



# Python Multithreading

Speed up your Scripts 1,000%

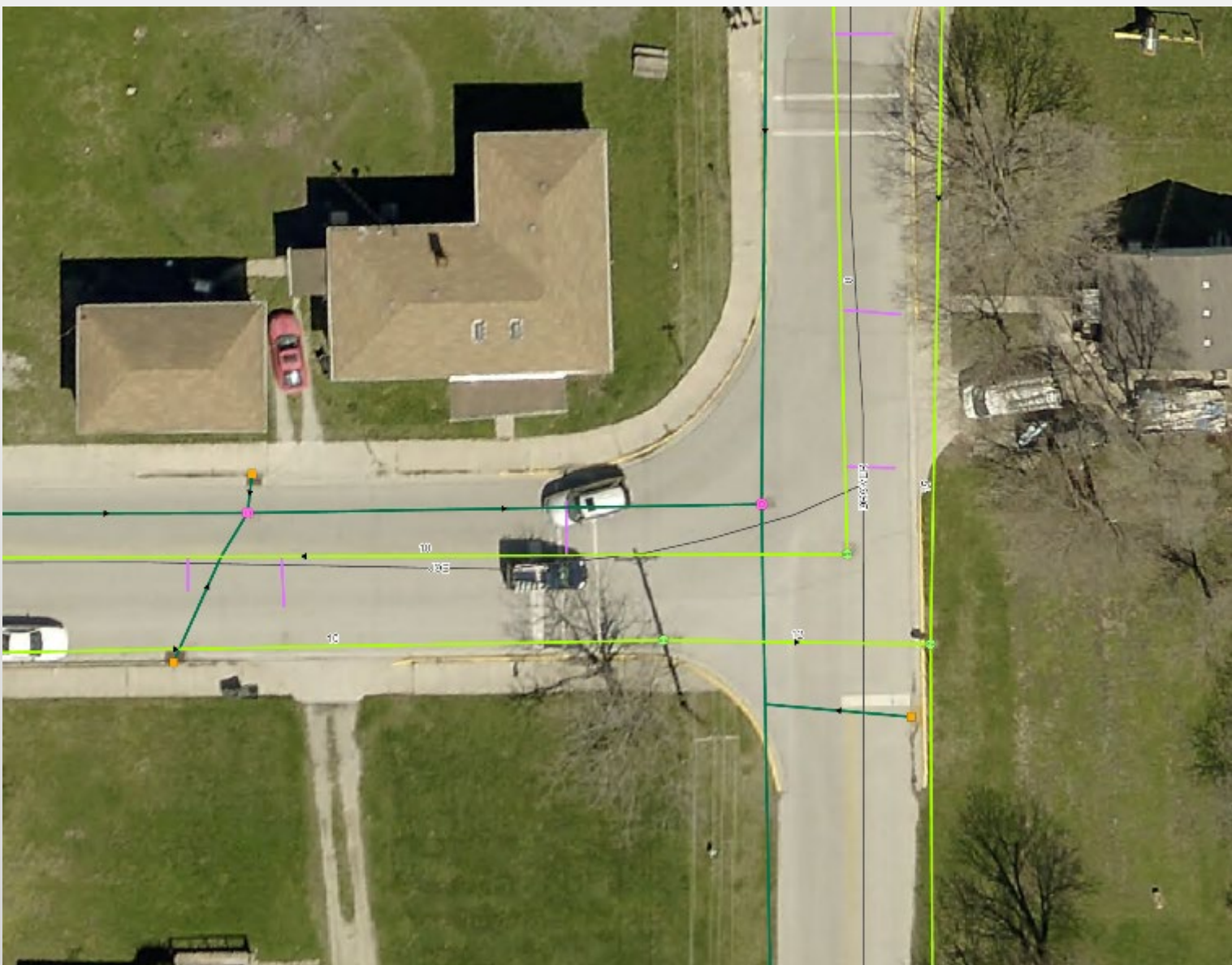




A red rectangular stamp with a distressed, ink-like texture. The word "EXAMPLE" is written in a bold, sans-serif font within the stamp. The stamp is tilted at an angle. In the bottom-left corner of the image, there are decorative diagonal stripes in blue and grey.

**EXAMPLE**







▲ You can use [map-reduce](#) for this task.

1

▼ Pseudo code:

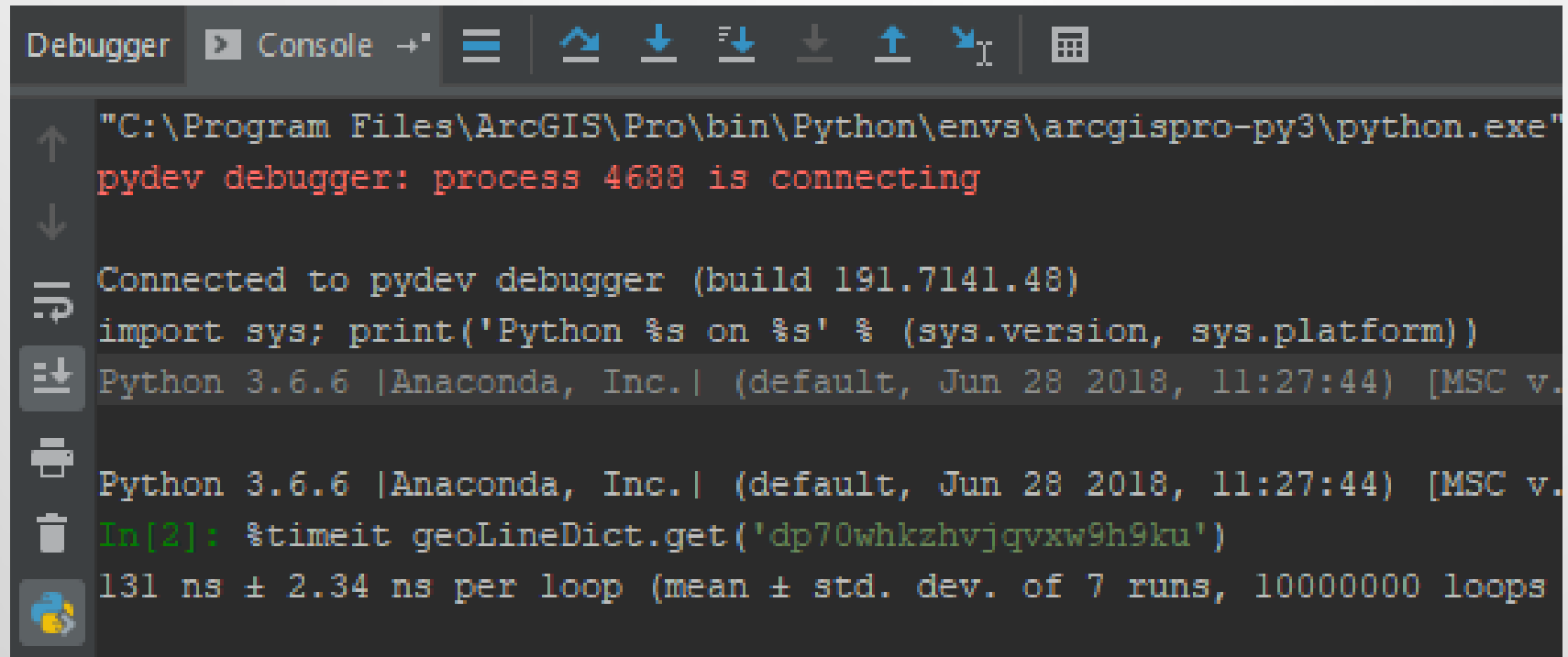
```
map1(list): #runs on first file
    for each (i,x,y) in list:
        emit ((x,y),(1,i))
map2(list): #runs on 2nd file
    for each (x,y,temp) in list:
        emit ((x,y),(2,temp))
reduce((x,y),list): #runs on output of both mappers
    for each (aux, val) in list:
        if aux == 1:
            i = val
        else:
            temp = val
    if both i and temp initialized:
        emit(i,temp)
```

The idea is first to map each of the files into some kind of hash table (this is done internally by the framework), and you have two hash tables:

1. key=(x,y) value = id
2. key=(x,y) value = temprature

Once you have both hash tables, it is easy to find which id is connected to which temprature in a single pass, and once a connection is made -output it.

Complexity of this code is  $O(n)$  average case.



```
Debugger Console → [Icons]

"C:\Program Files\ArcGIS\Pro\bin\Python\envs\arcgispro-py3\python.exe"
pydev debugger: process 4688 is connecting

Connected to pydev debugger (build 191.7141.48)
import sys; print('Python %s on %s' % (sys.version, sys.platform))
Python 3.6.6 |Anaconda, Inc.| (default, Jun 28 2018, 11:27:44) [MSC v.
Python 3.6.6 |Anaconda, Inc.| (default, Jun 28 2018, 11:27:44) [MSC v.
In[2]: %timeit geoLineDict.get('dp70whkzhvjqvwx9h9ku')
131 ns ± 2.34 ns per loop (mean ± std. dev. of 7 runs, 10000000 loops
```



# The Great Trash Cart Debacle of 2018



# Geocoder

```
# Single Line Geocoder Indiana Composite Locator
# Author: R. Stevens; B. Bond
import pandas
import requests
import time
from tqdm import tqdm

# Variables
inputTable = input("InputTable> ")
outputTable = input("Output Locaation> ")
addressColumn = "Address" # from input table
idColumn = "ID" # from input table
```



```
def geocode(address, id_num):
    geocode_url =
        "https://gis.in.gov/arcgis/rest/services/Indiana_Compos
        ite_Locator/GeocodeServer/findAddressCandidates?SingleL
        ine={}&f=json&outSR=4326&maxLocation=1&outFields=*".format(address)
    results = requests.get(geocode_url)
    results_dict = results.json()
    candidates = results_dict['candidates'][0]
    # Gather attributes and Sub Dicts
    location, attributes, score = candidates['location'],
        candidates['attributes'], candidates['score']
    # build output
    output = {
        "UniqueID": id_num,
        "input_string": address,
        "lat": location.get('y'),
        "lon": location.get('x'),
        "match_addr": attributes.get('Match_addr'),
        "match_type": attributes.get('Addr_type;'),
        "score": score
    }
    return output
```





GdalTranslateEshn.py  
geoHash.py  
geoHash.spec  
GPSAggregate.py  
HRRRetireeFormGenerator.py  
Indiana\_Geocoder\_RyanOrig.py  
IndianaGeocoder.py  
IndianaGeocoderAsync.py  
IndianaGeocoderConcurrent.py  
IndianaGeocoderConcurrent2.py

23  
24  
25  
26  
27  
28  
29  
30  
31

```
if addressColumn not in df.columns:  
    raise ValueError("Missing Address Column in input data")  
  
print("DataFrame Created") # use tqdm.write if progress bar messed up, its 5  
# for progress bar limits and worker count  
number_of_items = len(df.index)  
# Calculate Number of Workers any more than number of records to geocode is a  
workers = min(maxWorkers, number_of_items)
```

Run: IndianaGeocoder x

"C:\Program Files\ArcGIS\Pro\bin\Python\envs\arcgispro-py3\python.exe" L:/Scripts/Python3/IndianaGeocoder.py

2: Structure  
2: Favorites

4: Run 6: TODO 9: Version Control Terminal Python Console

4:16 CRLF ↵ U

**SLOTH**  
**RUNNING TEAM**



**WE'LL GET THERE**  
**WHEN WE GET THERE**



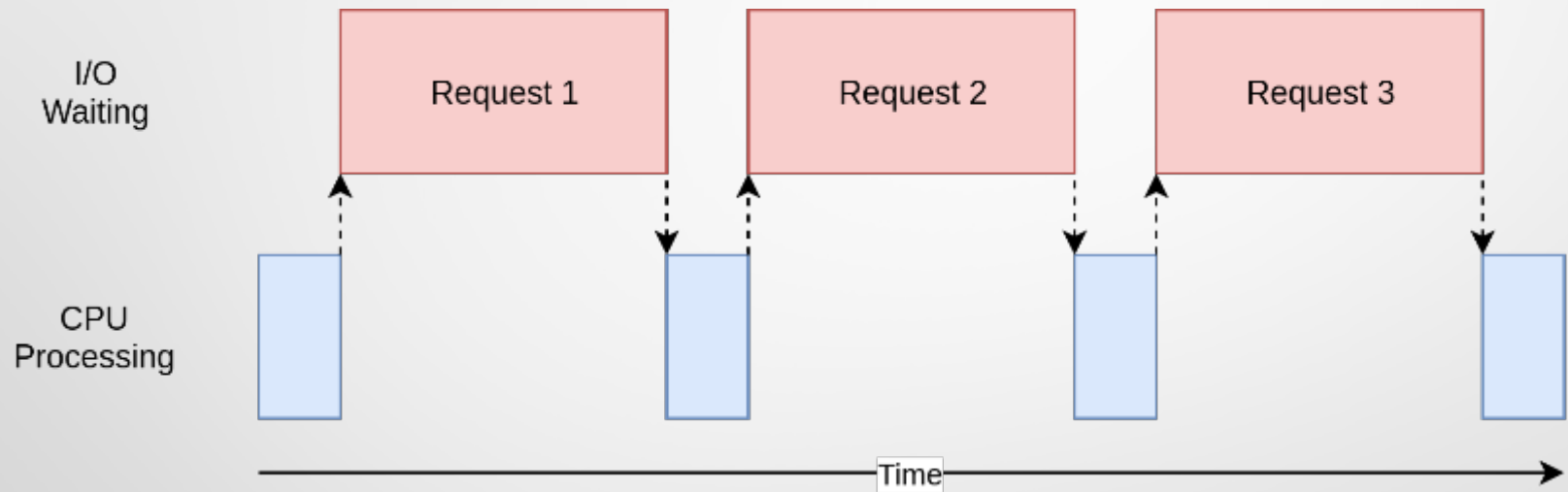




## 2 Types of Problems

- I/O bound
- CPU bound

# I/O Bound



# I/O Speeds

Device	CPU cycles	Proportional “human” scale
L1 cache	3	3 seconds
L2 cache	14	14 seconds
RAM	250	250 seconds
disk	41,000,000	1.3 years
network	240,000,000	7.6 years









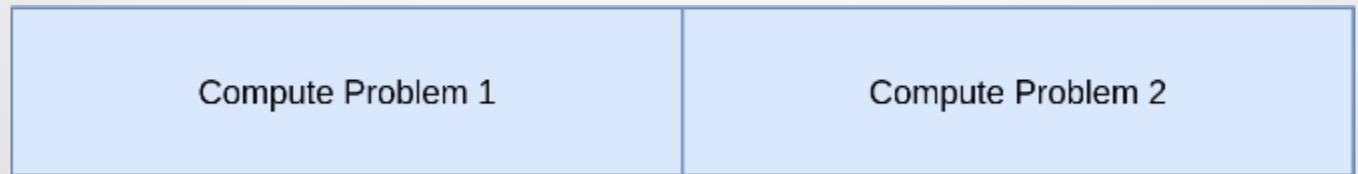


**SHE SAID SHE'D BE  
5 MINUTES**

# CPU Bound

I/O  
Waiting

CPU  
Processing



Time







# DID YOU KNOW?

To make your car go faster you can just shift into "race" gear after fifth.





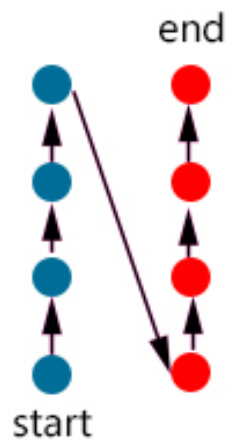
A Venn diagram consisting of two concentric circles. The outer circle is dark blue and contains the text 'Concurrency' and '(Threading, Async IO)'. The inner circle is light blue and contains the text 'Parallelism' and '(Multiprocessing)'. The inner circle is entirely contained within the outer circle, indicating that parallelism is a subset of concurrency.

Concurrency  
(Threading, Async IO)

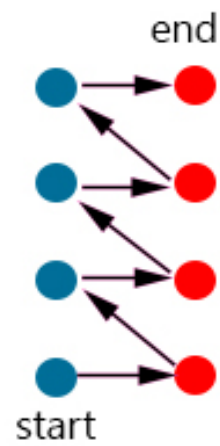
Parallelism  
(Multiprocessing)



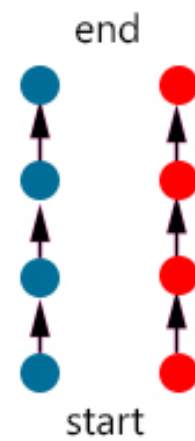
Sequential



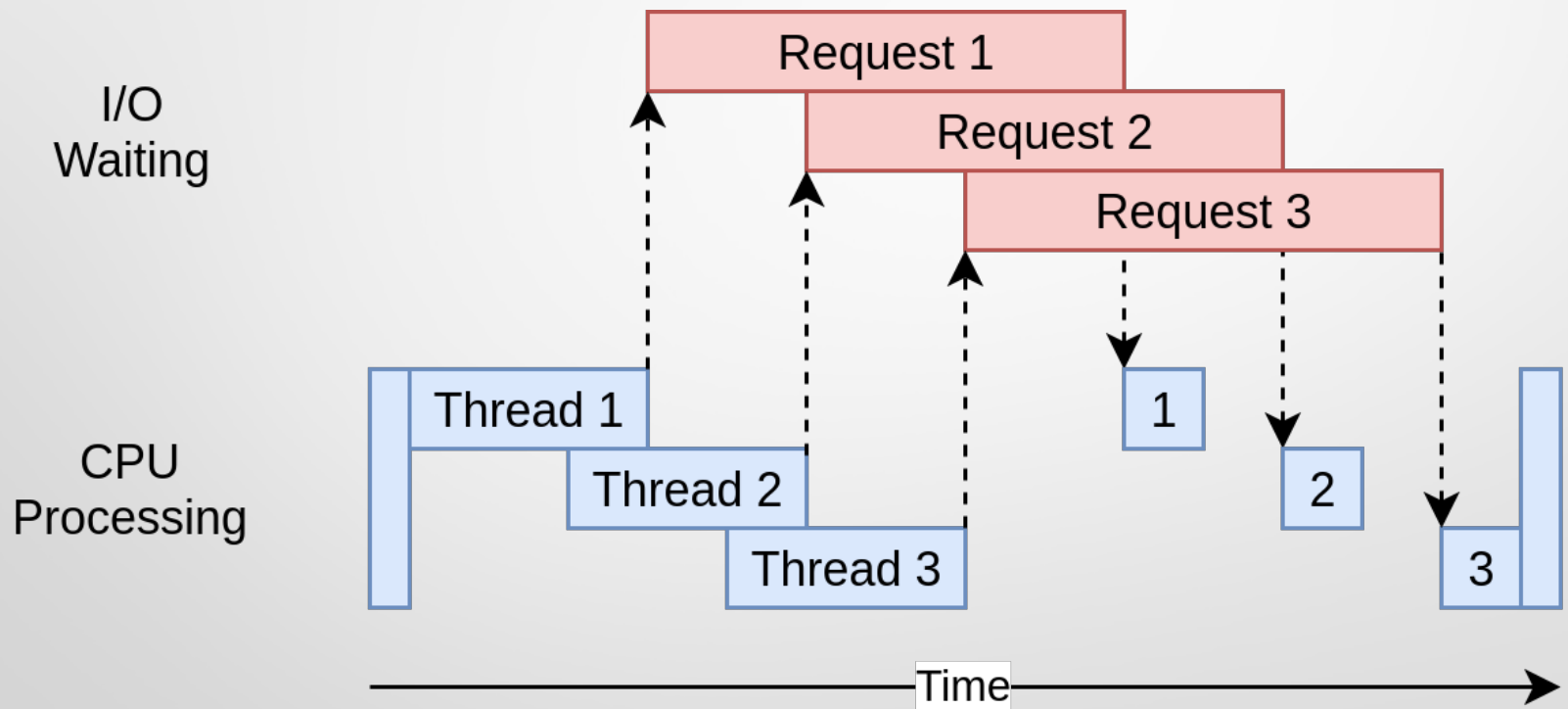
Concurrent



Parallel



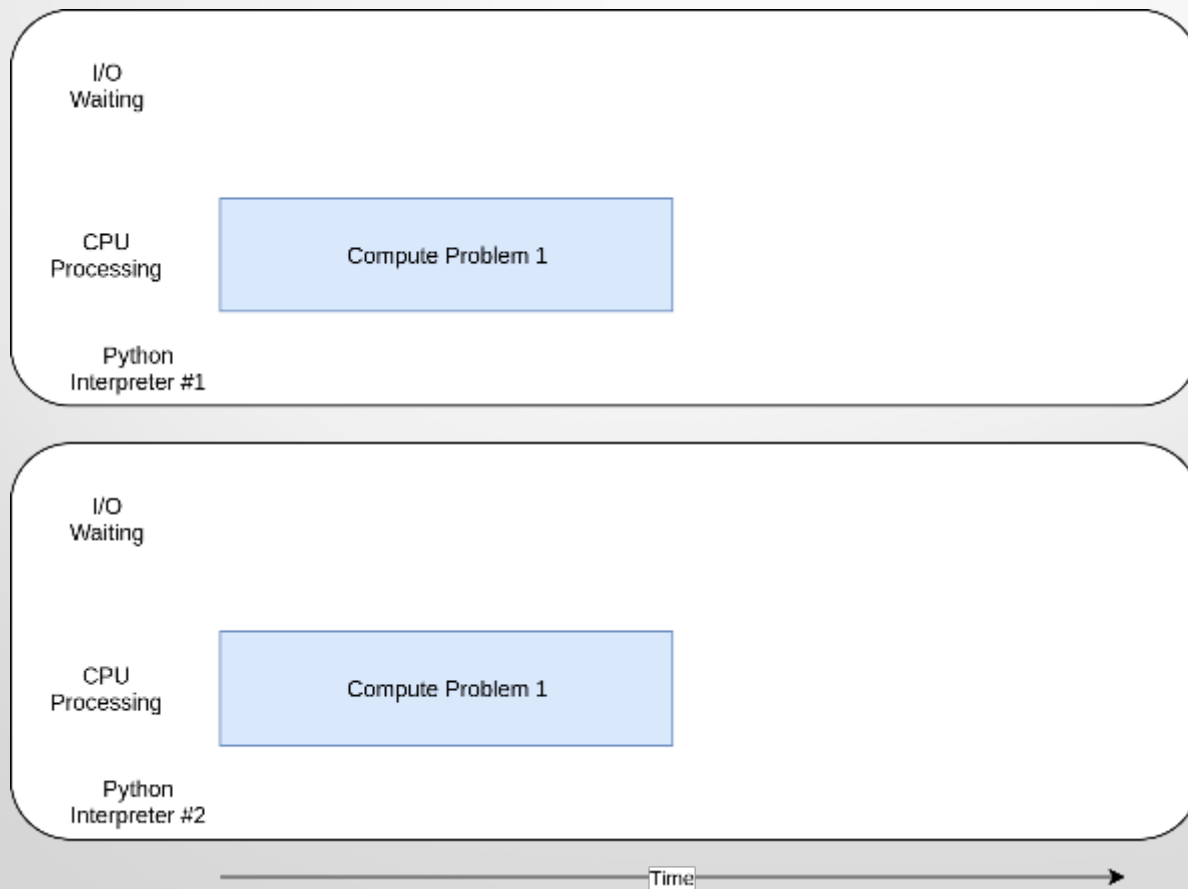
# Multithreading





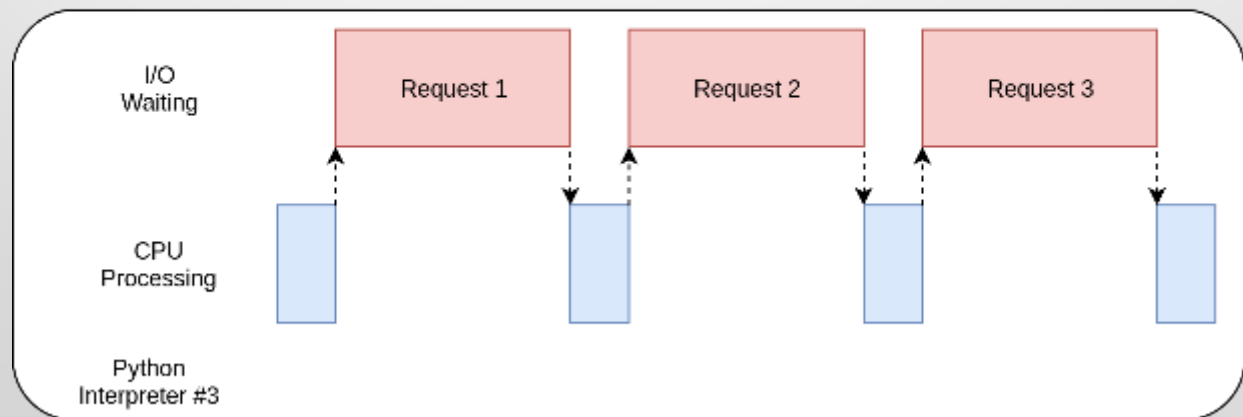
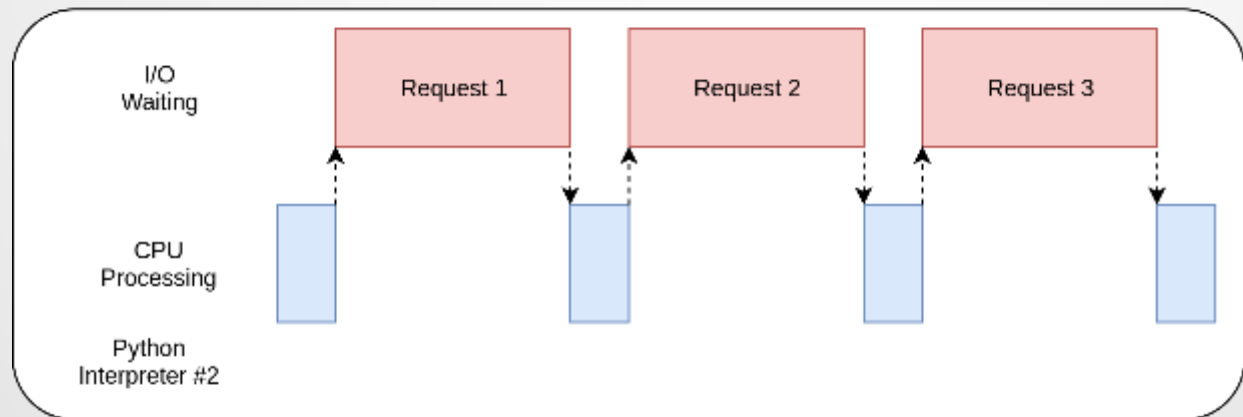
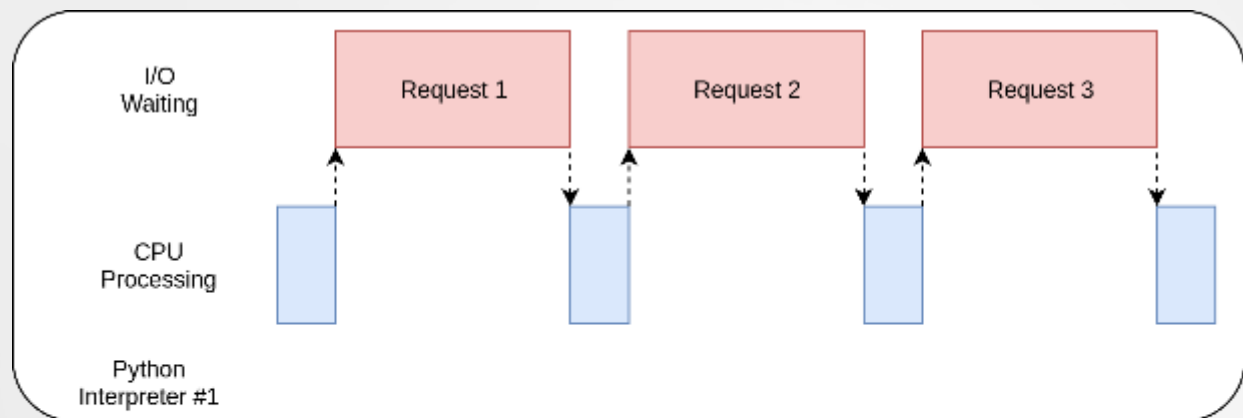


# Multiprocessing









Time →







# Quiz: How can we make the geocoder faster

- Buy it a Ferrari!
- Buy a bigger computer
- Multithreading (Concurrent)
- Asynchronous (Concurrent)
- Multiprocessing (Parallel)



# Concurrent Package

- Used to be Hard, threading, multiprocessing
- Concurrent
  - Multithreading
    - `Concurrent.futures.ThreadPoolExecutor(workers)`
  - Multiprocessing
    - `Concurrent.futures.ProcessPoolExecutor()`
- 2 functions
  - `Futures.map()`
  - `Executor.submit()`

# Multithreading (Map)

```
from concurrent import futures

maxWorkers = 20

def func():
    # do something to a single file or response

def concurrent_threading(process_list):
    workers = min(maxWorkers, len(process_list))
    with futures.ThreadPoolExecutor(workers) as executor:
        response = executor.map(func, process_list)
```

# Multiprocessing (Map)

```
from concurrent import futures

def func():
    # do something to a single file or response

def concurrent_processing(process_list):
    with futures.ProcessPoolExecutor() as executor:
        response = executor.map(func, process_list)
```



# Multithreading (submit)

```
from concurrent import futures

maxWorkers = 20

def func():
    # do something to a single file or response
    pass

def concurrent_threading(process_list):
    workers = min(maxWorkers, len(process_list))
    futures_list = []
    with futures.ThreadPoolExecutor(workers) as executor:
        for item_to_process in process_list:
            future = executor.submit(func, item_to_process)
            futures_list.append(future)

    done_iter = futures.as_completed(futures_list)
    for itm in done_iter:
        # Do stuff
```



**You got a problem with  
my driving, punk?**





GdalTranslateEshn.py  
geoHash.py  
geoHash.spec  
GPSAggregate.py  
HRRRetireeFormGenerator.py  
Indiana\_Geocoder\_RyanOrig.py  
IndianaGeocoder.py  
IndianaGeocoderAsync.py  
IndianaGeocoderConcurrent.py  
IndianaGeocoderConcurrent2.py

23  
24  
25  
26  
27  
28  
29  
30  
31

```
if addressColumn not in df.columns:  
    raise ValueError("Missing Address Column in input data")  
  
print("DataFrame Created") # use tqdm.write if progress bar messed up, its 5  
# for progress bar limits and worker count  
number_of_items = len(df.index)  
# Calculate Number of Workers any more than number of records to geocode is a  
workers = min(maxWorkers, number_of_items)
```

Run: IndianaGeocoderConcurrent2 x

"C:\Program Files\ArcGIS\Pro\bin\Python\envs\arcgispro-py3\python.exe" L:/Scripts/Python3/IndianaGeocoderConcurrent2.p

Structure  
Favorites



4: Run 6: TODO 9: Version Control Terminal Python Console

14:20 CRLF ↵ U



CONTINUE?

YES/NO

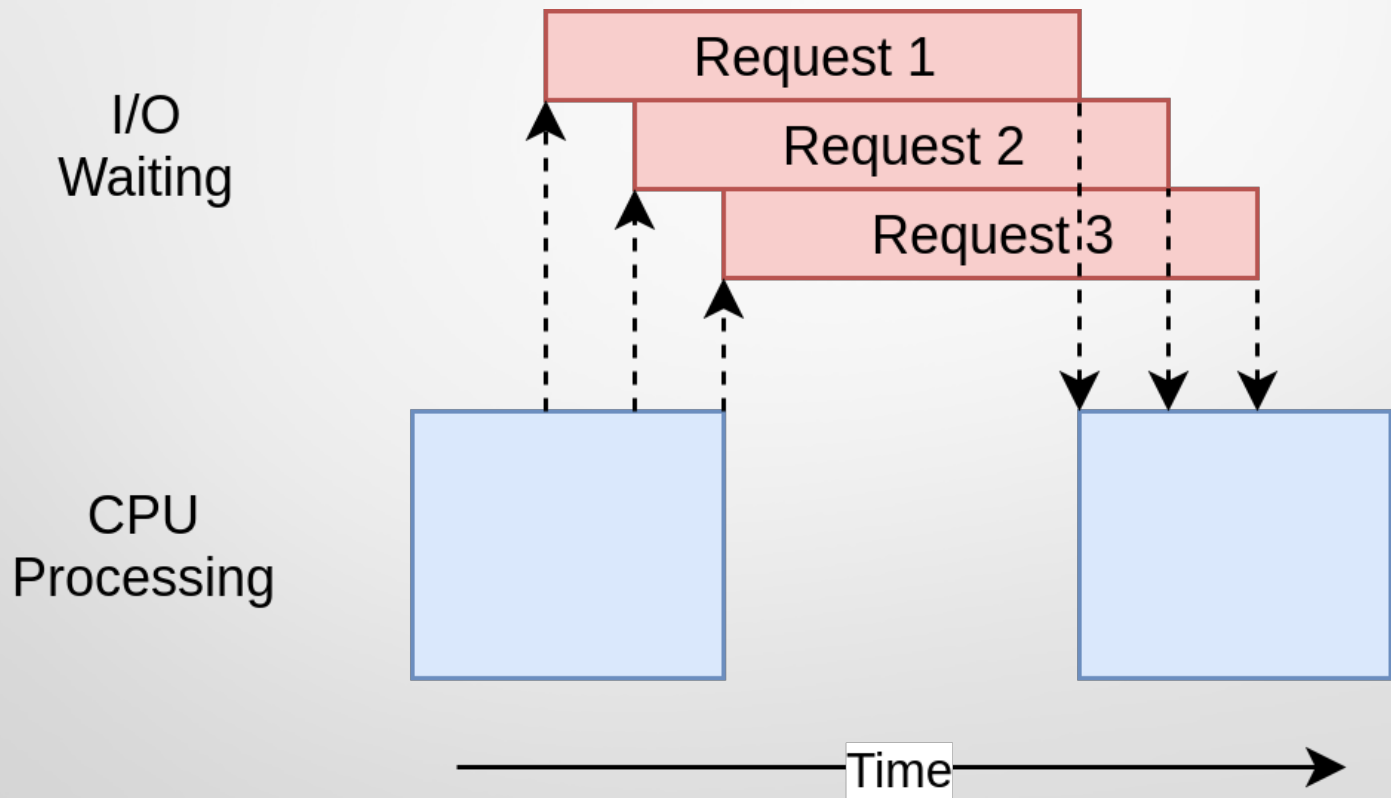
CREDIT 0

time  
machine



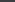
ASYNCIO

# Asynchronous





- GeoTranslator.py
- geoHash.py
- geoHash.spec
- GPSAggregate.py
- HRRRetireeFormGenerator.py
- Indiana\_Geocoder\_RyanOrig.py
- IndianaGeocoder.py
- IndianaGeocoderAsync.py
- IndianaGeocoderConcurrent.py
- IndianaGeocoderConcurrent2.py

Run:  IndianaGeocoderAsync ✕

- ★ 2: Favorites
- ■ 7: Structure

▶ 4: Run    ⚙ 5: Debug    ☰ 6: TODO    ↶ 9: Version Control    📄 Terminal    📄 Python Console

Push successful: Pushed 1 commit to origin/master (41 minutes ago)

16:1 CRLF ⇄ U

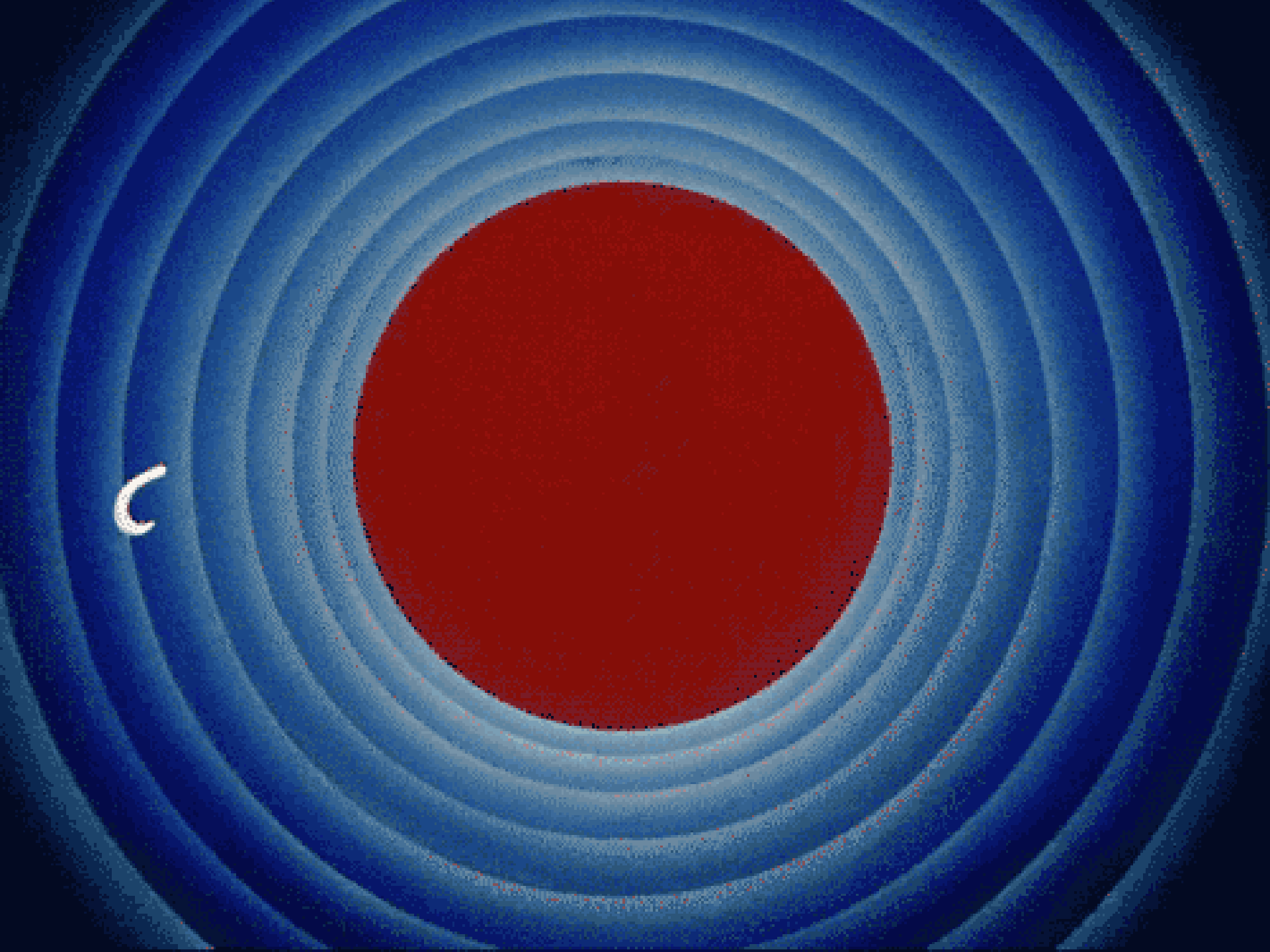






# Food for Thought

- Know your problem
- Make sure you need to speedup your script
- Write sequentially before going concurrent or parallel





# Questions?

[Ben.bond@huntington.in.us](mailto:Ben.bond@huntington.in.us)