4.p example5.p example6.p example7.p example8.p #! /usr/bin/env python | example9.p | raw times: 0.00844 0.00876 0.00837 5 loops, best of 3: 1.67 msec per loop |
| | # encoding UIF—8 | | |
| 4 5 | import numpy | | |
| | <pre>def sum_lots_of_random_nums(): s = numpy.sum(numpy.random.randint(0, 100000, (1000, 100))) print("The sum is: {:,}".format(s))</pre> | | Select an array of 100 rows of 1000 random numbers and numpy.sum those |
| | <pre>def main(): sum_lots_of_random_nums()</pre> | ← | raw times: 0.00731 0.00662 0.00728 5 loops, best of 3: 1.32 msec per loop |
| | <pre>ifname == 'main': main()</pre> | | <u>Better!</u> |
| | | | 1.2x speedup (over 2nd numpy version) |
| | | | 90x speedup (over 1st python version) |
| | | | |
| | | | |
| | | | |

Now, I've shown you all of those to show you that <u>how</u> you write your code is very, very important to how fast it runs

We achieved ~90x speedup by optimizing our code! (that's pretty darn amazing)



So being efficient about how you write your code can make a **huge** difference

Write more efficient code

We can also achieve speedups by taking advantage of multiple CPUs on many (current) computers



One important thing to remember is that <u>physical CPUs</u> are what's (mostly) important And most Intel chips have both <u>physical</u> and virtual CPUs (via "hyperthreading")

The number of physical CPUs you have is usually 1/2 that of physical + virtual CPUs

In Python, the **multiprocessing** module helps us parallelize code And, the **multiprocessing**. Pool class is the easiest way to use **multiprocessing**



Python creates copies of your function and "loads" one copy on each CPU



In:

Each CPU processes the incoming data, and writes the results to a list



Execution moves to the next data items, and those results are written to the output list



Execution moves to the last data item, and those results are written to the output list



| • • • | example10.py - /Users/bcf/Dropbox/Classes/BIOL7800/temp/lecture | 21 |
|---|---|--|
| | example2.py example3.py example4.py example5.py example6.py example7.py | example8.py example9.py example10.py time to beat |
| | <pre>#! /usr/bin/env python # encoding UTF-8 import random</pre> | raw times: 0.601 0.6 0.598 5 loops, best of 3: 120 msec per loop |
| | import multiprocessing | |
| | <pre>def sum_lots_of_random_nums(): sum_holder = []</pre> | Make a list of lists |
| | <pre>for elem in range(100): random_numbers = [] for cnt in range(1000): # pick 1 random number</pre> | [[x, y,, z], [x y z] |
| 14 15 16 17 18 19 20 21 | <pre>random_number = random.randrange(0, 100000) # sum those random_numbers.append(random_number) # sum the 1000 random numbers sum_holder.append(random_numbers) return sum_holder</pre> | [x, y,, z], [x, y,, z],] |
| | <pre>def get_sum(list_item): return sum(list_item)</pre> | This function just sums a list |
| | <pre>def main(): random_numbers = sum_lots_of_random_nums() pool = multiprocessing.Pool(4) sum_of_lists = pool.map(get_sum, random_numbers) s = sum(sum_of_lists) print("The sum is: {:,}".format(s)) pool.close() ifname == 'main':</pre> | We Pool.map the list of lists onto get_sum. Each CPU sums one sublist and returns the sum. |
| | main() | raw times: 0.71 0.706 0.705 5 loops, best of 3: 141 msec per loop |
| File 0 | Project 0 Volssues /Users/bcf/Dropbox/Classes/BIOL7800/temp/lecture23/example10.py* 14:56 | LF UTF-8 Python * Much slower! |

No multiprocessing

Multiprocessing

time to beat

raw times: 0.601 0.6 0.598
5 loops, best of 3: 120 msec per loop

raw times: 0.71 0.706 0.705
5 loops, best of 3: 141 msec per loop

Much slower!

Why, if we are using more CPUs is multiprocessing much slower?

(and we're nowhere near results we saw from numpy)

Making things faster Why, if we are using more CPUs is multiprocessing much slower?

There's a fair amount of "overhead" setting up multiprocessing



Making things faster Why, if we are using more CPUs is **multiprocessing**

much slower?

We often need job components to be more intensive in order to see benefits



Getting lots of data in to multiprocessing.Pool.map() can be tricky You need to "package" your data into lists/tuples

Let's say you have 4 units of "work"



Getting lots of data in to multiprocessing.Pool.map() can be tricky You have 4 units of "work"



In: def my_function(work): sum = var1 + var2 + var3 + var4 return sum So, how can we make 4 work

units of 4 variable into 1 work unit?

In: my_result = pool.map(my_function, <something here>)

Getting lots of data in to multiprocessing.Pool.map() can be tricky



In: my_result = pool.map(my_function, work)

Getting lots of data out of multiprocessing.Pool.map() can also be tricky

```
In: def my_function(work):
    var1, var2, var3 = work
    var1 += 1
    var2 += 2
    var3 += 3
    return ????
    The mapped function can only return one result
        How can we deal with this?
In: work = (
        (var1, var2, var3, var4), unit 1
        (var1, var2, var3, var4), unit 2
        (var1, var2, var3, var4), unit 3
        )
```

In: my_result = pool.map(my_function, work)

Getting lots of data out of multiprocessing.Pool.map() can also be tricky

```
In: def my_function(work):
       varl, var2, var3 = work
       varl += 1
       var2 += 2
                              The mapped function can only return one result
       var3 += 3
       return [var1, var2, var3]
                                   ----- We can package those data into a single tuple or list
In: work = (
             (var1, var2, var3, var4),
            (varl, var2, var3, var4),
                                     unit l
            (var1, var2, var3, var4),
                                     unit 2
                                     unit 3
In: my_result = pool.map(my_function, work)
In: for elem in my_result:
      var1, var2, var3 = elem
                                                 – And then unpack them
      # do stuff
```

How do you determine the number of CPUs on a system?

It's a little hard...

Because hyperthreading is hard to detect

Approach #1

```
def get_args():
    parser = argparse.ArgumentParser(
        description="""Do stuff in parallel""",
    )
    parser.add_argument(
    "--cores",
    type=int,
    default=1,
    help="Process in parallel using --cores"
    )

def main():
    args = get_args()
    pool = multiprocessing.Pool(args.cores)
    my_result = pool.map(my_function, [2, 4, 6, 8, 10])
```

How do you determine the number of CPUs on a system?

It's a little hard... Because hyperthreading is hard to detect

Approach #2