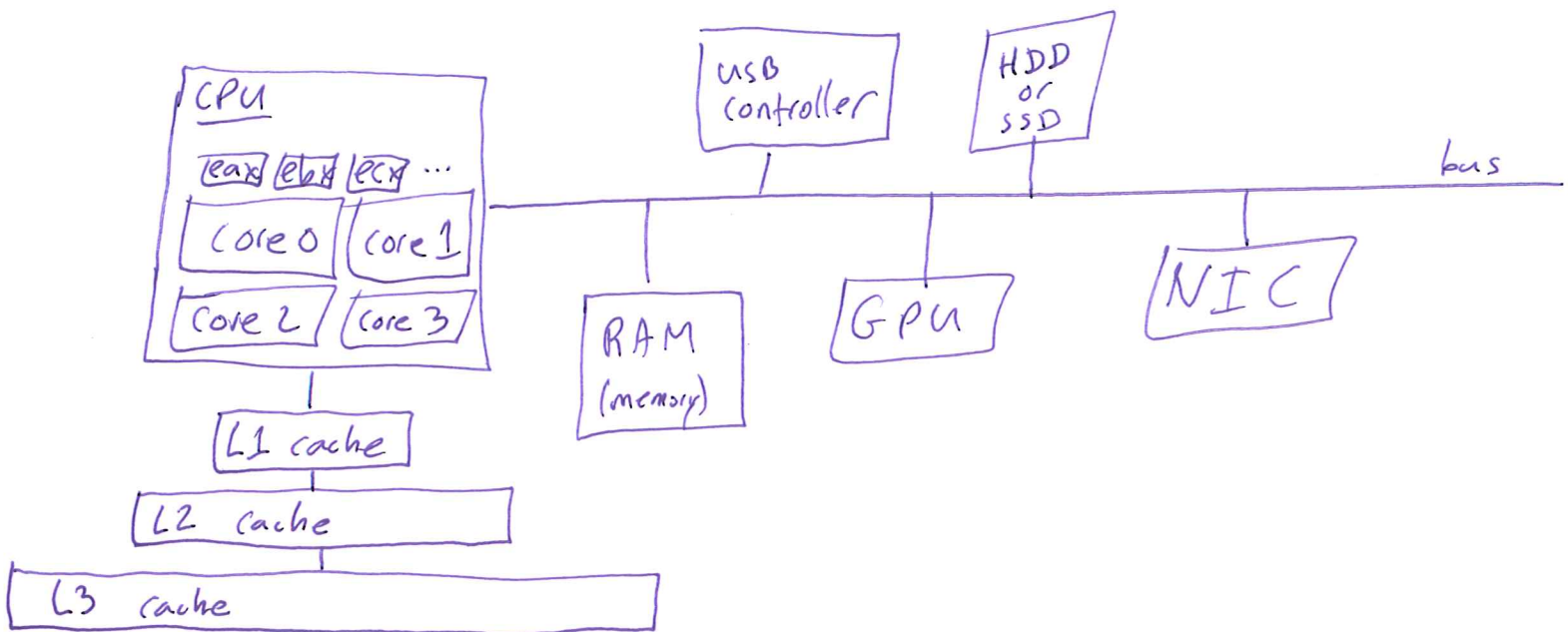# Parallel Programming

Let's draw a cartoon of our computer:



Notes:
- <u>Usually</u> each core has its own $L1$ and $L2$ cache.
- $L3$ is <u>usually</u> shared by all cores.
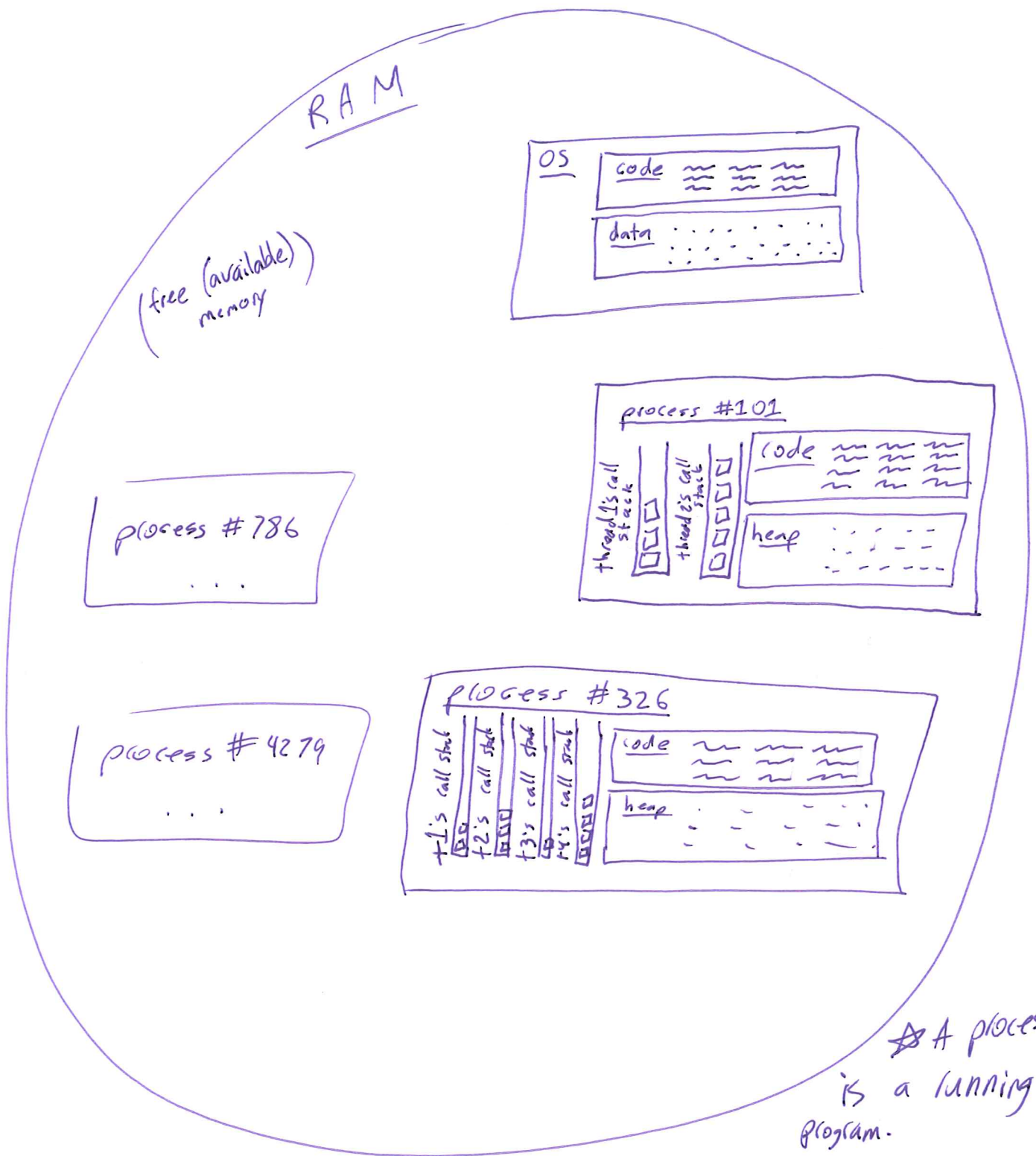- On an Intel i7 ~~CPU~~ CPU:

    - $L1$: 32kB
    - $L2$: 256kB
    - $L3$: 8MB

- Every program is bottlenecked by one of these components!
  - ↳ CPU-bound or IO-bound

Let's zoom in on the RAM and draw another cartoon.

RAM

OS | code ~~~ ~~~ ~~~
data · · · · ·

(free (available) memory)

process #786
. . .

process #4279
. . .

process #101
thread 1's call stack | thread 2's call stack | code ~~~ ~~~ ~~~ ~~~ ~~~
heap · · · · ·

process #326
t1's call stack | t2's call stack | t3's call stack | t4's call stack | code ~~~ ~~~ ~~~ ~~~ ~~~
heap · · · · ·

☆ A process is a running program.

Notes:
- The OS is just code and data just like any other program.
  → The os juggles all the threads and juggles itself in the mix!
- One of the OS's jobs is to schedule threads onto CPU cores.

Now we'll do some programming.

First, let's see how many cores, how much cache, and how much memory our Mac has.

Run 0.py → run do_stuff() from main thread
                      → check Activity Monitor

      → lanch two threads for do_stuff()

           → check A.M.

      → lanch 4 threadds ....

           → check A.M.

      → Stop using numpy → just in put a while true:pass

           → check A.M → what happened? ... the FIL.

      → Switch to using processes.

           → check A.M.

      → Promise to show a threading gotcha and a multi-processing gocha.

Run 1.py → mention that printing the PID is often useful so that you know which process to kill if needed.

Run 2.py → Show how to pass arguments to threads, and how to join w/ threads.

  → Piont out that the PID is the same for all threads in a given process.

Run 3.py → look at what ~~is~~ doing... looks right?
3.serial.py then 3.threads.py

  → Run repeatedly. What happens? → This is the Threading gotcha!!

  → Let's switch to using processes. What happens? This is the multi processing gotcha!!!!
    → Explain fork().

  → Fix it to use a queue for results.
    → live code it.

  → Compare runtimes of 3_*.py.

# Final remarks:

- How many threads / processes should you launch?

   A: the # of cores you have → look at your # core options on Aws.

- Why?

   A: - because there's an overhead to context-switching between threads
      ↳ the more threads, the more context switching is needed

   - Also, cache hits/misses. More misses w/ more threads.

- In Python, you'll probably choose processes over threads.

- Don't prematurely optimize. See 4.py for Donald Knuth's quote.

- This is only scratching the surface of parallel programming. Here are other technologies we won't touch, but that are widely used:

   - MPI: a very generic (flexible) distributed programming communication framework.

   - OpenMP: for parallelizing loops. (used in numpy very extensively)

   - GPU: a single GPU has thousands of cores!

# Final Remarks (continued):

- Parallel programming is **hard**. Threads make it easy to corrupt memory when you don't lock critical areas properly. Processes have large overhead in communication — it's easy to actually **slow** down your program if you are communicating between processes inefficiently.

- There are times when it's easy to parallelize your code. In cases where your threads/processes don't need to communicate much, then it's very easy to parallelize. ~~These~~ Programs of this type are called "embarrassingly parallelizable". Eg. grid search. E.g. K-fold cross-validation

- Let's consider grid search... and let's consider using aws. Say one train-test iteration with your model takes 70 seconds. You want to search a hyper-parameter space of 5x5x5 (grid search). But you want this to be finished in 20 minutes (you're leaving for lunch and will be back in 20 minutes and you want the search to be finished). How many cores do you need?
  ↳ |8 cores| ← buy on ~~aws~~ AWS. How much will that cost you?

# Final Remarks (con't):

   — Beware of fork(). It copies all the data in your ~~program~~ current process!

Don't do this:

    1: Load all you data.

    2: Fork

    3: In each process, use part of ~~that~~ the data.

Instead, do this:

    1: Fork

    2: Load only the data each process need.

    3: Use the data...

---

~~Apple~~ Appendix A:

| | Threads | Processes |
|---|---|---|
| quick creation/ destruction | ✓ | |
| easy and fast access to shared mem. | ✓ | |
| robust to bugs/ errors | | ✓ |
| not affected by Python's GIL | | ✓ |