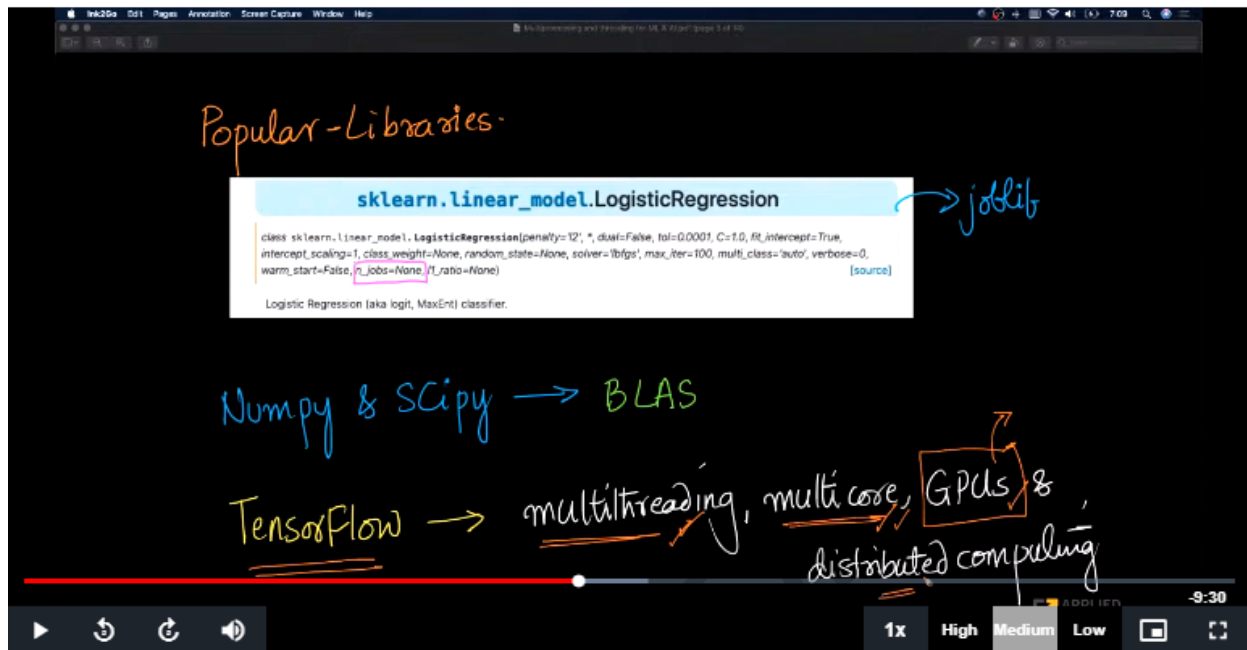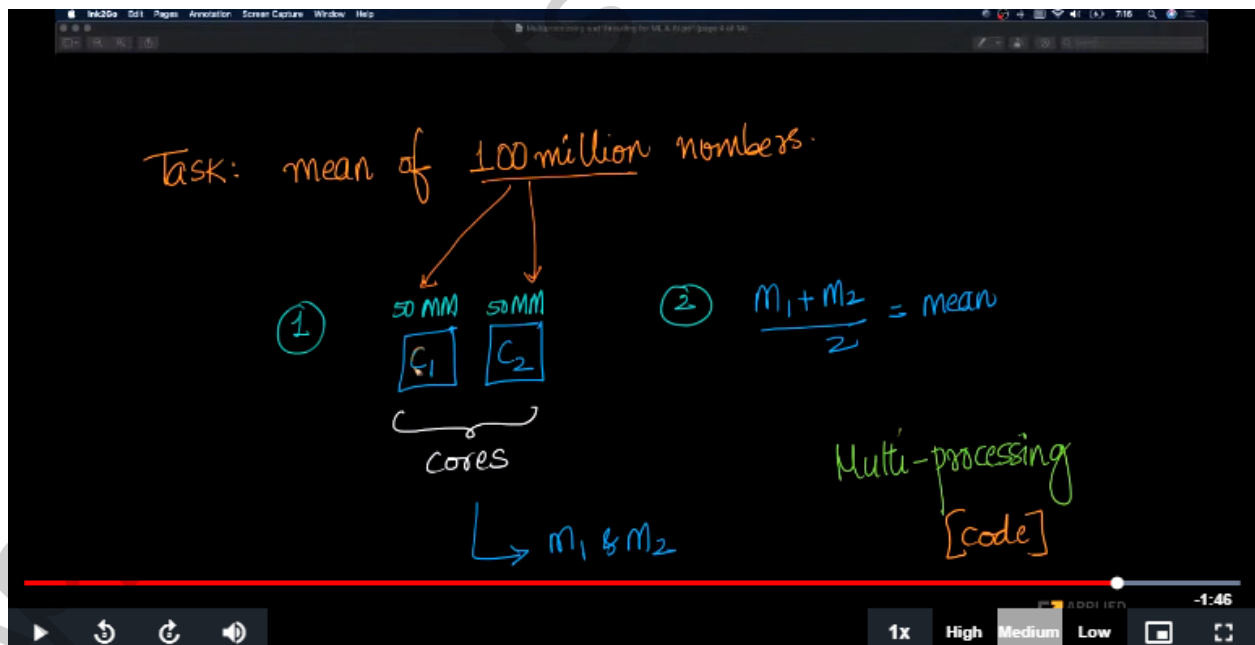# 17.1 Applied Introduction and Overview



Timestamp: 8:02

Many popular machine learning libraries like sklearn, numpy, scipy and tensorflow uses some form of parallelization libraries. Sklearn uses joblib, numpy and script internally uses BLAS. Tensorflow can perform tasks using multithreading, multicore, it can even use GPUs and distributed computing.
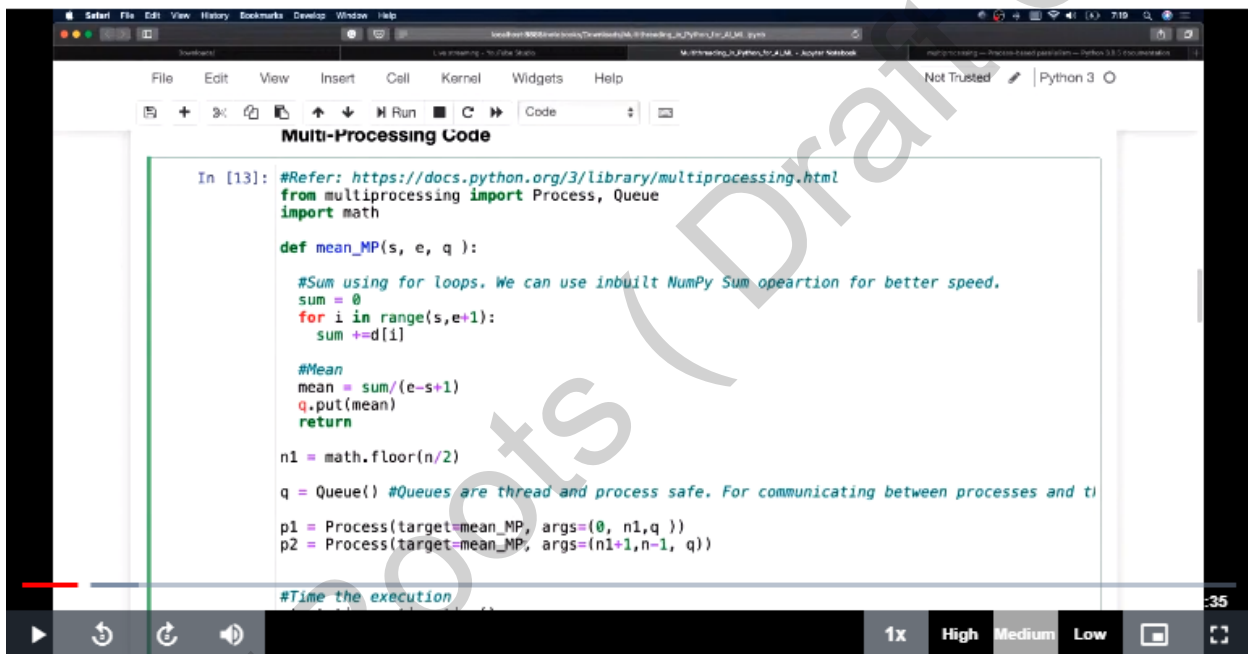
Timestamp: 14:58

As shown in the above figure a program to calculate mean of 100 million numbers takes 20.98 sec. Note that this program runs on a single core of the processor with this code.



Timestamp: 15:46

We can parallelize the code as shown in the above figure by splitting the 100 million numbers into 50M and 50M and compute the means m1,m2 of both the parts using cpu cores c1 & c2 respectively. Finally, we can combine both the means as (m1+m2)/2. Since each core can run independently, both run parallely to calculate respective means, hence the time taken will be reduced largely.

# 17.2 Multiprocessing: Compute the mean of 100 Million numbers

Timestamp: 1:10

As shown in the above figure, we import Process, Queue from multiprocessing library which will help us create processes and the Queue using which the processes can communicate. We define a function mean_Mp that calculate the mean of the array elements from the start index s to the index e and store the mean corresponding to the part of the array in the queue.

We define the process p1 and p2 to compute the means of the first and second part of the array d respectively.
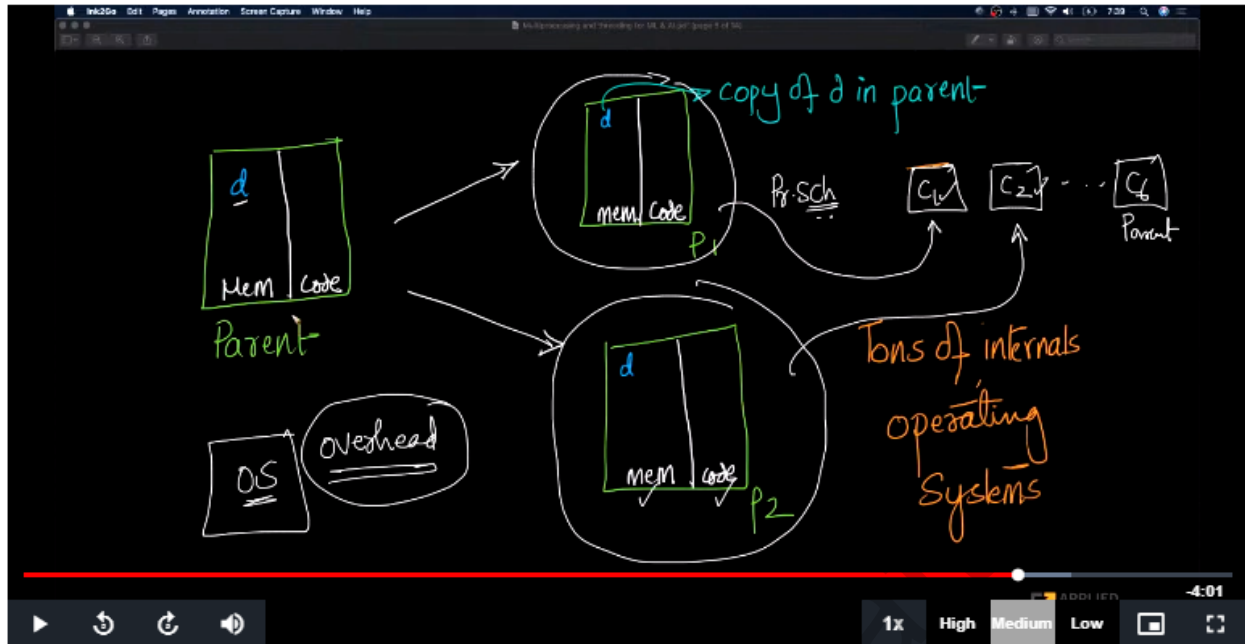
```
q = Queue() #Queues are thread and process safe. For communicating between processes and t

p1 = Process(target=mean_MP, args=(0, n1,q ))
p2 = Process(target=mean_MP, args=(n1+1,n-1, q))

#Time the execution
start_time = time.time()

p1.start()
p2.start()

p1.join() # Wait till p1 finishes
p2.join()

m=0;
while npt q.empty():
    m += q.get()

m /= 2;

end_time = time.time()
print (end_time-start_time)
print(m)
```

```
11.001178979873657
```

Timestamp: 14:37

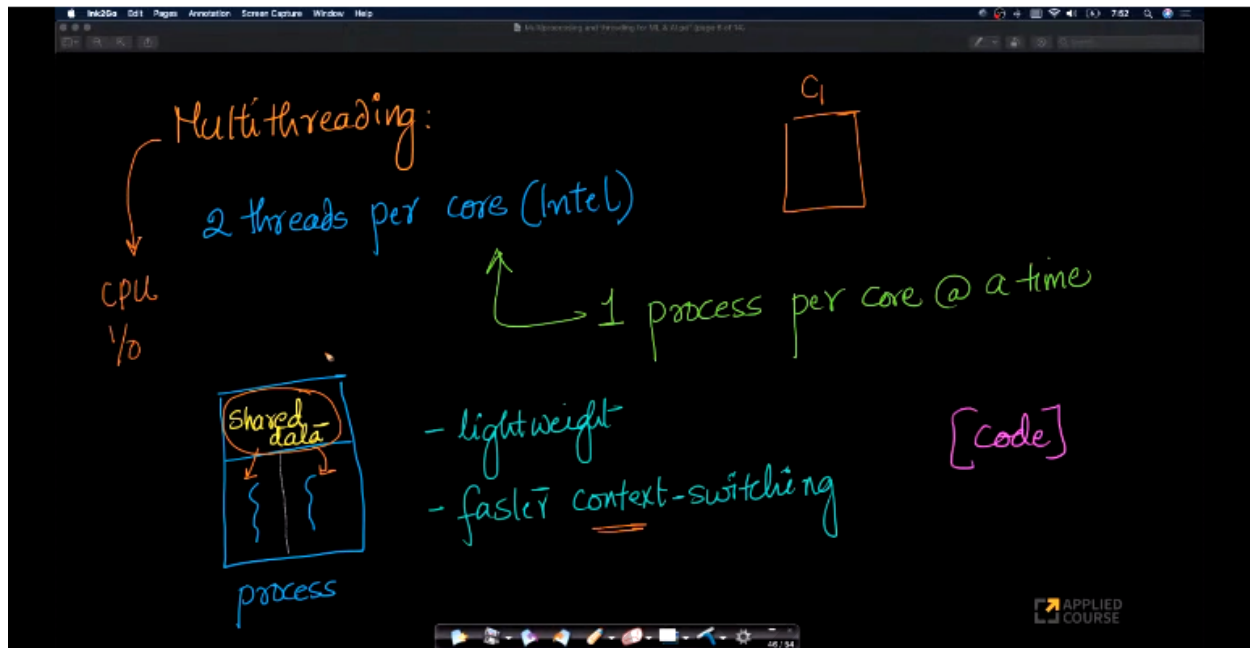Notice in the above figure, we start the processes using start() and we combine them using join(). join() is used, so that the parent process which actually runs this code waits for the two child processes responsible for calculating the means. So once the child processes complete calculating means and storing in the queue, we then calculate the mean of the entire array in the parent process using get method of queue. Notice that the time taken for calculating the mean of the array using two process is 11.00 sec. This is a significant improvement over time taken by 1 process.

Timestamp: 18:44

Notice that while creating these child processes, the parent process copies d to both processes P1 and P2. Hence both these processes has access to d for calculating the means. If we notice, with 2 processes it took 11.00 sec instead of 20.98/2=10.47 sec (due to two cores). This additional time came from the OS overhead of creating two processes with the required code, memory and allocating them to individual cores. Note that during the entire time when both the child processes is computing the means, the parent process sleeps since it cannot proceed further unless the child processes finishes executing their code.

# 17.3 Multi-threading with code



Timestamp: 9:00

Using multi-threading, we can run upto 2 threads per core. Threads are known as light weight processes with faster context switching. Context switching is the time to create two process with all the required code and memory. Since threads are executed on the same core and use same shared data they have faster context-switching times.

Multithreading is helpful when we have two threads both executing CPU bound and I/O bound operations, so that when one thread is performing the I/O operations, the other thread can run its code on the cpu.

Timestamp: 12:56

Multithreading can be done in python using the threading library. Similar to how we have created processes we have created threads but in addition we give the thread number 0 and 1 to the function mean_MT and store the means in 0 and 1 st index of the means since both the threads share common memory.



Timestamp: 14:53

Notice that the speed from 1 thread to 2 threads for the given problem is very less. This is since both are CPU bound process there is not much parallelism going on. It is also due to a problem called GIL to be explained below.



Timestamp: 20:14

While using multithreading in python, both the threads need access to the python interpreter to execute their code. But in python only one thread can access the python interpreter and the other thread is left waiting for the other thread to release the python interpreter. This is called global interpreter lock.This releasing happens when the thread is waiting for some I/O or disk operations.

Note that python's memory management is not thread safe meaning if different threads write to a single memory then it can get messed up.

# 17.4 Joblib: Simplified Parallel Programming in python



Timestamp: 3:14

Instead of using the raw multithreading and multiprocessing libraries we can work with joblib which is a simple parallel computing library in python. It can highly optimize for loops and it can use disk caching of function outputs.

So when we pass an input a to a function f and get output o, it store that information that a function f when given input a gives output o. So the next time when someone calls the same function with the same input a, it searches whether it has seen the same input for the same function, if found in cache it returns the corresponding output o without computing it again.

It is widely used in scikit learn.

Timestamp: 9:01

Notice in the above code with no parallelization, it took 17.11 seconds to run the function f(i) for i in range(n).



Timestamp: 10:10

Notice that using joblib with 2 cores specified by n_jobs it took only 9.58 sec. Please refer the below link to learn about delayed. But what it does is to say the function we want to run and the inputs to run it with and let all the inputs be ready before we parallelize.

https://stackoverflow.com/questions/42220458/what-does-the-delayed-function-do-when-used-with-joblib-in-python#:~:text=Delayed%20is%20a%20decorator%20that,and%20popped%20out%20as%20needed.



Timestamp: 17:41

We can also do multithreading using joblib as shown in the above figure. Note that the time taken for 2 threads is longer than the time taken with 1 thread. This is due to GIL and since both the threads are cpu bound.

Timestamp: 18:35

Notice that with 6 processes it took 4.55 which is longer than 17.12/6=2.85. This is because of the overhead of creating these 6 process and tracking them and content switch.