# [544] Processes and Threads

Tyler Caraza-Harter

# Motivation

Modern CPUs have many cores (maybe dozens)

Trend: **more** cores rather than **faster** cores

Problem: a simple Python program can use at most ONE core
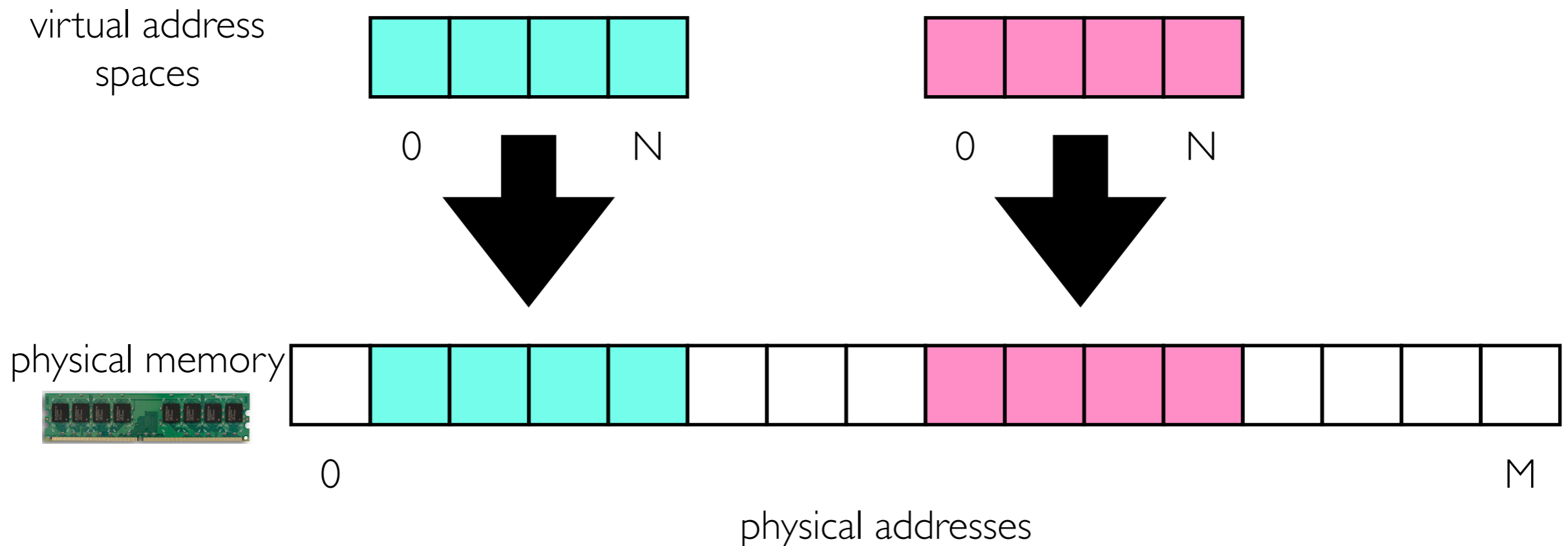(less if it accesses files or the Internet)

Understanding threads and processes will:
- let us write programs that fully utilize CPU resources
- decide the structure of our concurrent program (threads or processes) depending on the situation
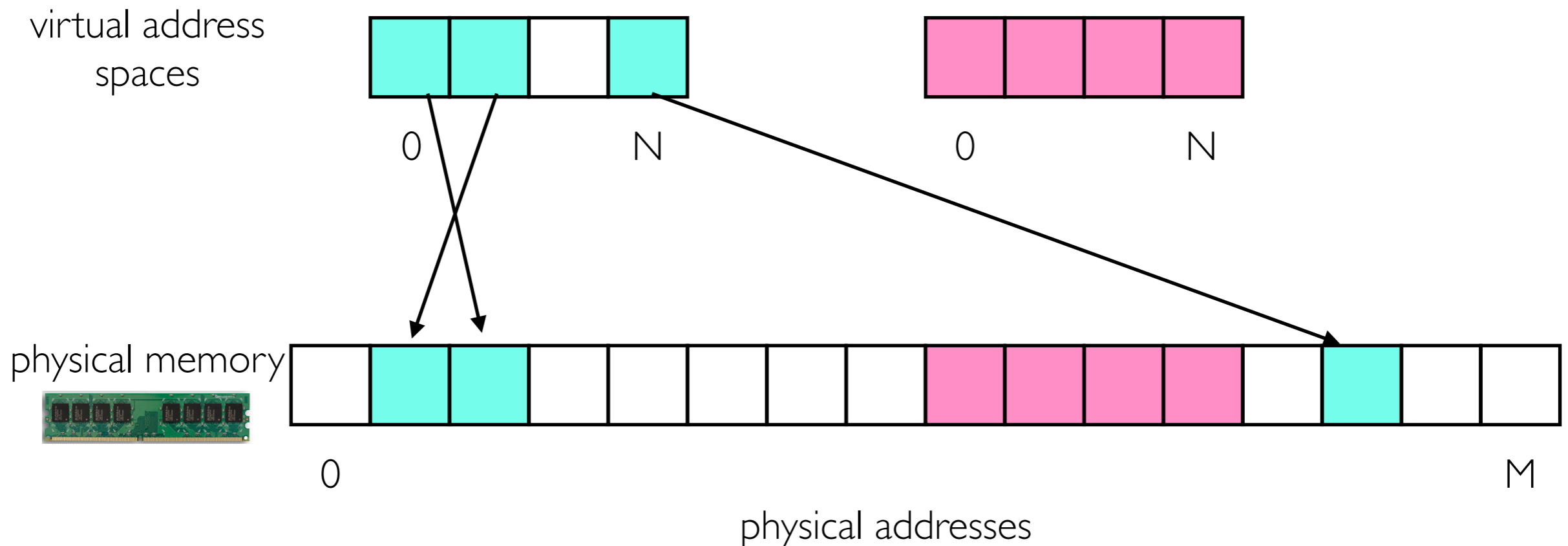
# Processes and Address Spaces

Address spaces

- A process is a running program
- Each process has it's own virtual address space
- The same virtual address generally refers to different memory in different processes
- Regular processes cannot directly access physical memory or other addr spaces

virtual address
spaces

0          N          0          N

physical memory

0                                                                M

physical addresses

# Processes and Address Spaces

Address spaces

- A process is a running program
- Each process has it's own virtual address space
- The same virtual address generally refers to different memory in different processes
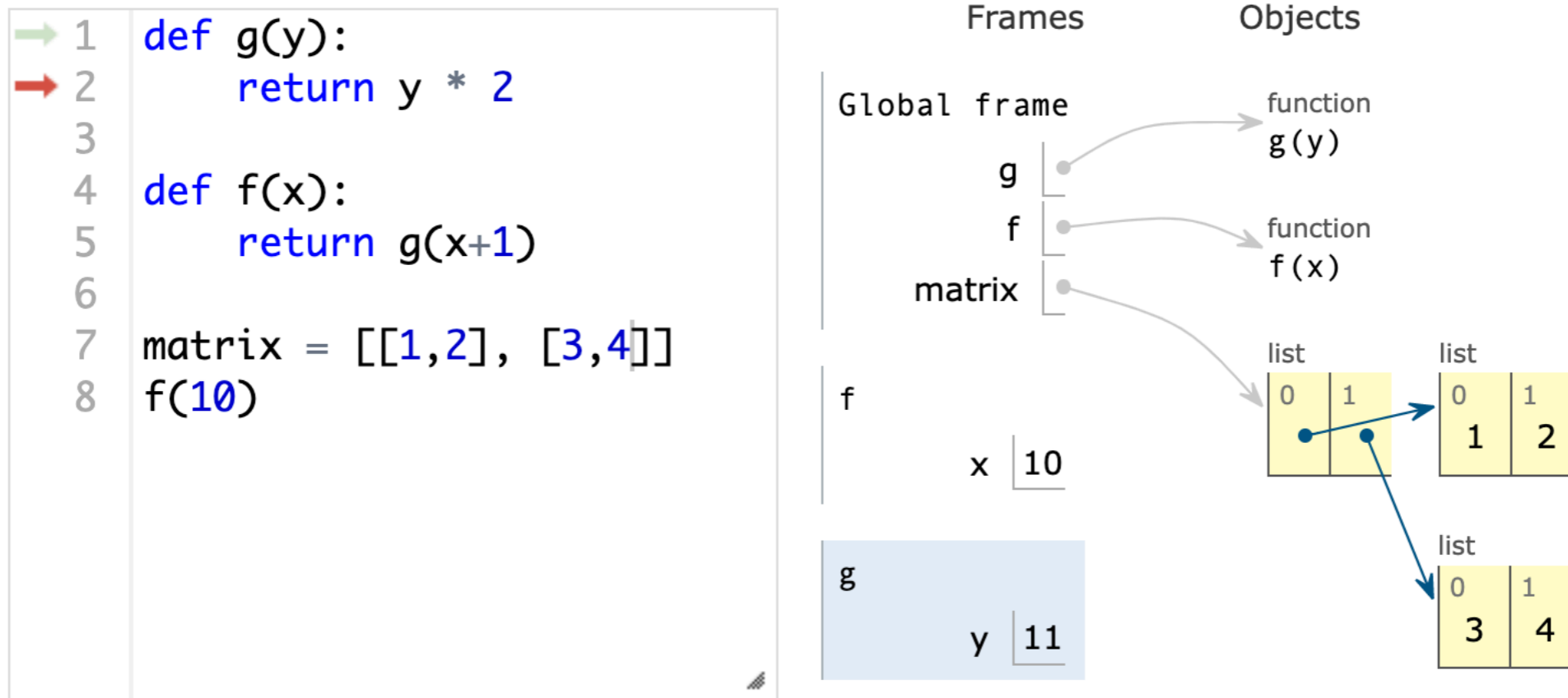- Regular processes cannot directly access physical memory or other addr spaces
- Address spaces can have holes (N is usually MUCH bigger than M)
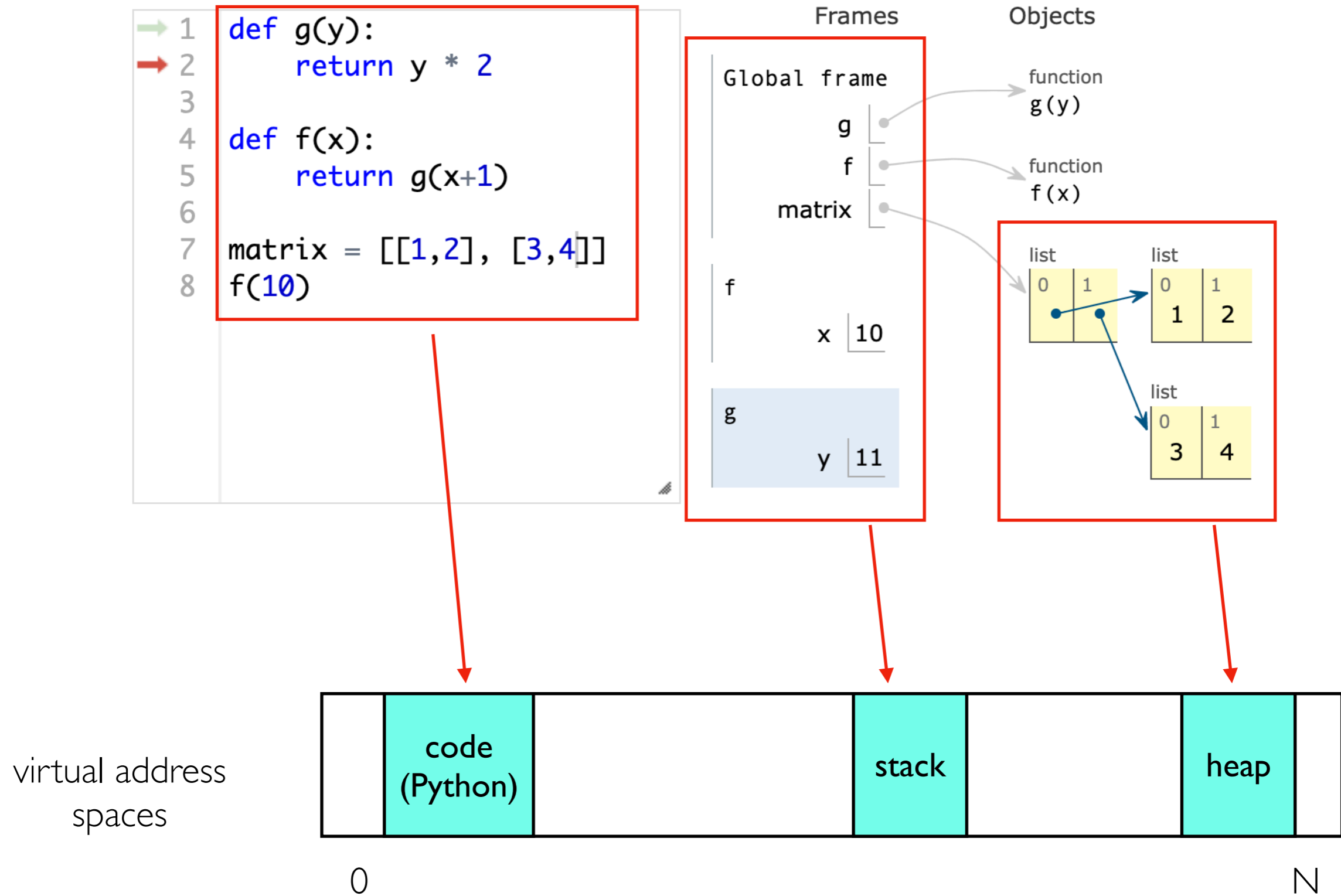- Physical memory for a process need not be contiguous

virtual address
spaces

0      N      0      N

physical memory
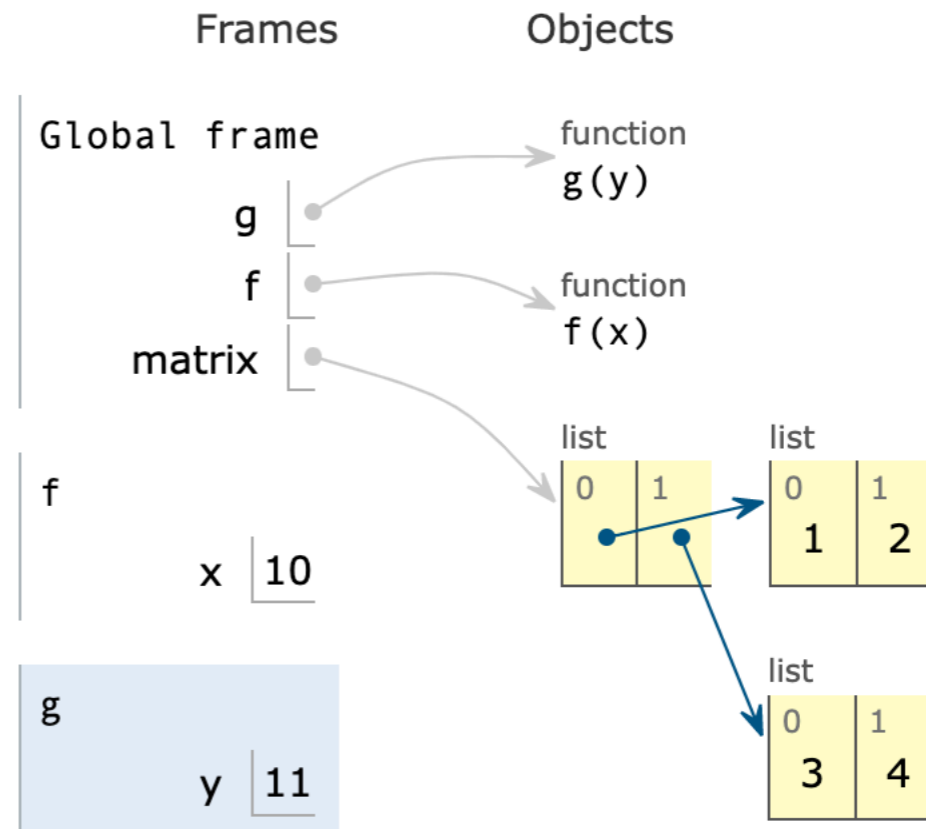
0      M

physical addresses

# What goes in an address space?

```
1  def g(y):
2      return y * 2
3
4  def f(x):
5      return g(x+1)
6
7  matrix = [[1,2], [3,4]]
8  f(10)
```

Frames                Objects

Global frame          function
                      g(y)
        g
        f             function
                      f(x)
    matrix
                      list       list
    f                 0  1       0  1
                                 1  2
        x  10
                                 list
    g                            0  1
                                 3  4
        y  11

https://pythontutor.com/

virtual address
spaces

**what goes here?**
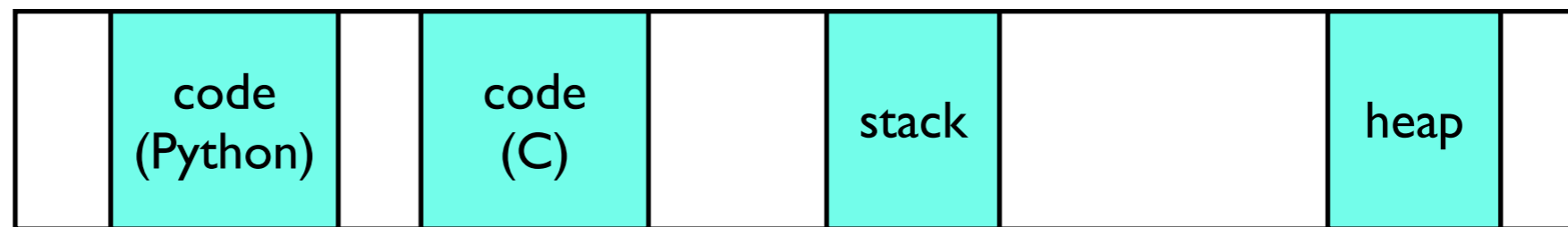
0                                              N

# What goes in an address space?

```
1  def g(y):
2      return y * 2
3
4  def f(x):
5      return g(x+1)
6
7  matrix = [[1,2], [3,4]]
8  f(10)
```

Frames

Objects

Global frame

g → function g(y)

f → function f(x)

matrix

f

x | 10

g

y | 11

list
| 0 | 1 |

list
| 0 | 1 |
| 1 | 2 |

list
| 0 | 1 |
| 3 | 4 |

virtual address spaces

| | code (Python) | | stack | | heap | |

0 ............................................................ N

**Note**: code and heap generally not contiguous

# What goes in an address space?

```
→ 1  def g(y):
→ 2      return y * 2
  3
  4  def f(x):
  5      return g(x+1)
  6
  7  matrix = [[1,2], [3,4]]
  8  f(10)
```

**Frames**

Global frame

g

f

matrix

f

x  10

g

y  11

**Objects**

function
g(y)

function
f(x)

list

| 0 | 1 |
|---|---|

list

| 0 | 1 |
|---|---|
| 1 | 2 |

list

| 0 | 1 |
|---|---|
| 3 | 4 |

virtual address spaces

some packages
(like numpy)

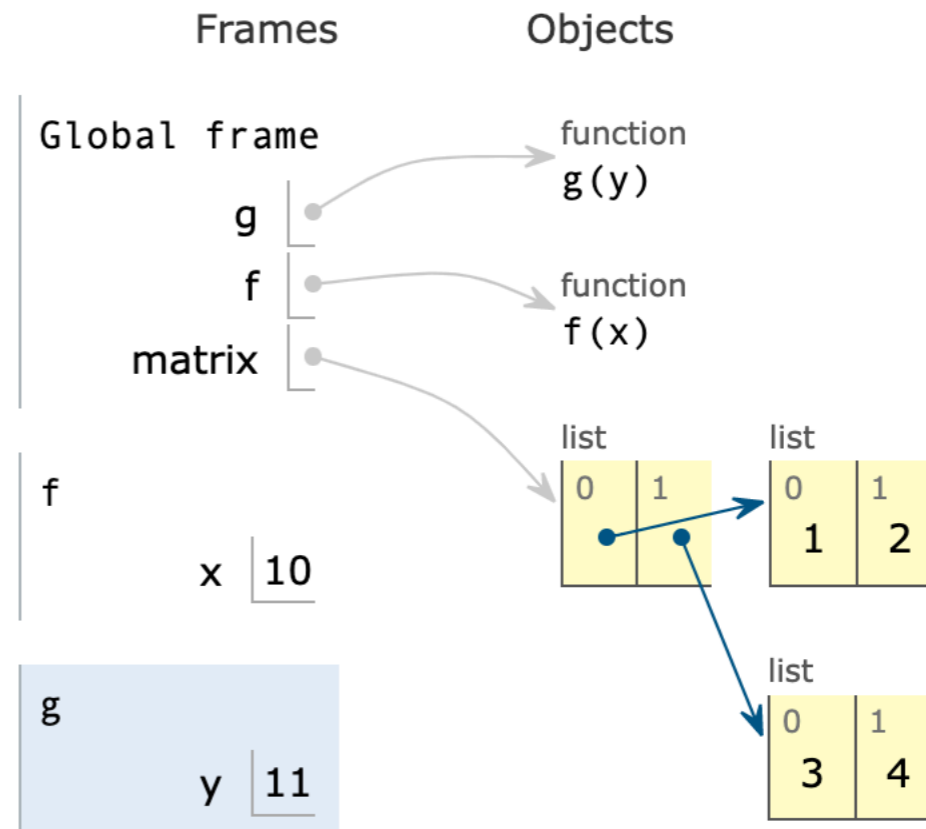| | code (Python) | | code (C) | | stack | | heap | |
|---|---|---|---|---|---|---|---|---|

0

N

# How does code execute?

```python
1  def g(y):
2      return y * 2
3
4  def f(x):
5      return g(x+1)
6
7  matrix = [[1,2], [3,4]]
8  f(10)
```

Frames | Objects

Global frame
- g → function g(y)
- f → function f(x)
- matrix → list [0, 1]

f
- x | 10

g
- y | 11

list
| 0 | 1 |
|---|---|
| 1 | 2 |

list
| 0 | 1 |
|---|---|
| 3 | 4 |

instruction pointer

virtual address spaces

| | code (Python) | | code (C) | | stack | | | heap | |
|---|---|---|---|---|---|---|---|---|---|

0          N

# How does code execute?

CPUs
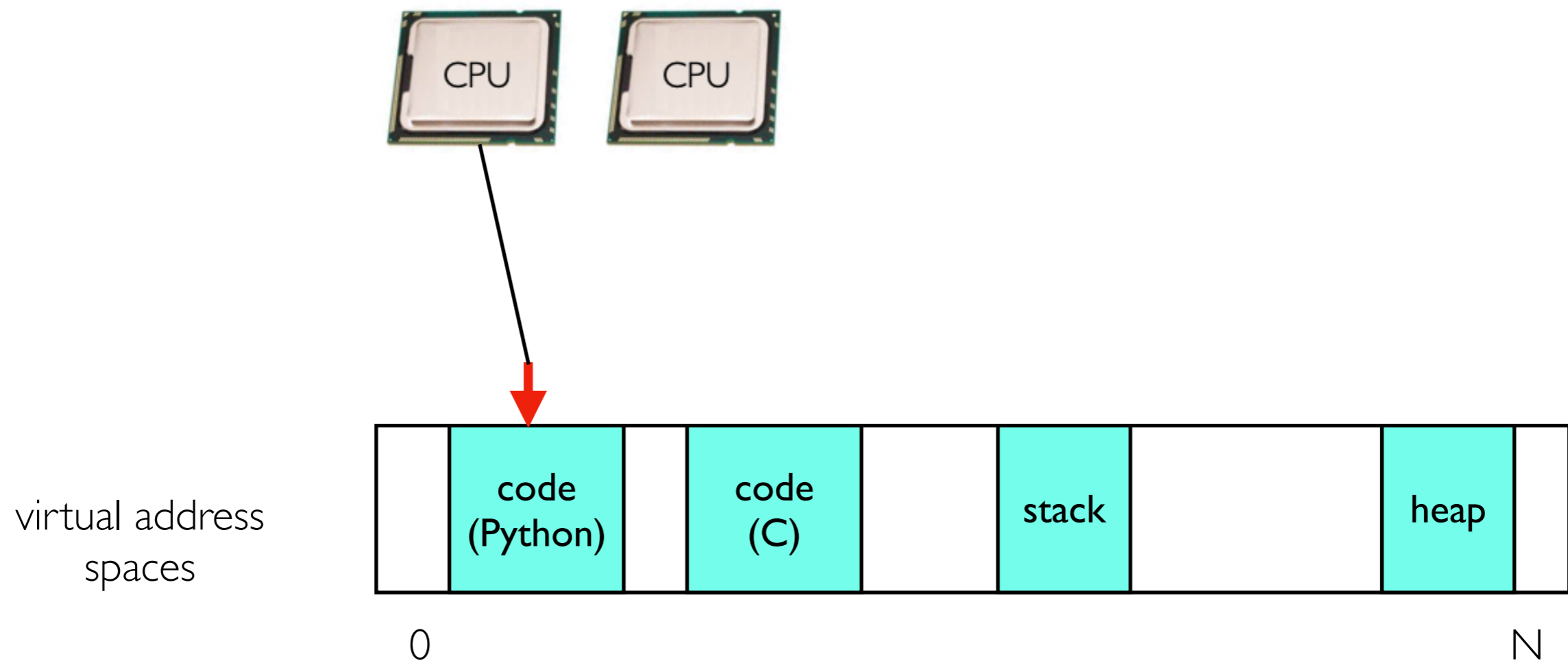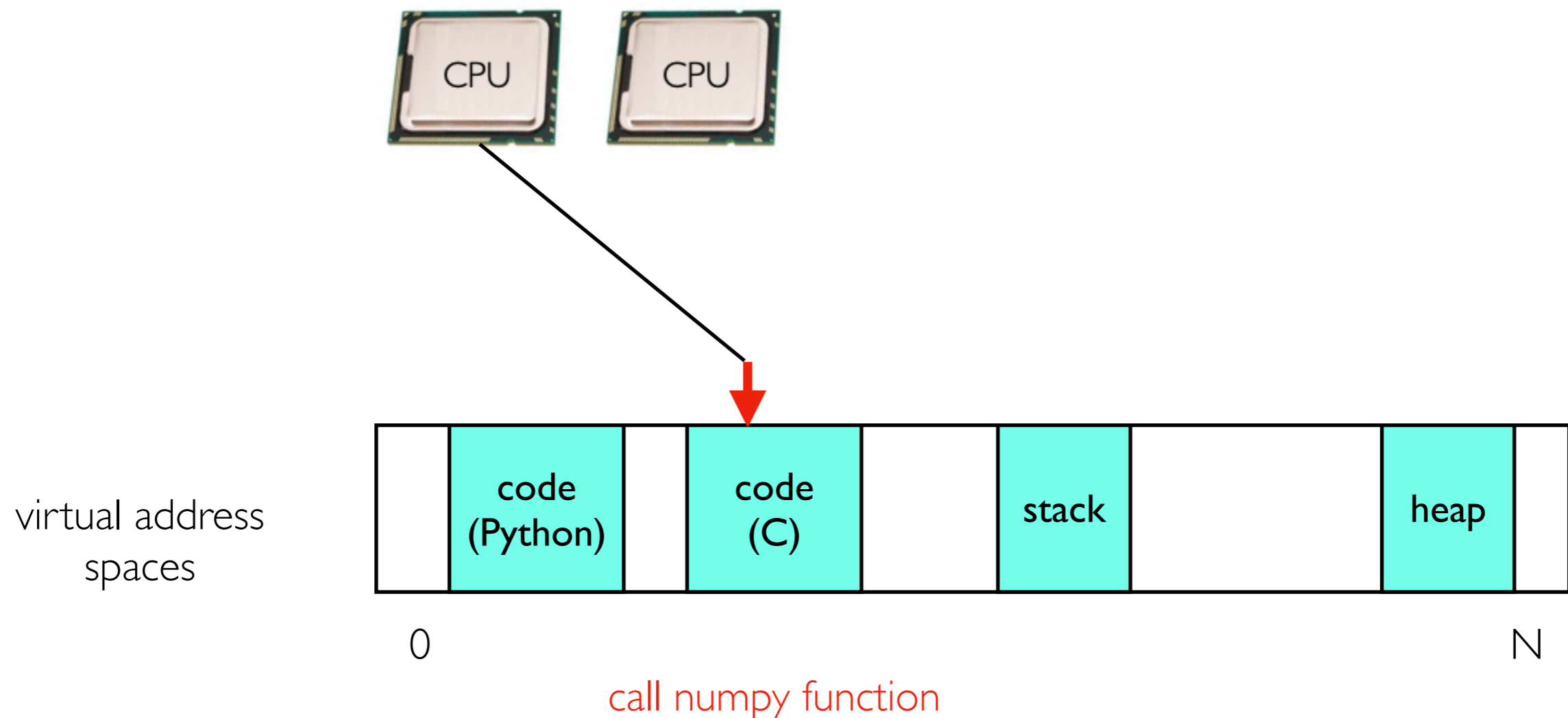- CPUs are attached to at most one instruction pointer at any given time
- they run code by executing instructions and advancing the instruction pointer
- **Note**: interpreter left out for simplicity (CPU points to interpreter code, which points to Python bytecode)



virtual address spaces

| | code (Python) | | code (C) | | stack | | heap | |

0                                                                                    N
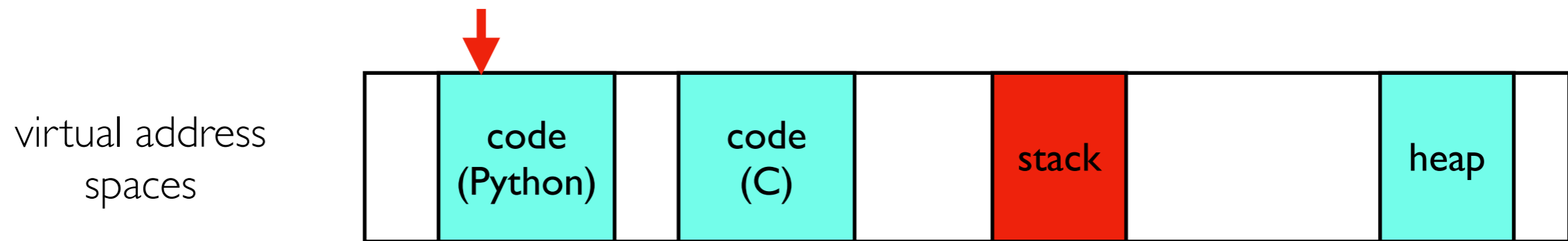
# How does code execute?

CPUs
- CPUs are attached to at most one instruction pointer at any given time
- they run code by executing instructions and advancing the instruction pointer
- **Note**: interpreter left out for simplicity (CPU points to interpreter code, which points to Python bytecode)



virtual address spaces

| | code (Python) | | code (C) | | stack | | heap | |
|---|---|---|---|---|---|---|---|---|

0                                                                         N

# How does code execute?

CPUs

- CPUs are attached to at most one instruction pointer at any given time
- they run code by executing instructions and advancing the instruction pointer
- **Note**: interpreter left out for simplicity (CPU points to interpreter code, which points to Python bytecode)
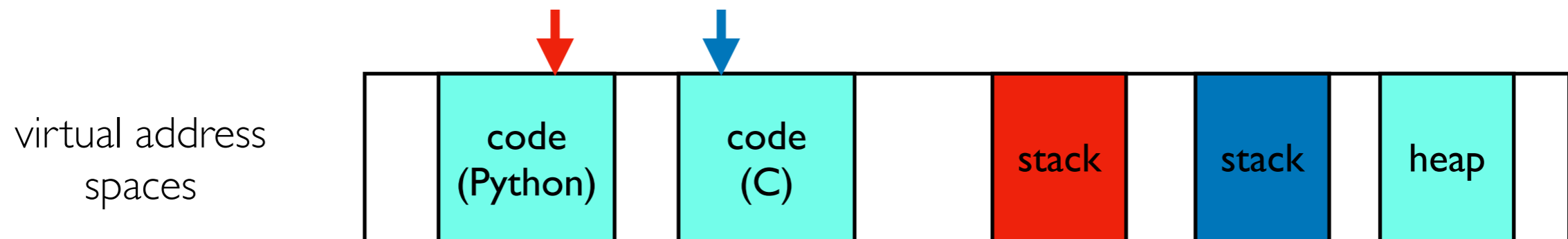


virtual address spaces

| | code (Python) | | code (C) | | stack | | heap | |
|---|---|---|---|---|---|---|---|---|

0                                                                                                    N

call numpy function

# Threads

Threads have their own instruction pointers and stacks, but share the heap.
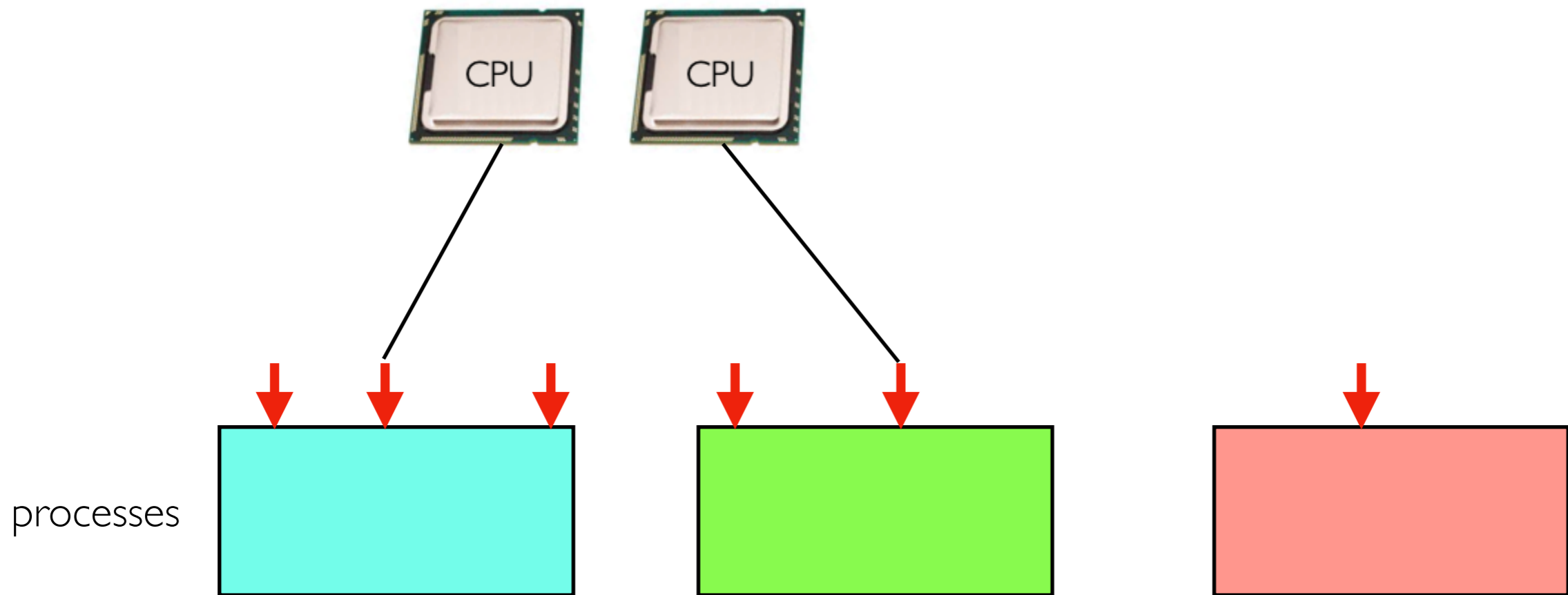
**Single-threaded** process:

virtual address spaces

| | code (Python) | | code (C) | | stack | | heap | |
|---|---|---|---|---|---|---|---|---|

**Multi-threaded** process:

virtual address spaces

| | code (Python) | | code (C) | | stack | | stack | | heap | |
|---|---|---|---|---|---|---|---|---|---|---|

# Context Switch

Schedulers
- CPU scheduler is an important sub system in an operating system
- schedulers decide when to context switch between threads
- context swich: change which thread a CPU is running



processes

# Context Switch

Schedulers

- CPU scheduler is an important sub system in an operating system
- schedulers decide when to context switch between threads
- context swich: change which thread a CPU is running



**context switch!**
same process, diff thread

**context switch!**
thread in diff process

processes

# Scheduling Restrictions: Blocked Threads

Threads can be in one of three states

- **running**: CPU is executing it

- **blocked**: waiting on something other than CPU (network, input, disk, etc)

- **ready**: scheduler can choose to context switch to it

CPU cannot advance instruction pointer until network request finishes

```
r = requests.get(URL)
total = sum(r.json())
print(total)
```

running          ready                              running

blocked

processes

# Efficient Use of Compute Resources

**Wasted cores:** (1) not enough threads (2) blocked threads

For 100% CPU utilization (difficult goal)
- need at least one ready/running thread for each CPU core
- generally need more threads than cores (threads are often blocked)
- **threads could be in one process (or many)**

Multi-threaded applications
- good when multiple threads need to access frequently modified data structures
- new kinds of bugs possible (race conditions, deadlock)

Multi-process applications
- easier to program (or just manually launch several processes in background)
- better at keeping multiple cores busy simultaneously (Python specific)

Both approaches work well for dealing with blocked threads

# Coding Demos, Worksheet

Thread operations
- t = threading.Thread()
- t.start(target=????, args=[????])
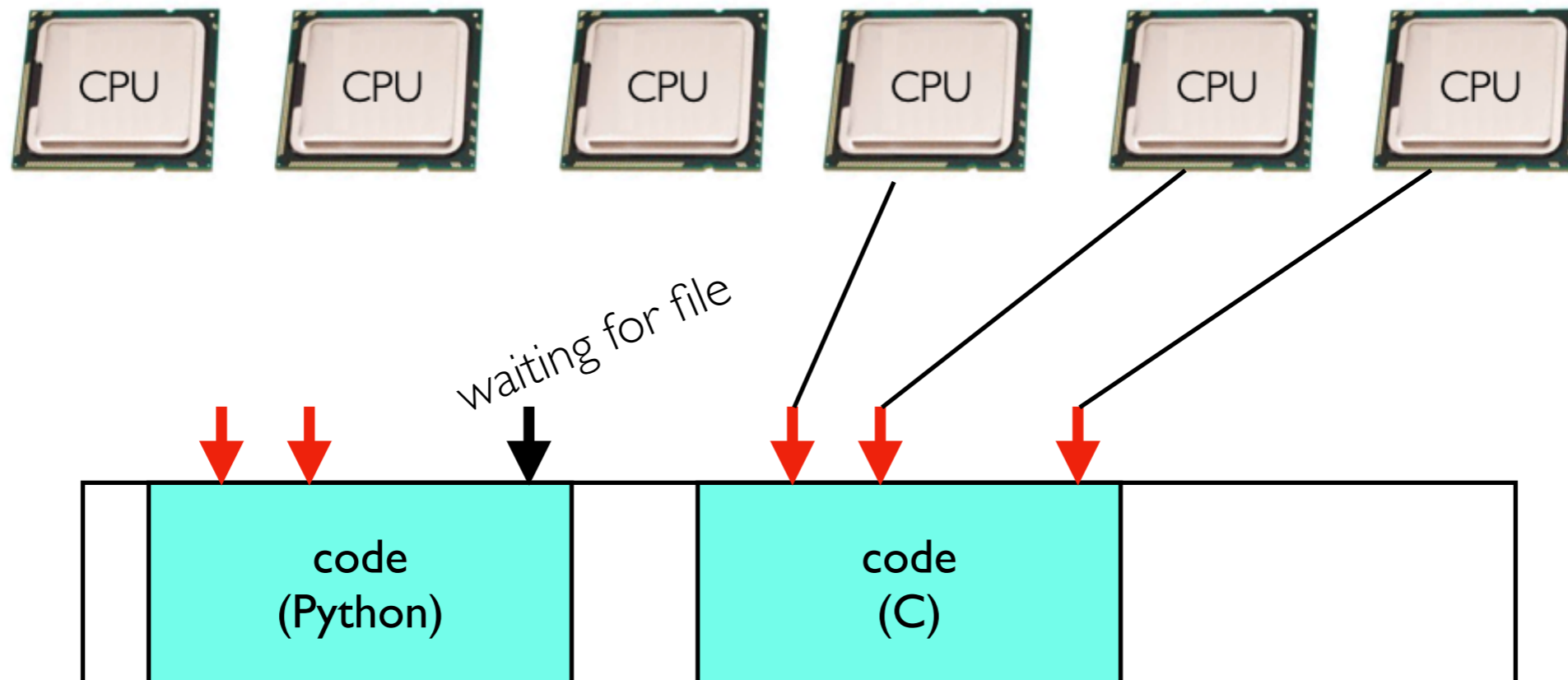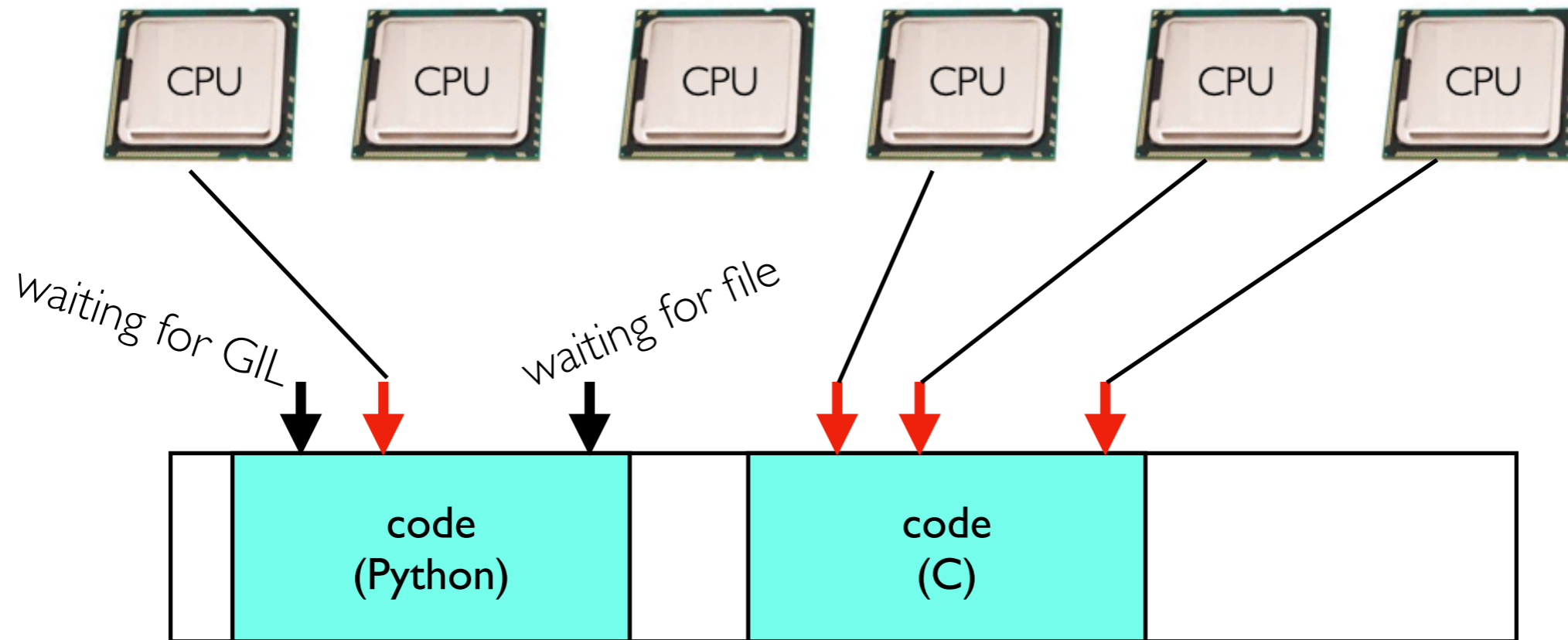- t.join()
- t.get_native_id()

# Python's GIL (Global Interpreter Lock)



Global Interpreter Lock
- Only one thread can be running Python code in a process at once
- Python threads are bad for using multiple cores
- They're still useful for threads blocked on I/O
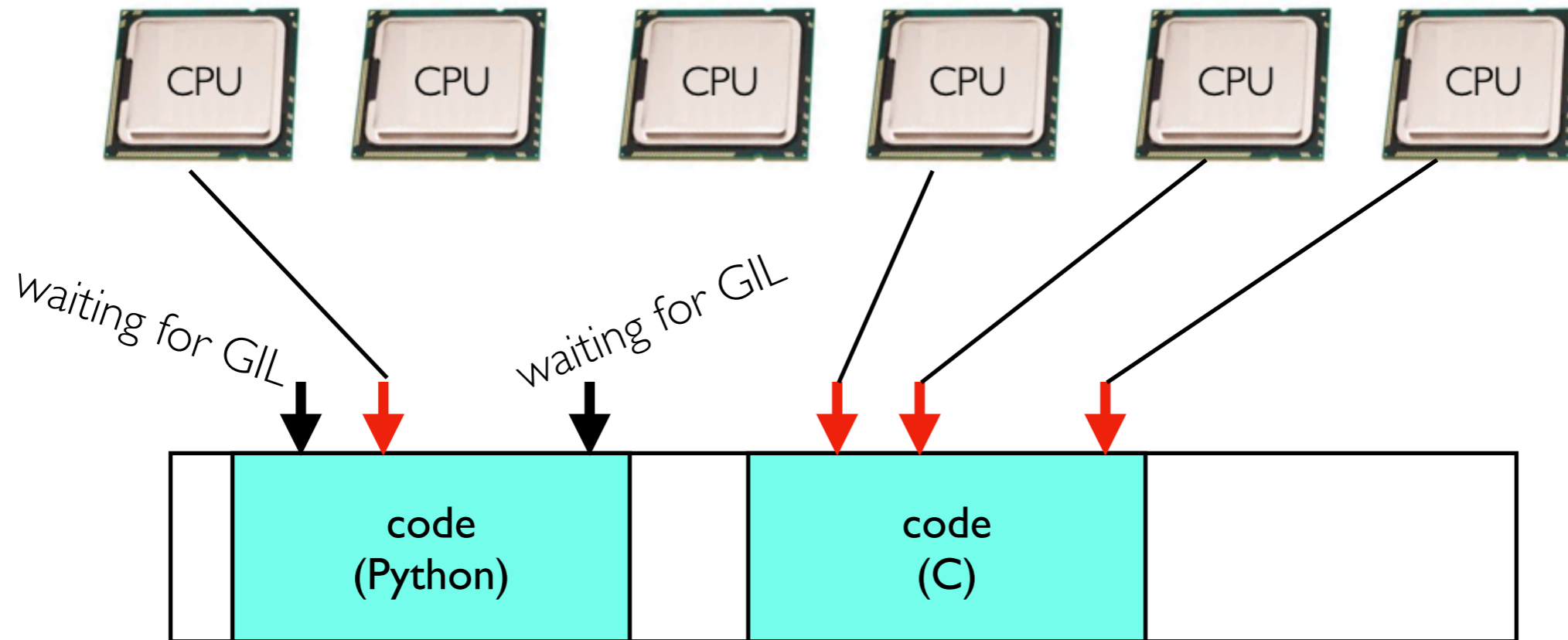- Some Python libraries using other languages allow parallelism

# Python's GIL (Global Interpreter Lock)



Global Interpreter Lock
- Only one thread can be running Python code in a process at once
- Python threads are bad for using multiple cores
- They're still useful for threads blocked on I/O
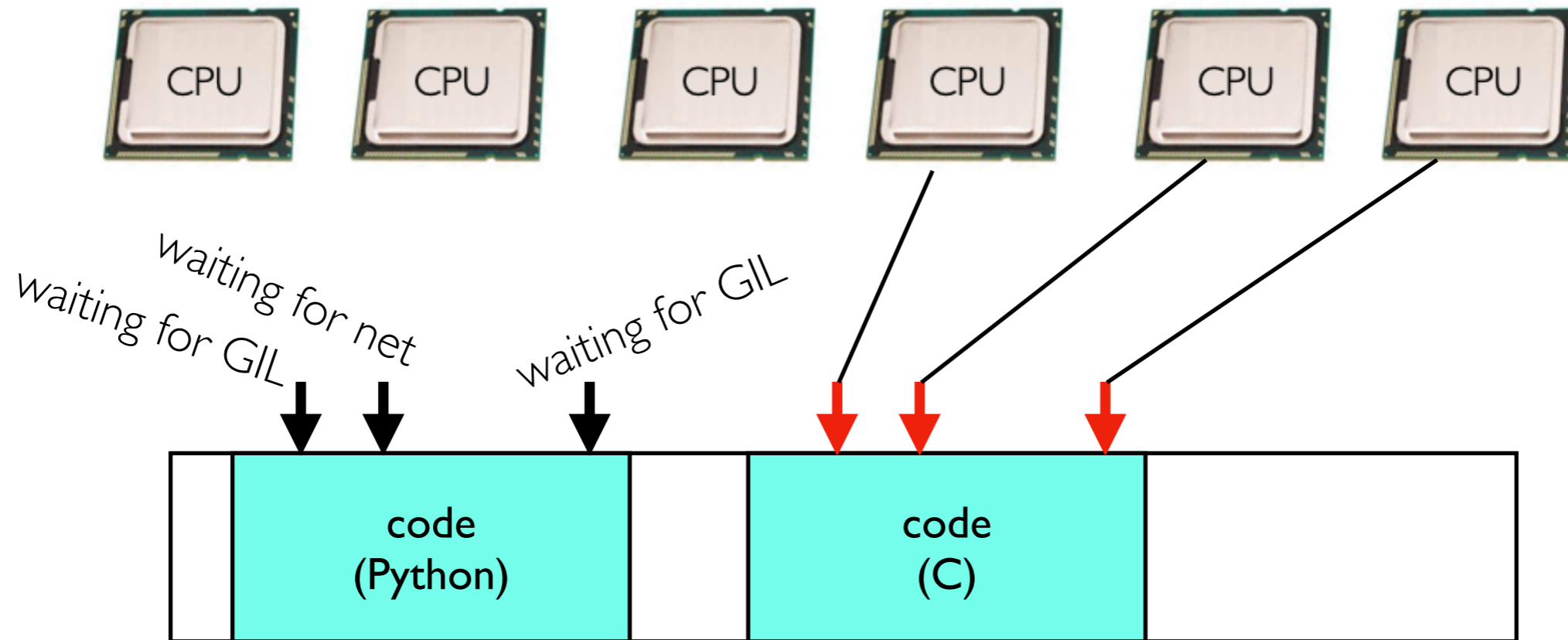- Some Python libraries using other languages allow parallelism

# Python's GIL (Global Interpreter Lock)



Global Interpreter Lock
- Only one thread can be running Python code in a process at once
- Python threads are bad for using multiple cores
- They're still useful for threads blocked on I/O
- Some Python libraries using other languages allow parallelism

# Python's GIL (Global Interpreter Lock)



Global Interpreter Lock
- Only one thread can be running Python code in a process at once
- Python threads are bad for using multiple cores
- They're still useful for threads blocked on I/O
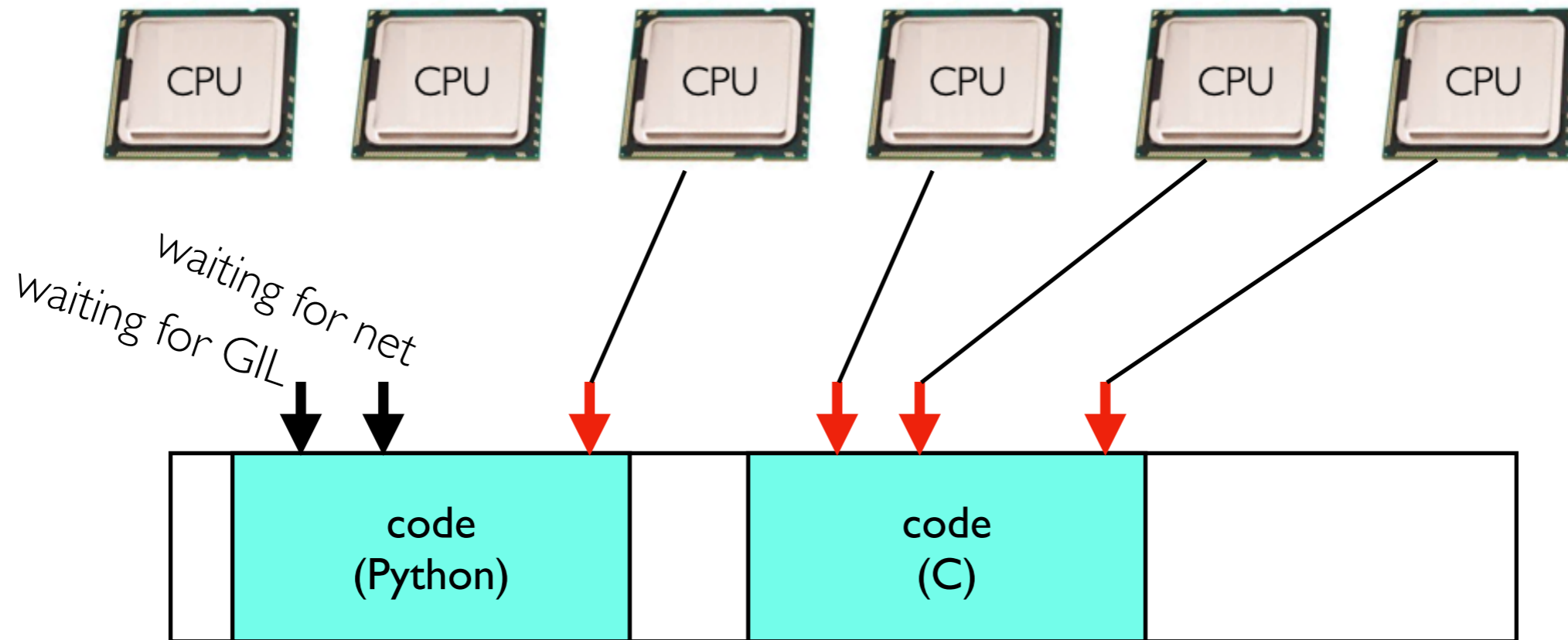- Some Python libraries using other languages allow parallelism

# Python's GIL (Global Interpreter Lock)



Global Interpreter Lock
- Only one thread can be running Python code in a process at once
- Python threads are bad for using multiple cores
- They're still useful for threads blocked on I/O
- Some Python libraries using other languages allow parallelism

# Python's GIL (Global Interpreter Lock)



Global Interpreter Lock
- Only one thread can be running Python code in a process at once
- Python threads are bad for using multiple cores
- They're still useful for threads blocked on I/O
- Some Python libraries using other languages allow parallelism