

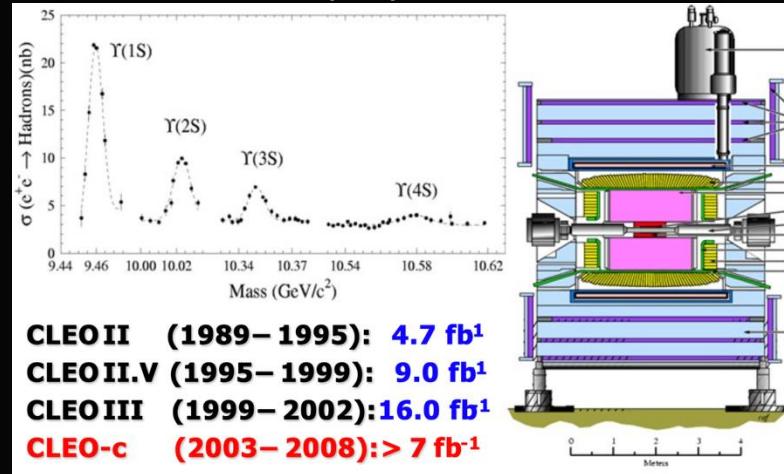
Lightning Introduction to Parallel Programming

David Lange

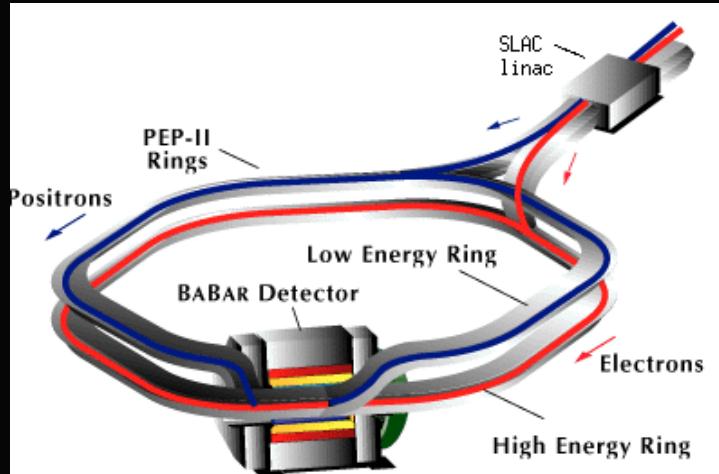
Princeton University

My career in high-energy physics

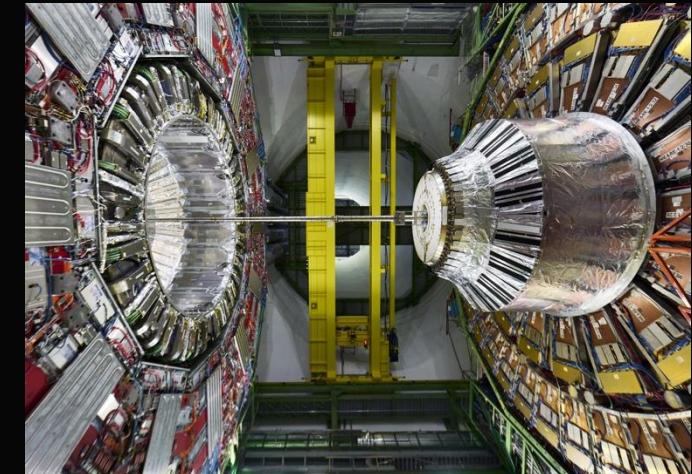
CLEO-II(.5) @ Cornell



BABAR @ SLAC



CMS @ CERN



- Silicon detector calibration + operations
- Form-factor analysis in semileptonic decays
- Event generators

- CP violation analysis
- RPC detector operations + software
- Event generators (“EvtGen”)
- Event reconstruction software

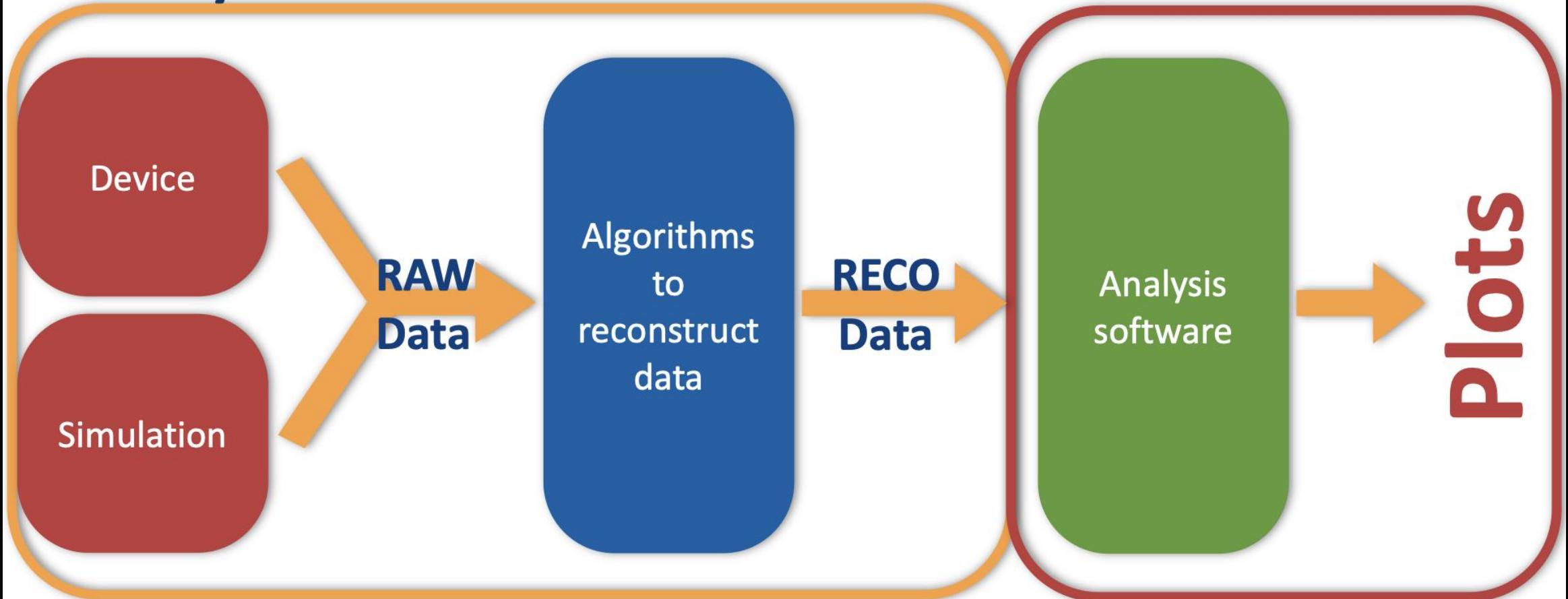
- Event reconstruction software
- Simulation techniques
- Computing resource projections

My Goal for Today

- Why worry about parallel programming for science?
- Considerations for high energy physics codes
- Methods , approaches and hardware
- Languages, tools and resources (likely to be skipped today)
- Using GPUs with Python

Bulk data processing vs time to insight

Analysis in CMS



Central

Hundreds of physicists analyze
the data with different goals
at the same time

How can we get our science done faster or cheaper?

- Reduce the complexity of algorithms in the application
 - Write better code
 - Consider approximations
 - Consider needed/appropriate numerical precision
- Increase the speed and capacity of the compute
 - Increase clock frequency
- Find tasks within the application that can be performed in parallel

How can we get our science done faster?

- Reduce the complexity of algorithms in the application
 - Write better code
 - Consider approximations
 - Consider needed/appropriate numerical precision
- Increase the speed and capacity of the compute
 - Increase clock frequency
- Find tasks within the application that can be performed in parallel

How can we get our science done faster?

- Reduce the complexity of algorithms in the application
 - Write better code
 - Consider approximations
 - Consider needed/appropriate numerical precision
- Increase the number of cores
 - Increase core speed
- Find tasks which can be parallelized

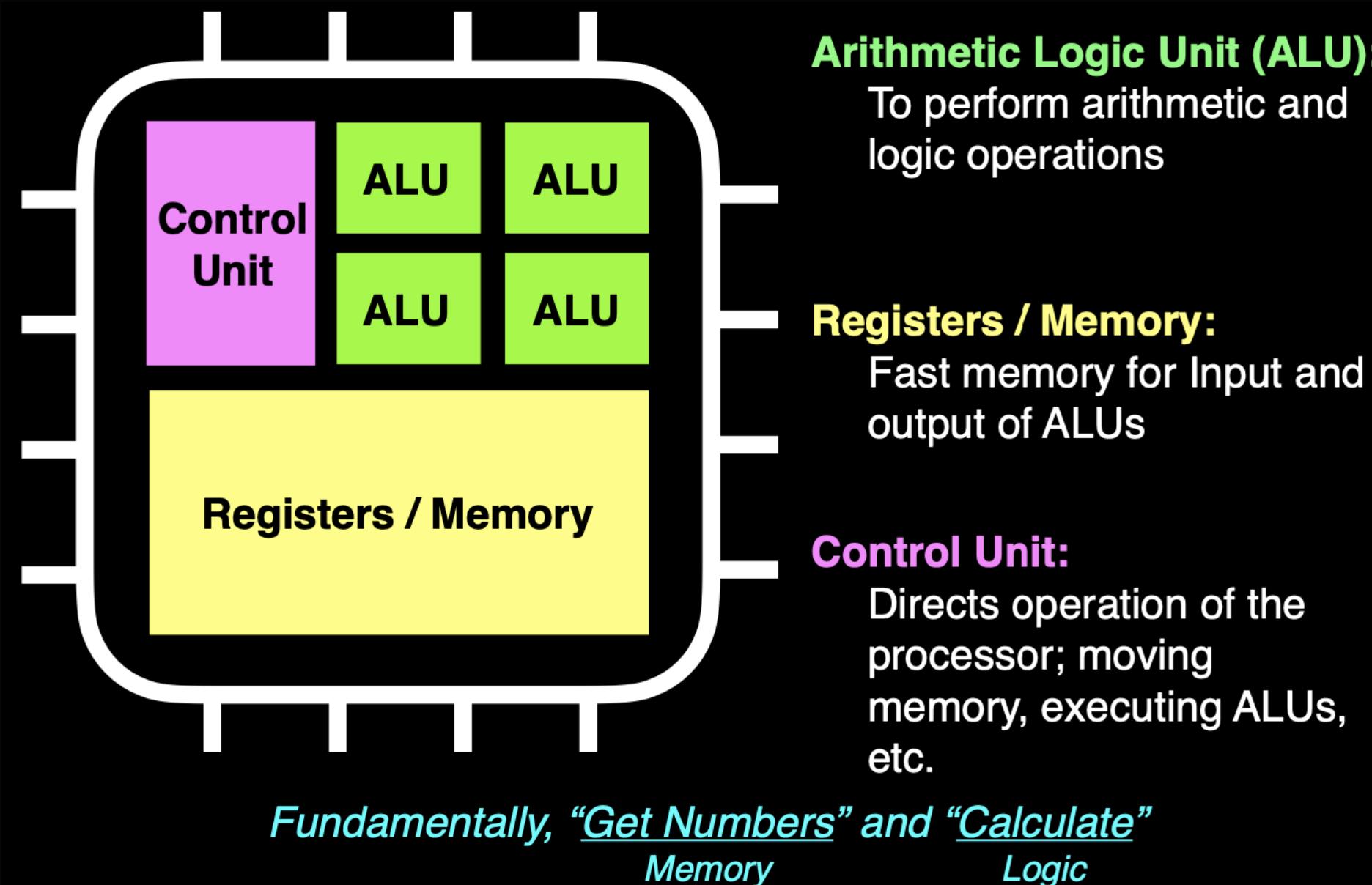
We won't talk about this today

Efficient coding and optimization techniques can easily fill an entire course (or an entire career...)

How can we get our science done faster?

- Reduce the complexity of algorithms in the application
 - Write better code
 - Consider approximations
 - Consider needed/appropriate numerical precision
- Increase the speed and capacity of the compute
 - Increase clock frequency
- Find tasks within the application that can be performed in parallel

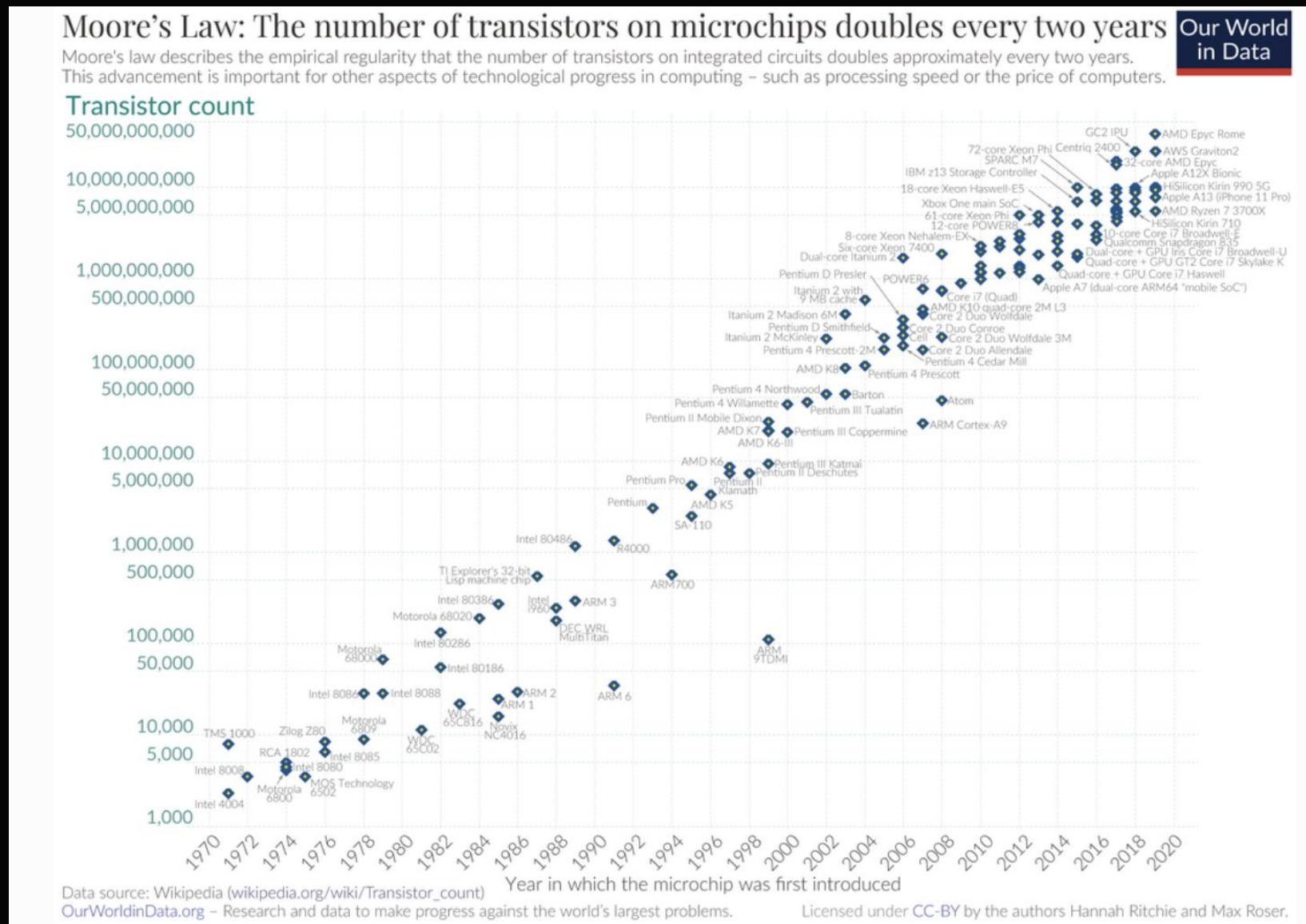
Central Processing Units (CPUs)



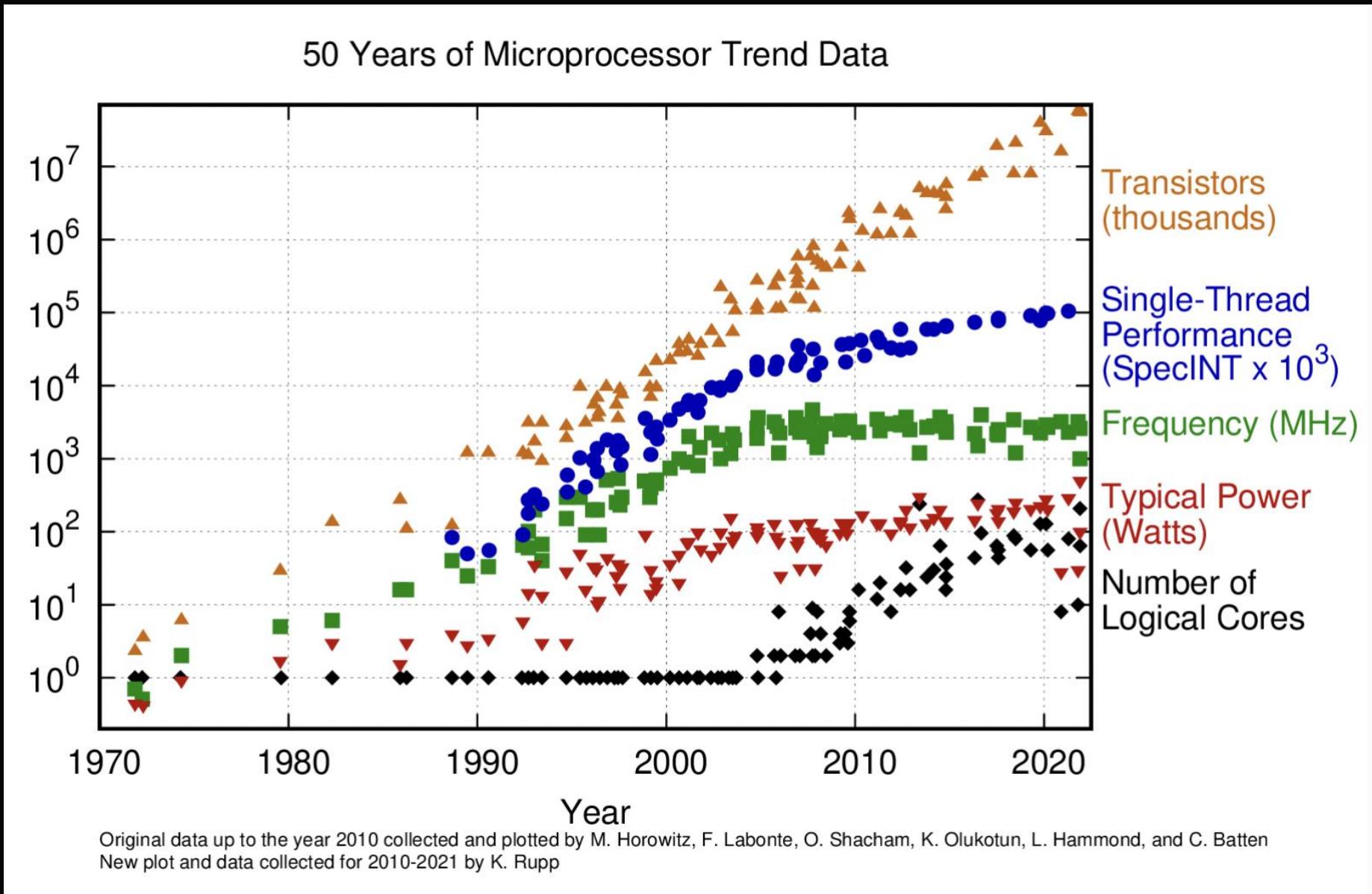
How to go faster?

- *More logic units in same space = more transistors*
- *Faster clocking = higher frequency*
- *I/O performance needs to keep up...*

More transistors? Moore's Law



Trends in transistors to compute architectures



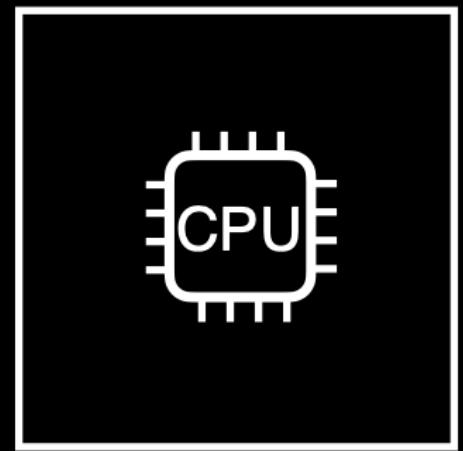
What happened in the early 2000s?

- **Dennard scaling** suggested that each individual transistor in a new generation would be cooler and draw less power
 - Broke down in the early 2000s. This effectively stopped the continuous trend of clock speed increase in CPUs

Why?

- More power per core brings two problems: more power (more money) to run the CPU and higher heat removal requirements (yet more power/money)

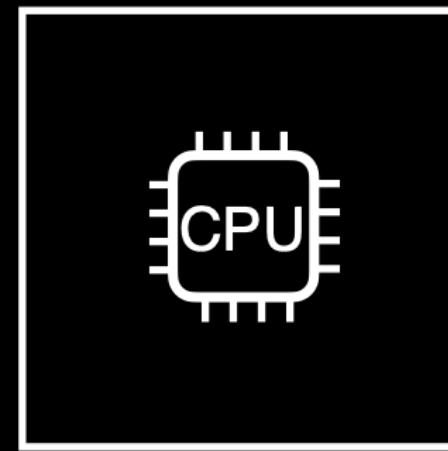
Single core 1 GHz



Power ~ 1W

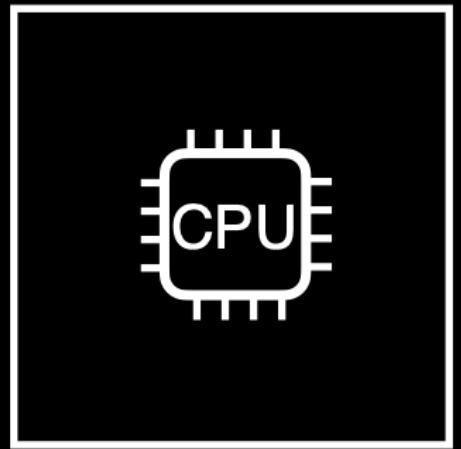
$$Power \sim (freq)^3$$

Single core 4 GHz



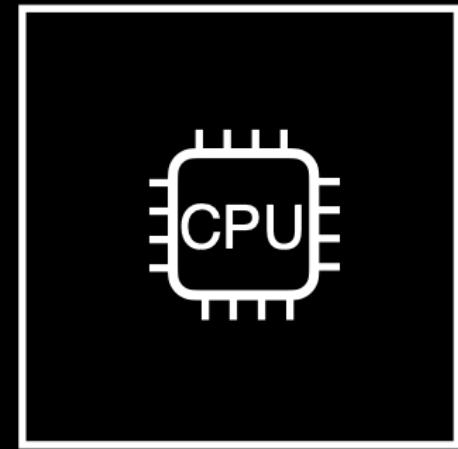
Power ~ 64W

Single core 1 GHz



Power ~ 1W

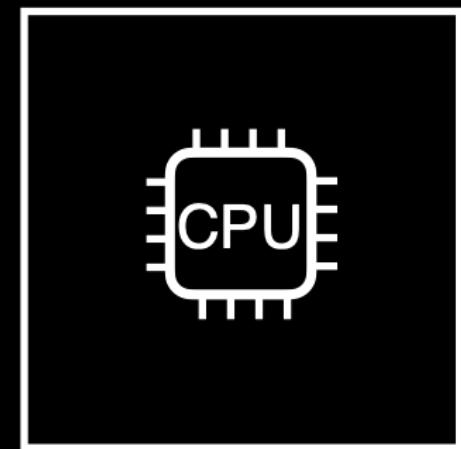
Single core 4 GHz



Power ~ 64W

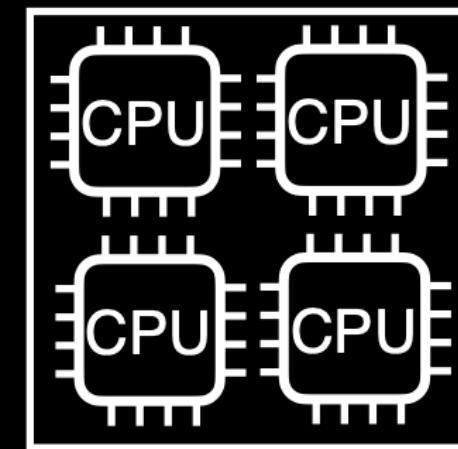
$$Power \sim (freq)^3$$

Single core 1 GHz



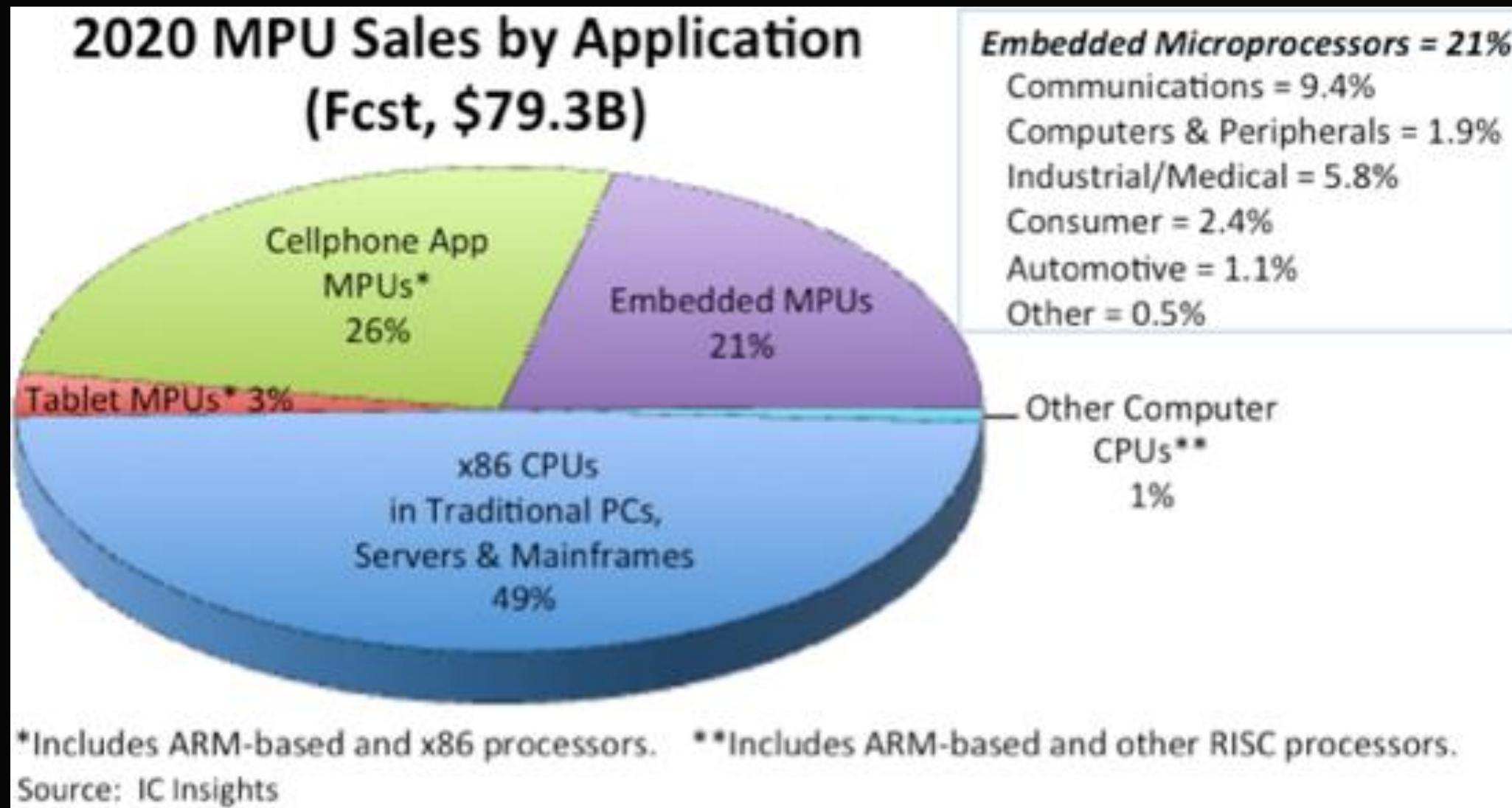
Power ~ 1W

Four cores 1 GHz

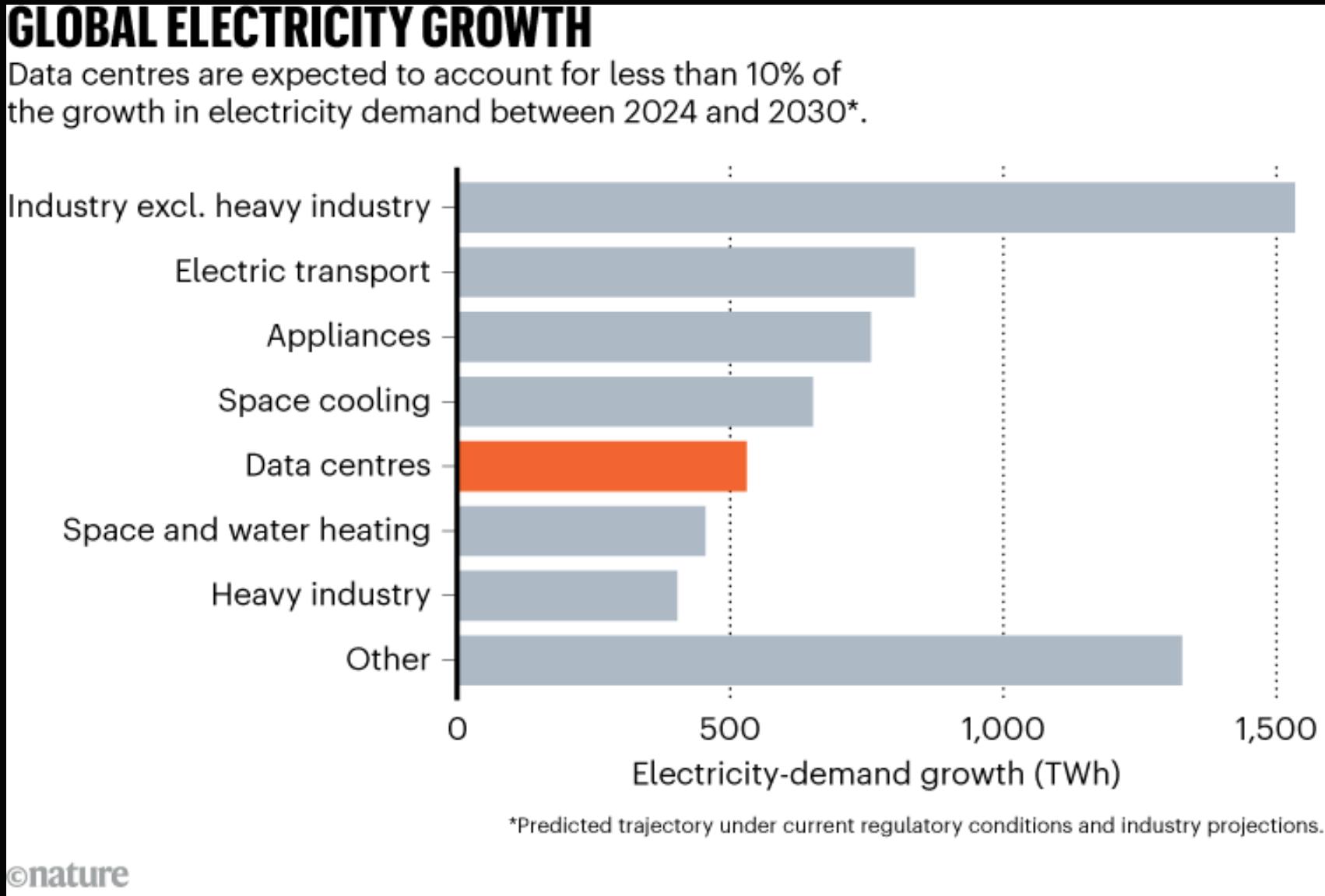


Power ~ 4W

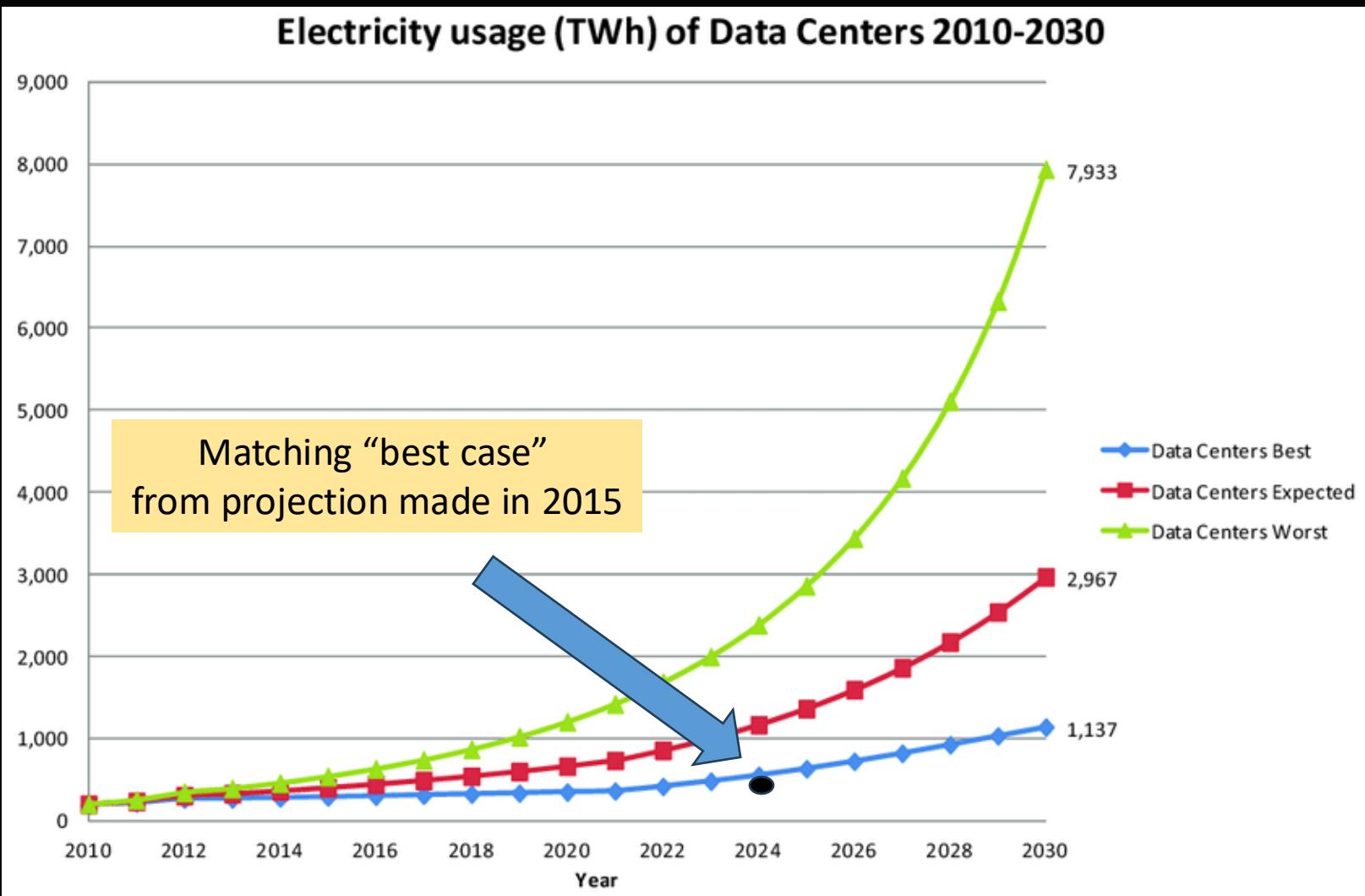
Data centers and low-power applications drive the market



Power is a leading consideration in data centers (too)



Just how important is energy in data centers?



How can we get our science done faster?

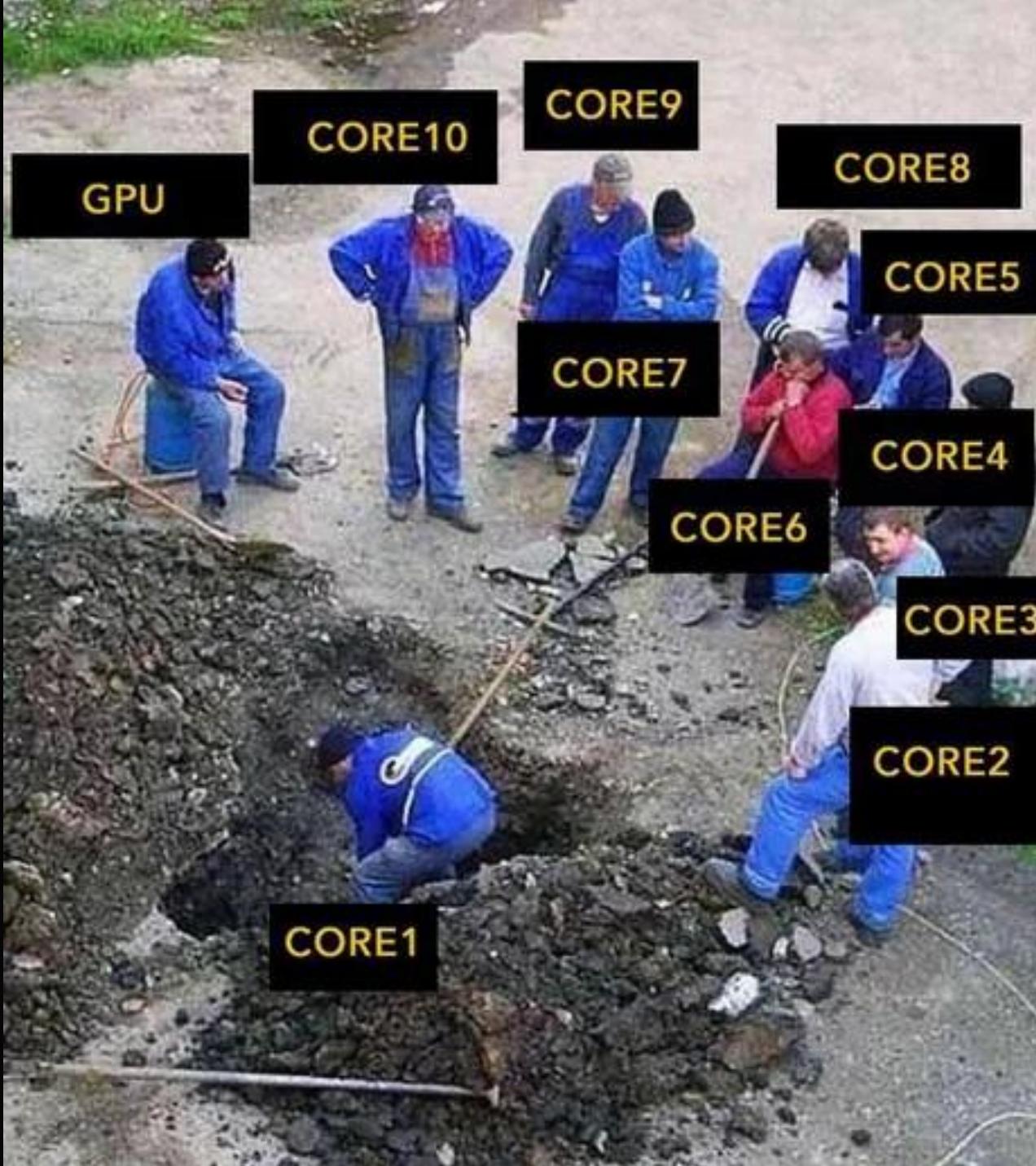
- Reduce the complexity of algorithms in the application
 - Write better code
 - Consider approximations
 - Consider needed/appropriate numerical precision
- Increase the speed and capacity of the compute
 - Increase clock frequency
- Find tasks within the application that can be performed in parallel

Parallelization of certain tasks of our program, so that they are executed simultaneously (threads)

For this it is necessary to:

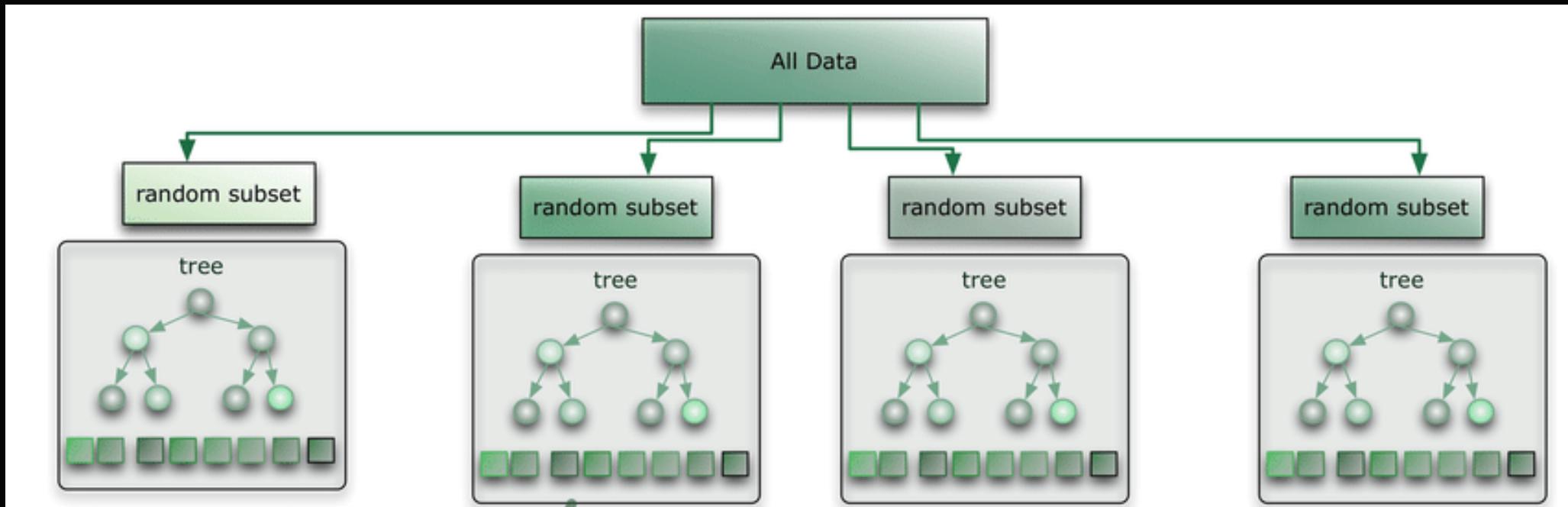
- Have parallel platforms where to run.
 - We just saw that this is not a problem. All modern compute is parallel to some extent
- Adapt your application to be **efficiently** executed in parallel.
 - How easy is this?

Parallel hole digging

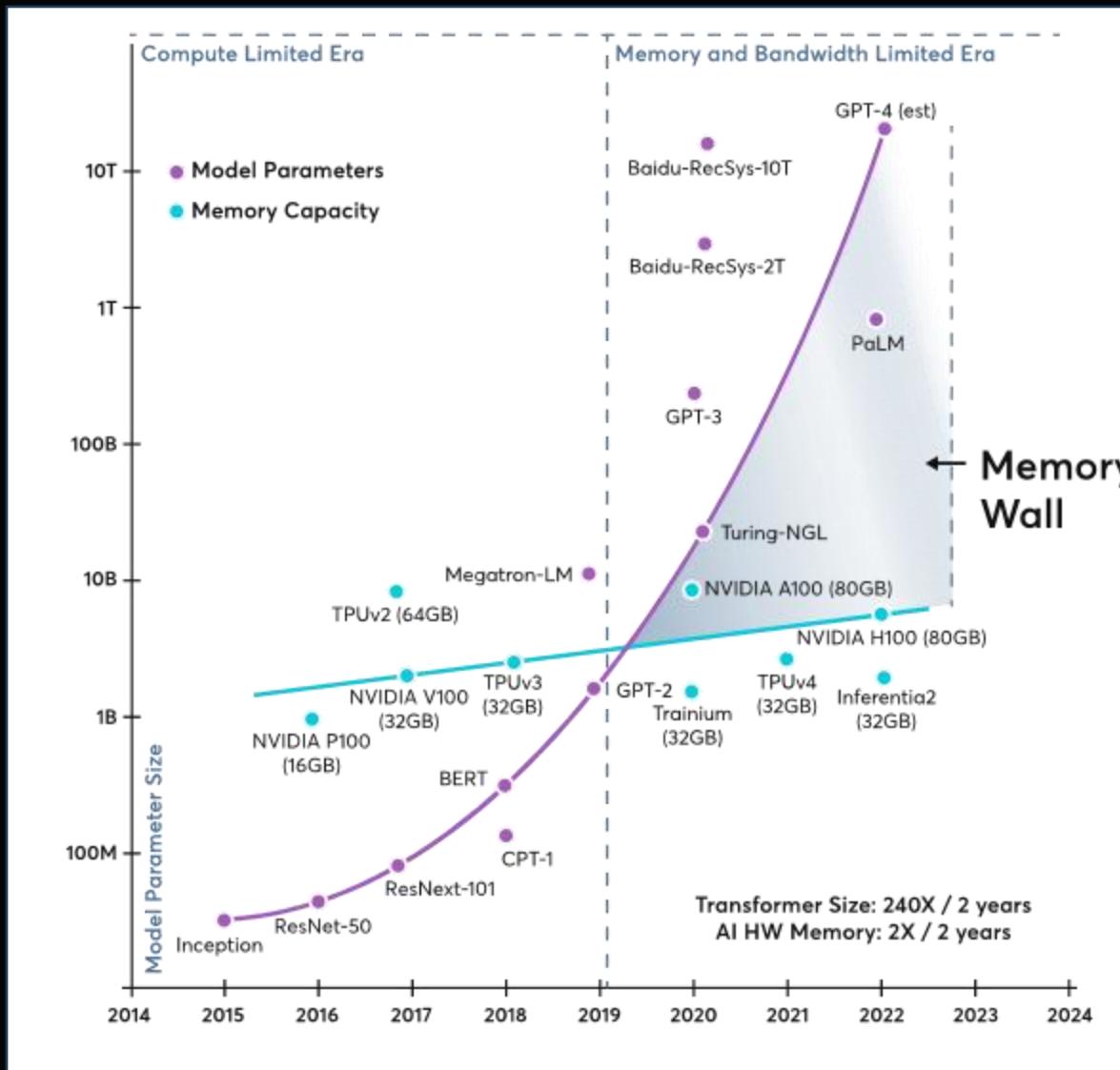


Why not just “embarrassingly” parallel?

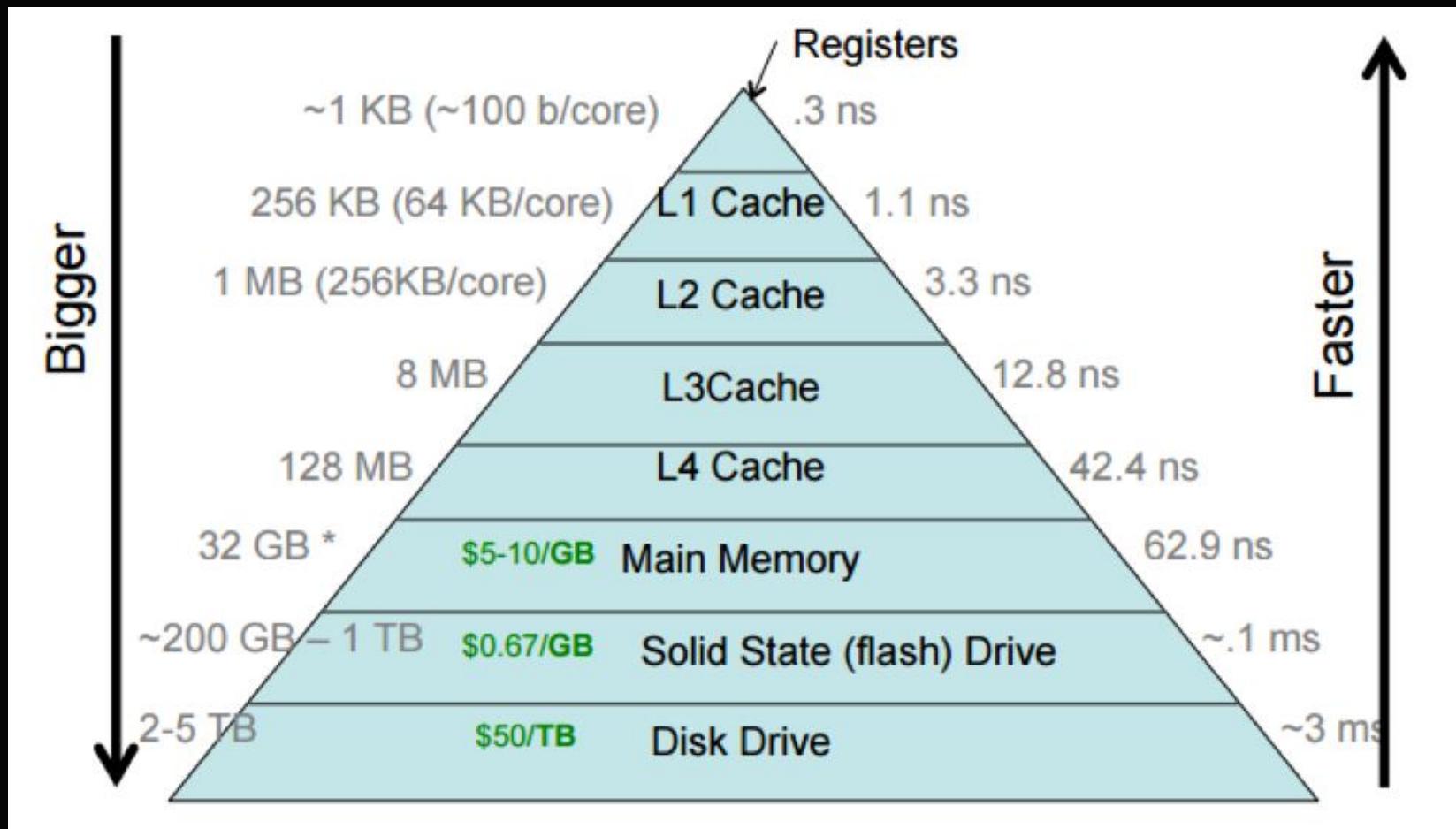
- HEP experimentalists have lots of events(100s of billions per year), and **events can be analyzed independently**.
 - Can't we just write applications that do the equivalent of digging lots of holes?



Memory demands out pace supply



Another way to look at the importance of memory footprint



Memory efficient algorithms will run faster!

What does all this mean?

Science does not drive the market for modern computing, so we must adapt to market trends

To take full advantage of modern computing architectures, we need to think about building applications from algorithms that are capable of exploiting parallelism opportunities that speed up the application and allow it to use less memory

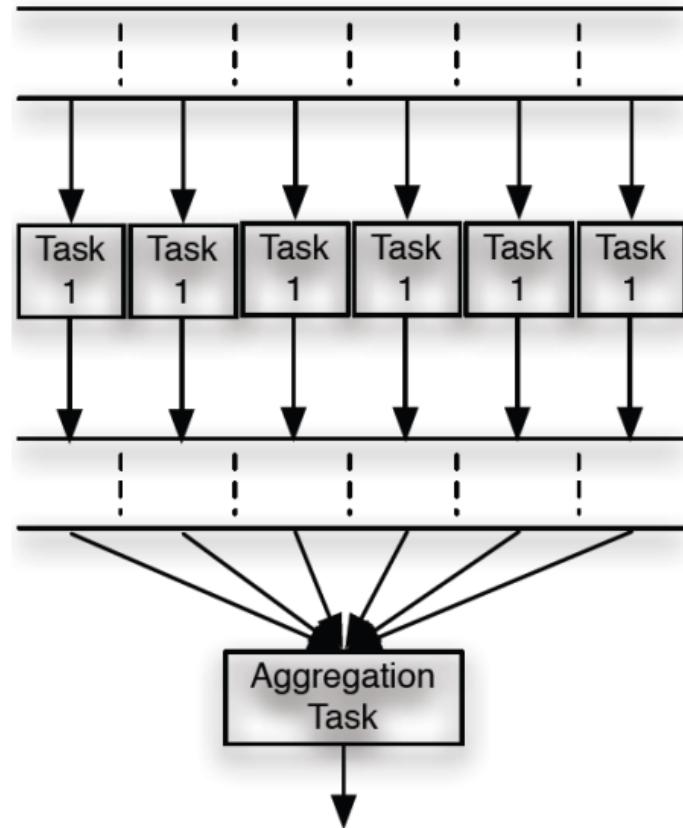
What does all this mean?

We have to "think parallel" when writing code

We should adopt parallel capable libraries to help us do that

Task Parallelism and Data parallelism examples

Data Parallelism

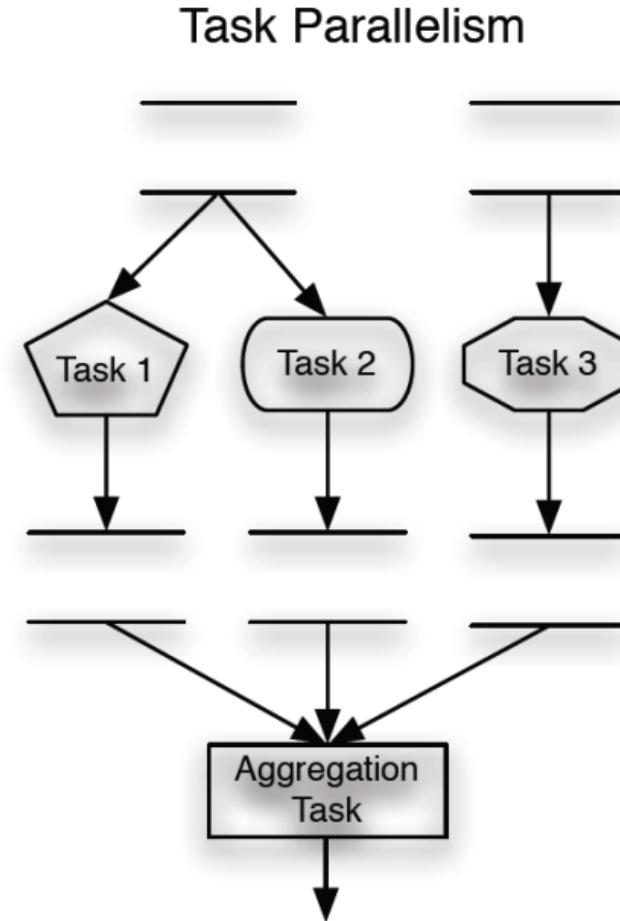


Task Parallelism

Input Data

Parallel Processing

Result Data



Programming considerations

- Tasks need to be able to run in parallel
- Aggregation tasks are potential bottlenecks
- Our approach may be different depending on if our application aims to minimize “time-to-insight” or to minimize “cost per event”

Amdahl's “Law”: (P)arallel code fraction determines how much faster your application can go using N cores/threads

Amdahl's law gives the theoretical speed up of a code when running on N processors, given the fraction of the code that is parallel, P :

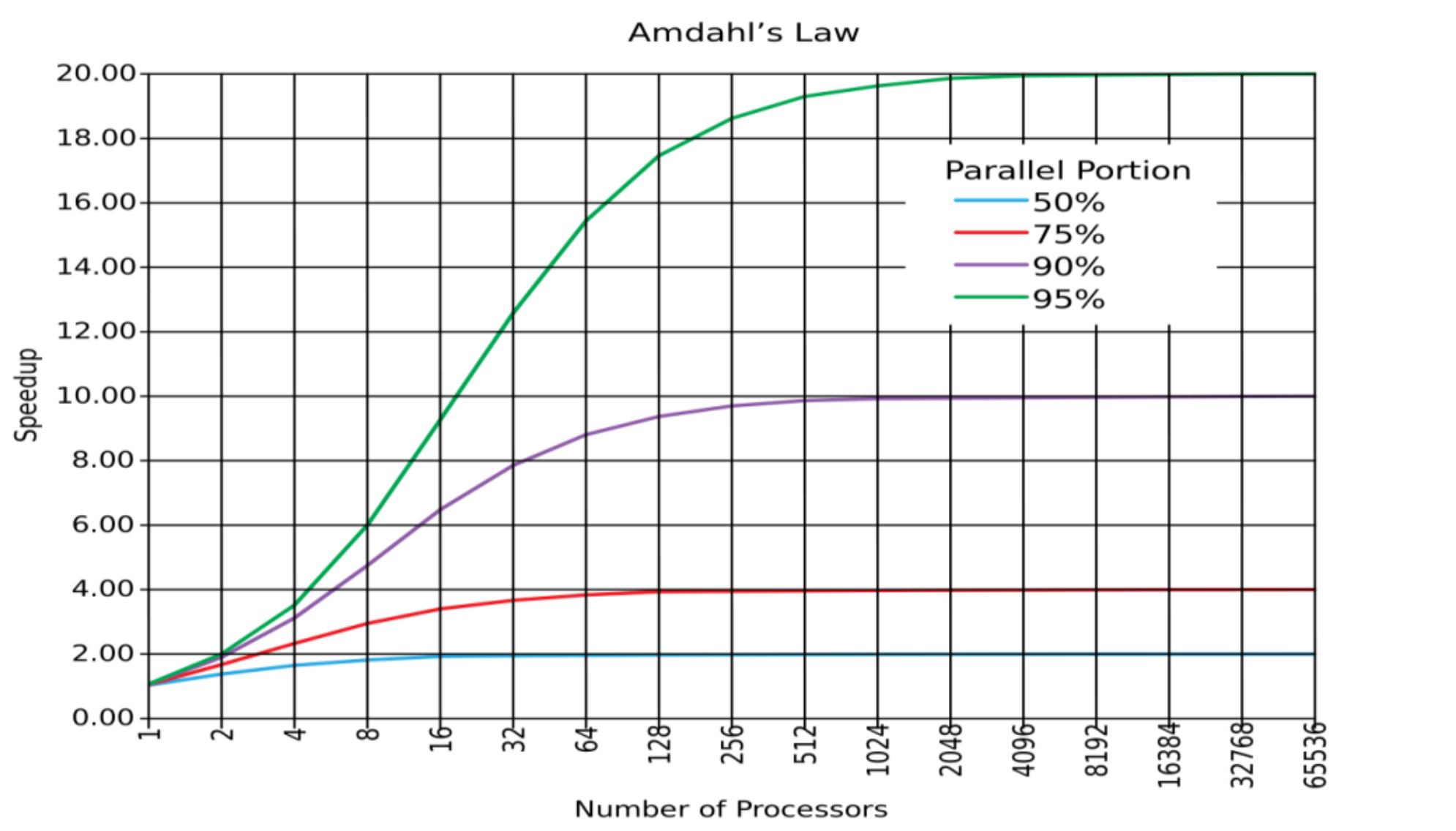
$$S = \frac{1}{(1 - P) + (P/N)}$$

To use a 64 core CPU at high efficiency, how much of your application must be parallel?

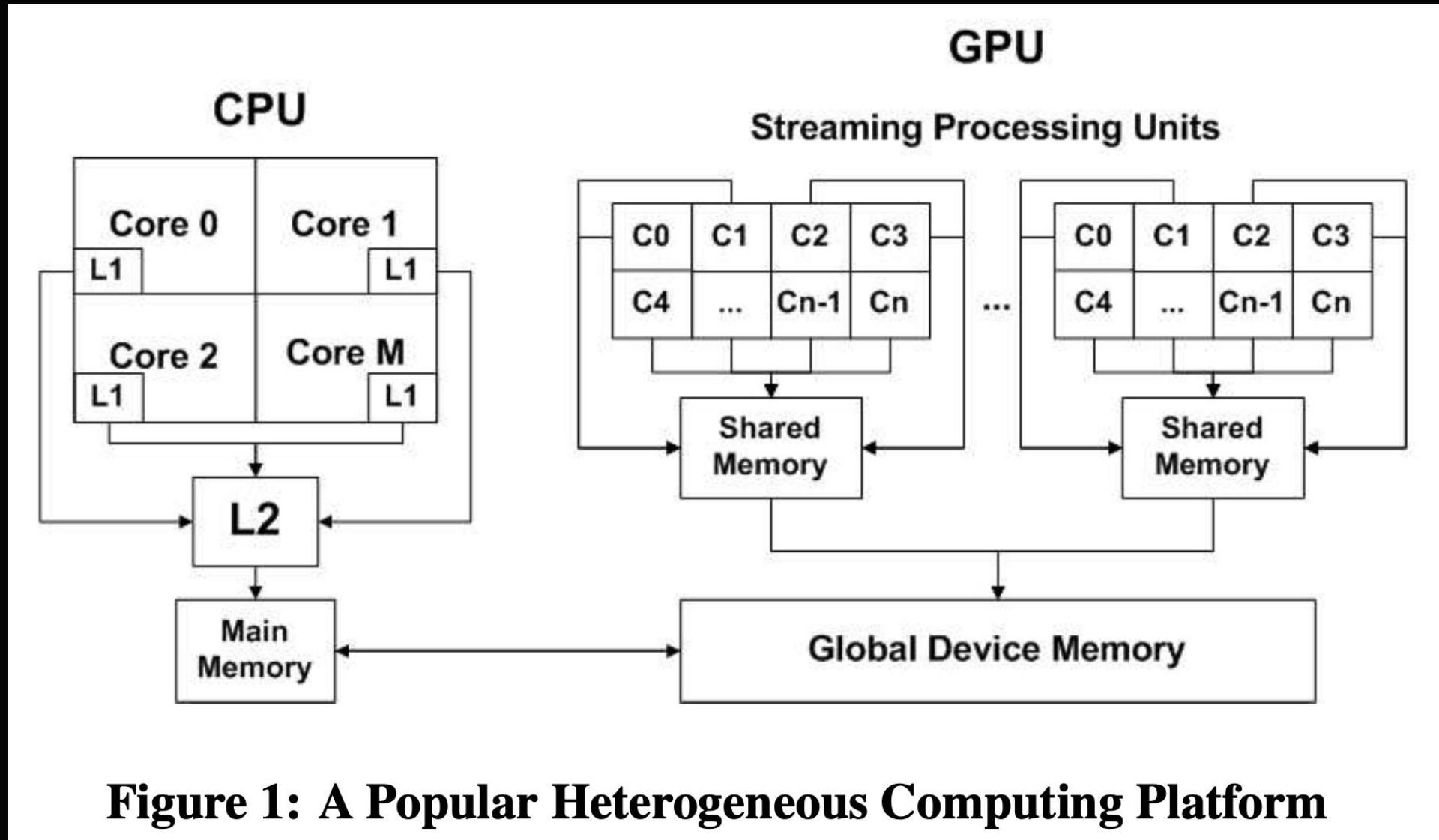
- a) 50%
- b) 90%
- c) 99%
- d) 99.9%
- e) 99.99%

If you want to use 90% of the cores
you need 99.9% parallel code in this case





Heterogeneous computing

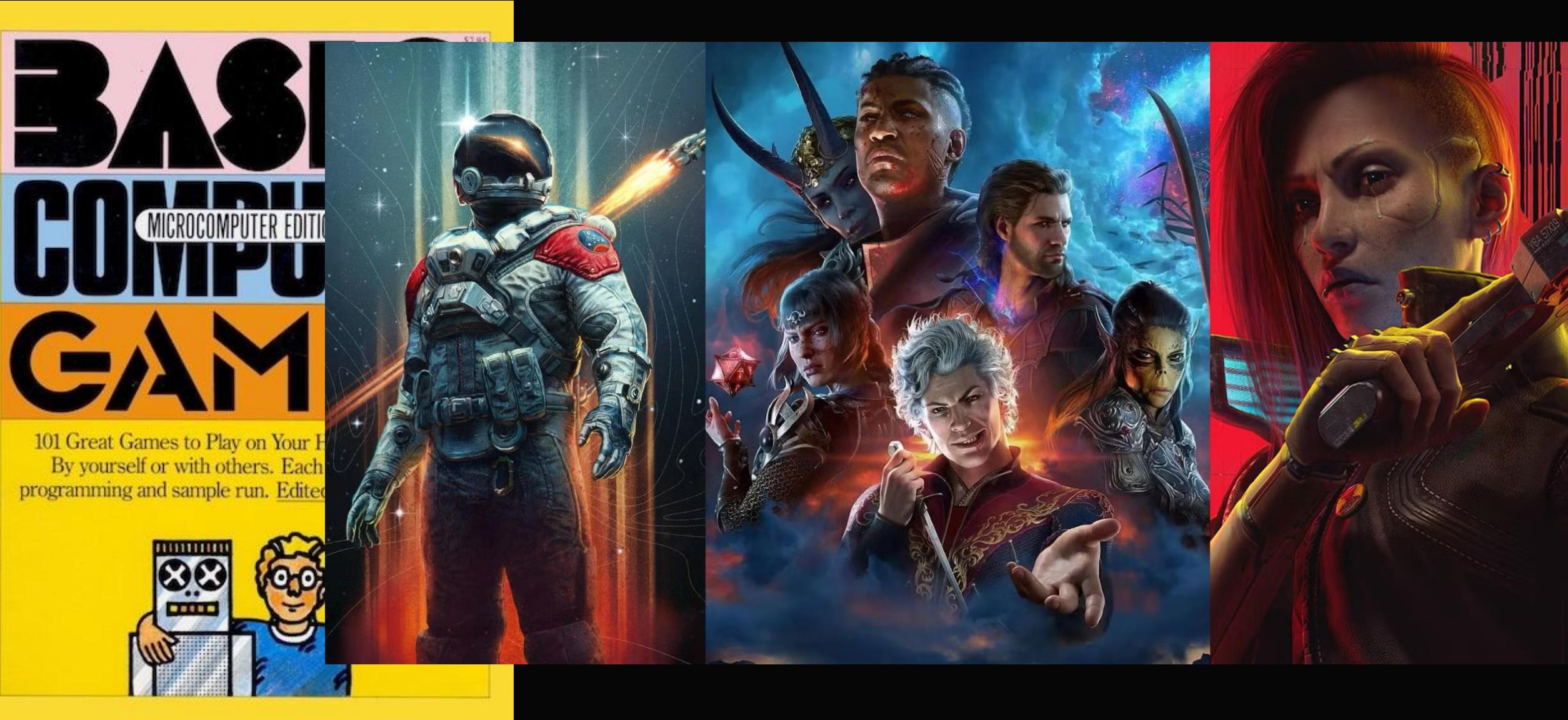


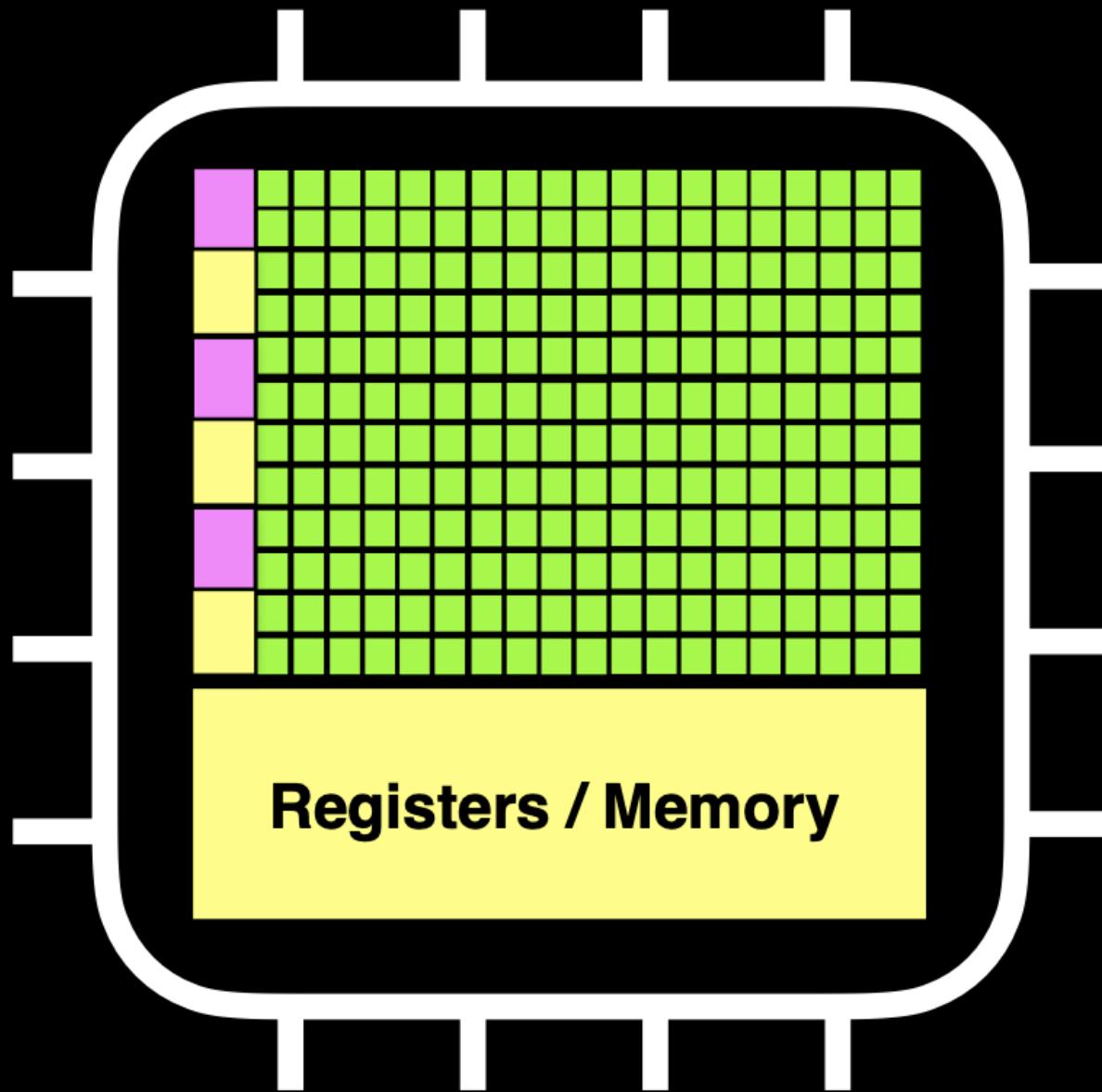
What is a GPU? Graphics Processing Units

- The processor inside the graphics cards was purpose-specific for graphics only.
- GPGPU: Programming techniques to enable the GPU as a parallel co-processor.
- Currently: Intuitive programming and high number of processors
(e.g. over 16000 cuda cores in RTX4090).



Wait aren't GPUs for gaming?





Arithmetic Logic Unit (ALU):

To perform arithmetic and logic operations

Registers / Memory:

Fast memory for Input and output of ALUs

Control Unit:

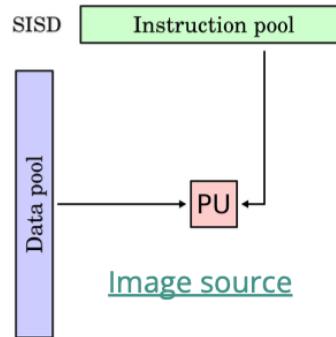
Directs operation of the processor; moving memory, executing ALUs, etc.

GPU devotes more silicons to computing

Flynn's classification of computer architecture

- Based on the number of instruction and data streams that can be processed simultaneously, computing systems are classified into four major categories

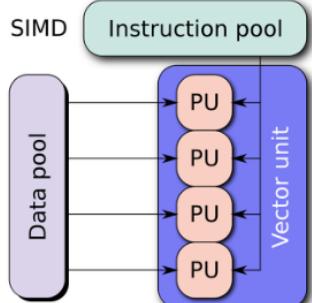
i.e. single cpu-core machine



Data streams

SISD (Single Instruction Single Data)	MISD (Many Instructions Single Data)
SIMD (Single Instruction Many Data)	MIMD (Many Instructions Many Data)

Instruction streams



i.e. vector processor cpus

i.e pipeline architectures - not commonly used

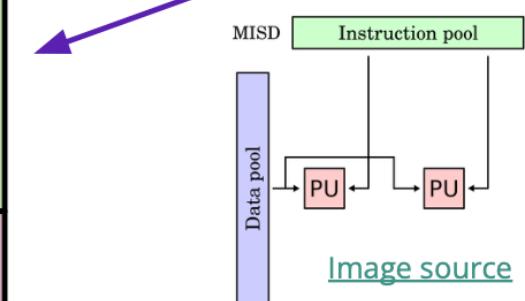


Image source

i.e multi cpu-core machines / grid computing etc.

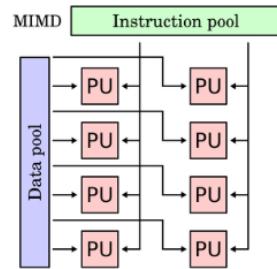
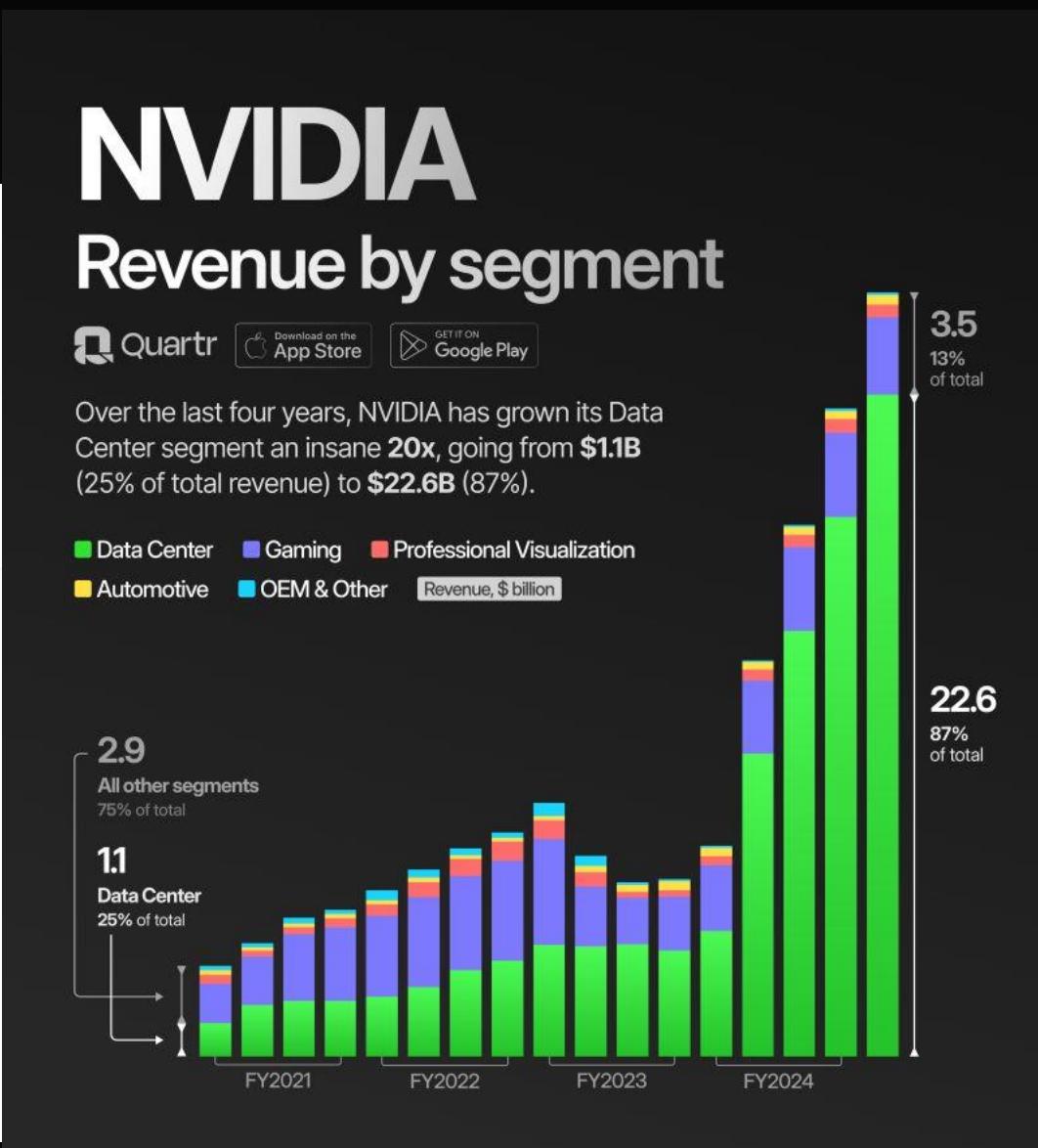
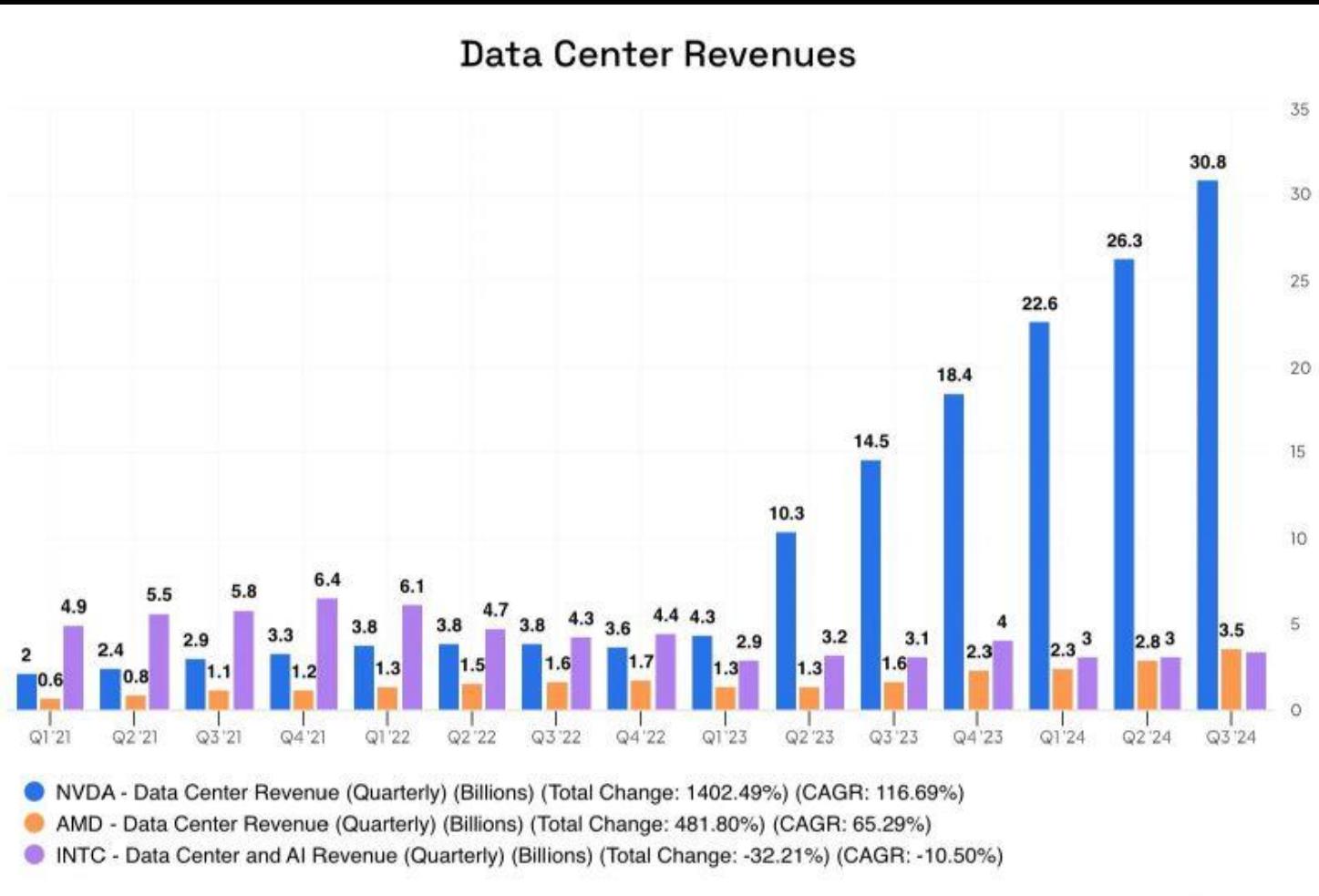


Image source

Other types of hardware accelerators

- FPGA (Field Programmable Gate Array)
 - FPGAs can be reprogrammed to suit the needs of an application or required functionality
 - More about FPGAs later this week
- ASIC: Integrated circuit customized for a specific purpose/application
- TPU: Optimized to perform matrix-multiplication operations

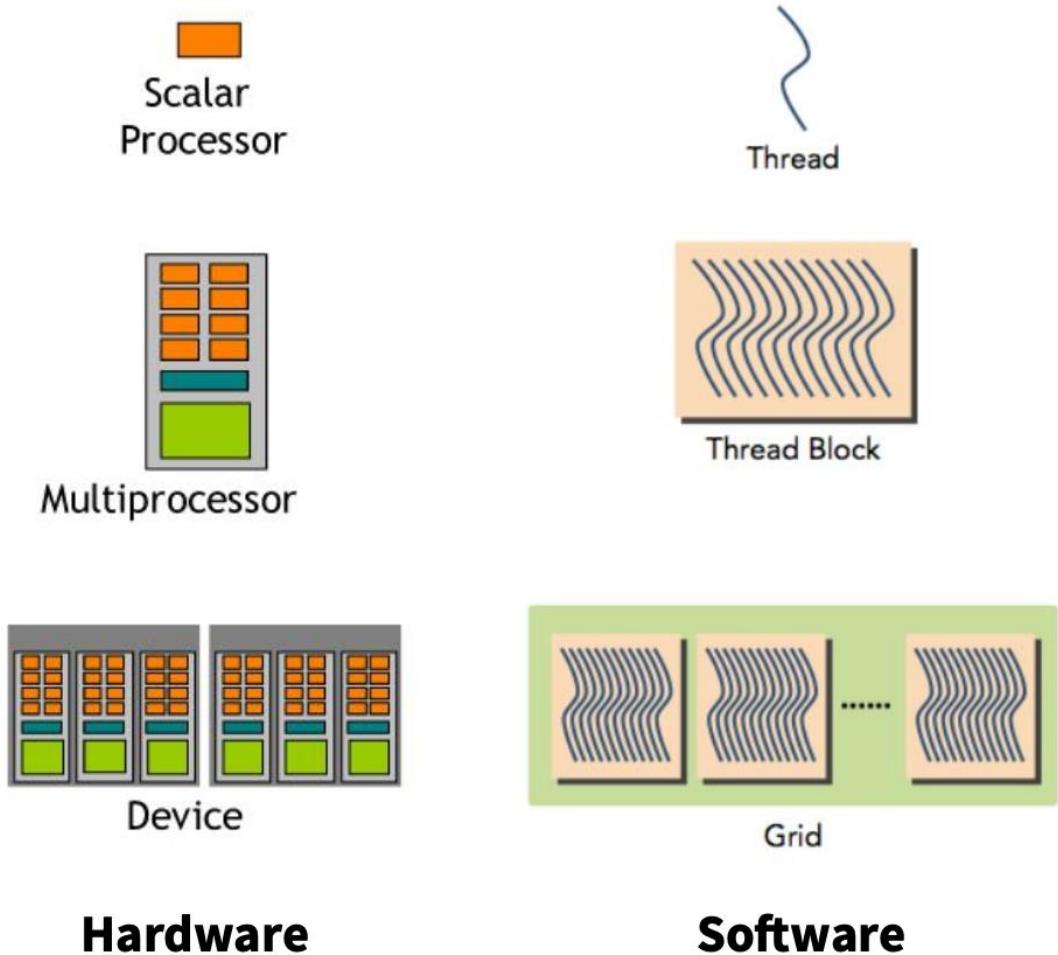
GPUs have risen together with AI



How does a GPU differ from a many-core CPU?

- CPU: Follows the Single instruction, multiple threads (SIMT) execution model. Asynchronous programming model where threads are executed independently
- GPU: High data throughput and massive parallel computing: a GPU consist of hundreds of cores performing the same operation on multiple data items in parallel.
 - GPU architectures are built around a scalable array of Streaming Multiprocessors (SM). Each SM in a GPU is designed to support concurrent execution of hundreds of threads
 - Lower clock speeds, higher latencies, higher throughput, lower power per performance

Hardware to software



- A **scalar processor or CUDA core is equivalent to a software thread**
- Scalar processors are grouped into a SM
- Each execution of a GPU function is done concurrently on a number of threads referred to as a **thread block**
- Each thread block is executed by one SM and cannot be migrated to other SMs in GPU
- The **set of thread blocks** executing the GPU function is called a **grid**.
- In CUDA terminology the GPU is referred to as the **device**

Software for accelerator hardware

Applications

Libraries

Easy to use
Most Performance

Compiler
Directives

Easy to use
Portable code

Programming
Languages

Most Performance
Most Flexibility

GPU Accelerated Libraries

Linear Algebra
FFT, BLAS,
SPARSE, Matrix

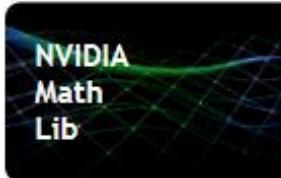


cULA|tools



C U S P

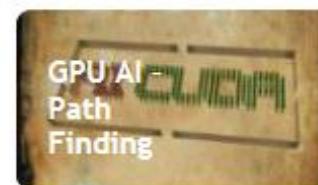
Numerical & Math
RAND, Statistics



ArrayFire



Data Struct. & AI
Sort, Scan, Zero Sum



cuDNN

Visual Processing
Image & Video



NVIDIA
Video
Encode

Sundog™
Software

Nvidia GPU programming: CUDA

CUDA Parallel Computing Platform
www.nvidia.com/getcuda



Programming Approaches

Libraries OpenACC Directives Programming Languages

“Drop-in” Acceleration Easily Accelerate Apps Maximum Flexibility

Development Environment

 Nsight IDE
Linux, Mac and Windows
GPU Debugging and Profiling

CUDA-GDB debugger
NVIDIA Visual Profiler

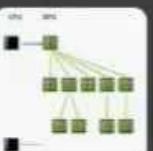
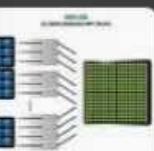
Open Compiler Tool Chain

 LLVM COMPILER INFRASTRUCTURE

Enables compiling new languages to CUDA platform, and CUDA languages to other architectures

Hardware Capabilities

SMX Dynamic Parallelism HyperQ GPUDirect


To effectively use such a high thread count, “kernels” must be identified and sent to the accelerator to compute

- The kernel expresses the portion of the application that is parallelizable
 - On a NVIDIA GPU: it will be executed in parallel by different CUDA threads
- Linear algebra is the classic example
 - Large arrays allocated in memory shared across multiple (many) cores/threads simultaneously
 - Each thread operates on a smaller portion of these arrays
 - Avoid conflicts and synchronization points between threads
 - Some intermediate variables may be computed on multiple threads
 - Extra calculations can keep threads synchronized

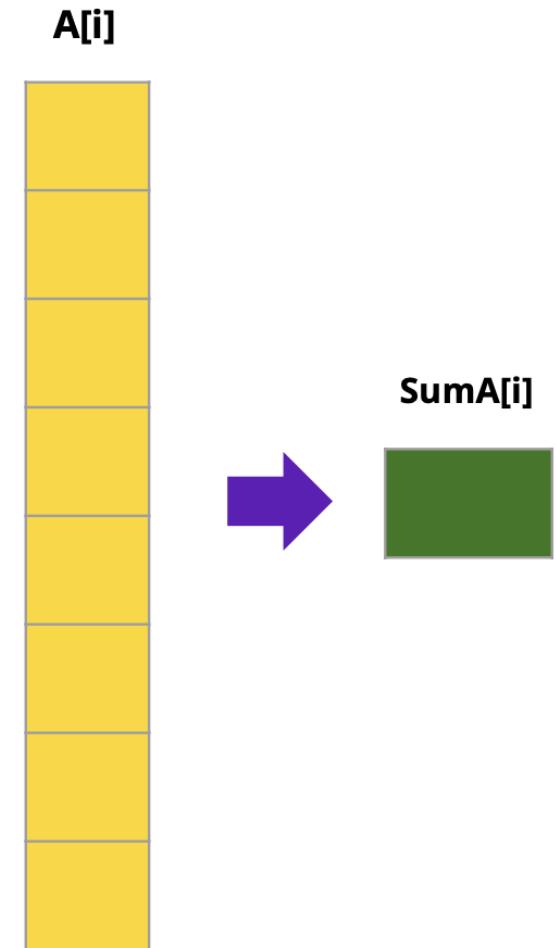
A simple CUDA example?

Adding elements in a vector

Let's start by writing a CUDA kernel that calculated the sum of the elements of a vector :

```
__global__ void add_array(float* A, float* sum) {
    int idx = threadIdx.x + blockIdx.x * blockDim.x;
    if (idx < N) {
        *sum += A[idx];
    }
}
```

- There are 3 instructions that will be executed :
 - **Load** the value of A for each thread
 - **Read** the value of c
 - **Modify** the value of c



What goes wrong here?

Extra care is needed any time a thread writes to shared memory

Adding elements in a vector

Using **atomicAdd** to sum the vector elements :

```
__global__ void add_array(float* A, float* sum) {  
    int idx = threadIdx.x + blockIdx.x * blockDim.x;  
    if (idx < N) {  
        atomicAdd(sum, A[idx]);  
    }  
}
```

Each read-modify-write access cannot
be interrupted

A[i]



SumA[i]



Now the sum will be correct!!

Programming model

- There is no magic compiler flag that makes your code run in parallel. So your approach really depends on what you are doing and where you intend to run.
- You must design from the start for parallelism
- I can't be exhaustive today about how to do this... Let me talk about a few examples

Memory bandwidth and data movement

- To do a calculation on a GPU, all needed data must be copied to the GPU (eg, host to device transfer)
- GPU memory bandwidth is 100s or 1000s of GB/second. However this can be a limitation:
 - If your compute kernel is too small
 - If your application needs to frequently move data back and forth between the device and host.
- Common tricks:
 - Dimensionality reduction
 - Efficiently representations (eg, sparse matrices)
 - Building a series of kernels that can be performed without intermediate data transfer (even if one might be less efficient on a GPU by itself)

Which code results in more a performant parallel application? (independently of if a GPU, vector CPU, etc)

```
row = blockIdx.x*blockDim.x + threadIdx.x;
for (col=0; col<N; col++)
    output[row][col]=2*input[row][col];
```

```
col = blockIdx.x*blockDim.x + threadIdx.x;
for (row=0; row<N; row++)
    output[row][col]=2*input[row][col];
```

For GPUs:

- Maximum bandwidth is when data for multiple threads can be loaded in a single transaction
- This happens if data access patterns are "close enough". Eg, if 16 consecutive threads all access data from within the same memory segment.
- This condition is met when consecutive threads read consecutive memory addresses. If not, memory accesses are serialized, significantly degrading performance.

Tools and Libraries: OpenMP

- The OpenMP API supports multi-platform shared-memory parallel programming in C/C++ and Fortran. The OpenMP API defines a portable, scalable model with a simple and flexible interface for developing parallel applications on platforms from the desktop to the supercomputer.
(<https://www.openmp.org>)

NTU Talk
January 14
2009

12

The OpenMP Execution Model

Fork and Join Model

The diagram illustrates the Fork and Join Model. It starts with a single blue square node at the top, representing the Master Thread. This node has four arrows pointing down to a horizontal bar labeled "Parallel region". From each end of this bar, two more arrows point down to another horizontal bar labeled "Worker Threads". These four "Worker Threads" are grouped together under the label "Synchronization". This entire process is repeated once more, creating a second "Parallel region" and another set of "Worker Threads" and "Synchronization" points. The Sun Microsystems logo is in the top right corner.

NTU Talk
January 14
2009

10

The OpenMP Memory Model

The diagram illustrates the OpenMP Memory Model. In the center is a large purple cloud labeled "Shared Memory". Six pink circles, each containing a capital letter "T" and labeled "private", are connected to the cloud by arrows. The Sun Microsystems logo is in the top right corner.

- ✓ All threads have access to the same, globally shared, memory
- ✓ Data can be shared or private
- ✓ Shared data is accessible by all threads
- ✓ Private data can only be accessed by the thread that owns it
- ✓ Data transfer is transparent to the programmer
- ✓ Synchronization takes place, but it is mostly implicit

RvdP/V1

An Overview of OpenMP

Tools and Libraries: std::thread

```
// thread example
#include <iostream> // std::cout
#include <thread> // std::thread

void foo() { // do stuff... }
void bar(int x) { // do stuff... }

int main() {
    std::thread first (foo); // spawn new thread that calls foo()
    std::thread second (bar,0); // spawn new thread that calls bar(0)
    std::cout << "main, foo and bar now execute concurrently...\n";

    // synchronize threads:
    first.join(); // pauses until first finishes
    second.join(); // pauses until second finishes

    std::cout << "foo and bar completed.\n";
    return 0;
}
```

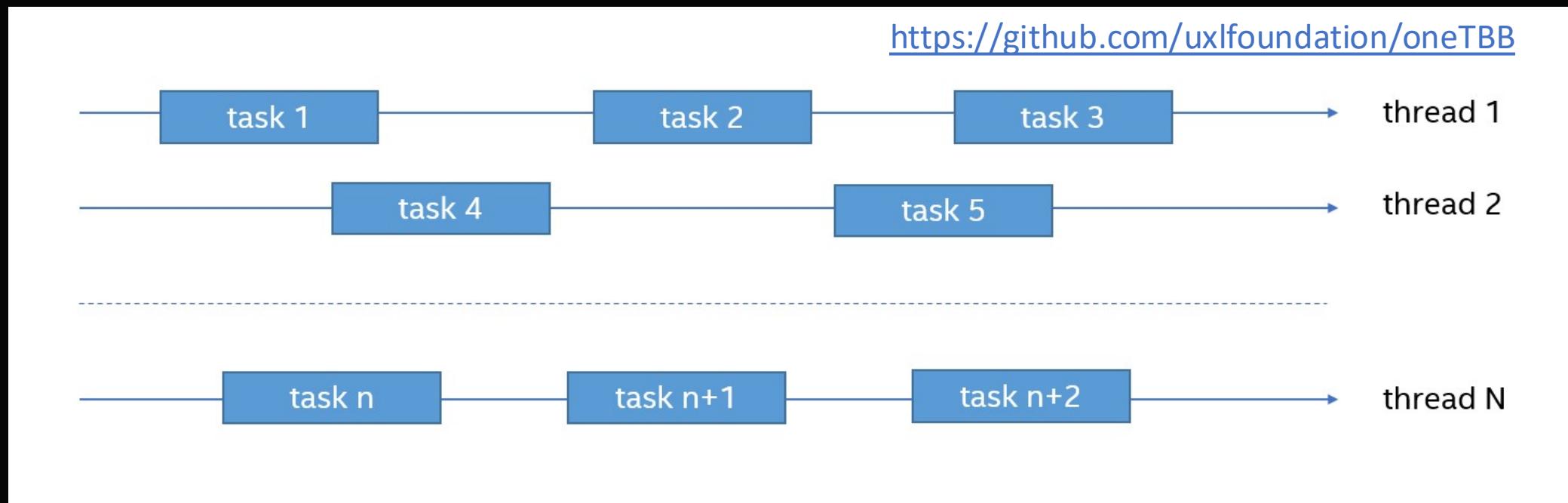
Similarly in python....

```
import threading  
  
def print_cube(num):  
    print("Cube: {}" .format(num * num * num))  
  
def print_square(num):  
    print("Square: {}" .format(num * num))  
  
  
if __name__ == "__main__":  
    t1 = threading.Thread(target=print_square, args=(10,))  
    t2 = threading.Thread(target=print_cube, args=(10,))  
    t1.start()  
    t2.start()  
    t1.join()  
    t2.join()  
    print("Done!")
```

Tools and Libraries: TBB

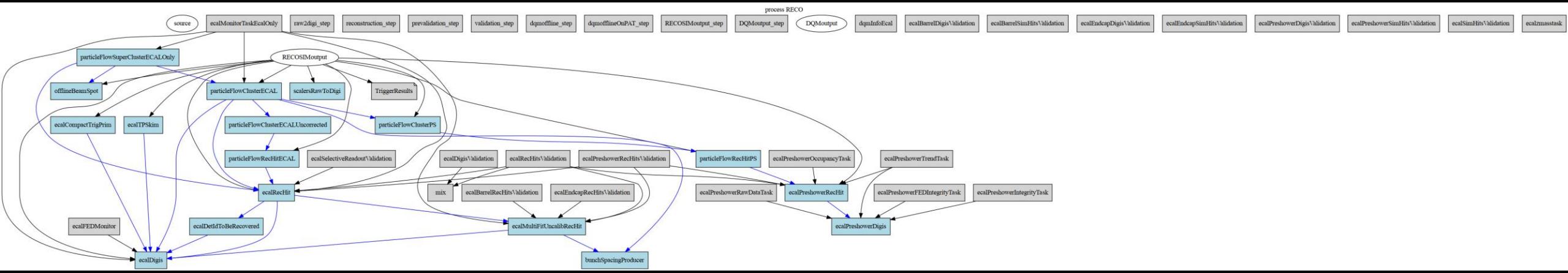
TBB is well suited for scalable applications that:

- **Specify logical parallel structure instead of threads.**
- Emphasize data-parallel programming.
- Take advantage of concurrent collections and parallel algorithms.



The CMS multi-threaded framework uses TBB to schedule 1000s of tasks per event and to allow multiple events to be processed concurrently

- Dependency (task) graph from a small set of CMSSW algorithms



Does heterogeneity mean I need to write code for each different architecture?

- Maybe...but doesn't each GPU vendor have their own API..
- Fortunately, there are a number of approaches that aim for performance portability. Eg, to enable portability across accelerators through the abstraction of the underlying levels of parallelism without giving much (too much) performance on any given backend

Example: Alpaka

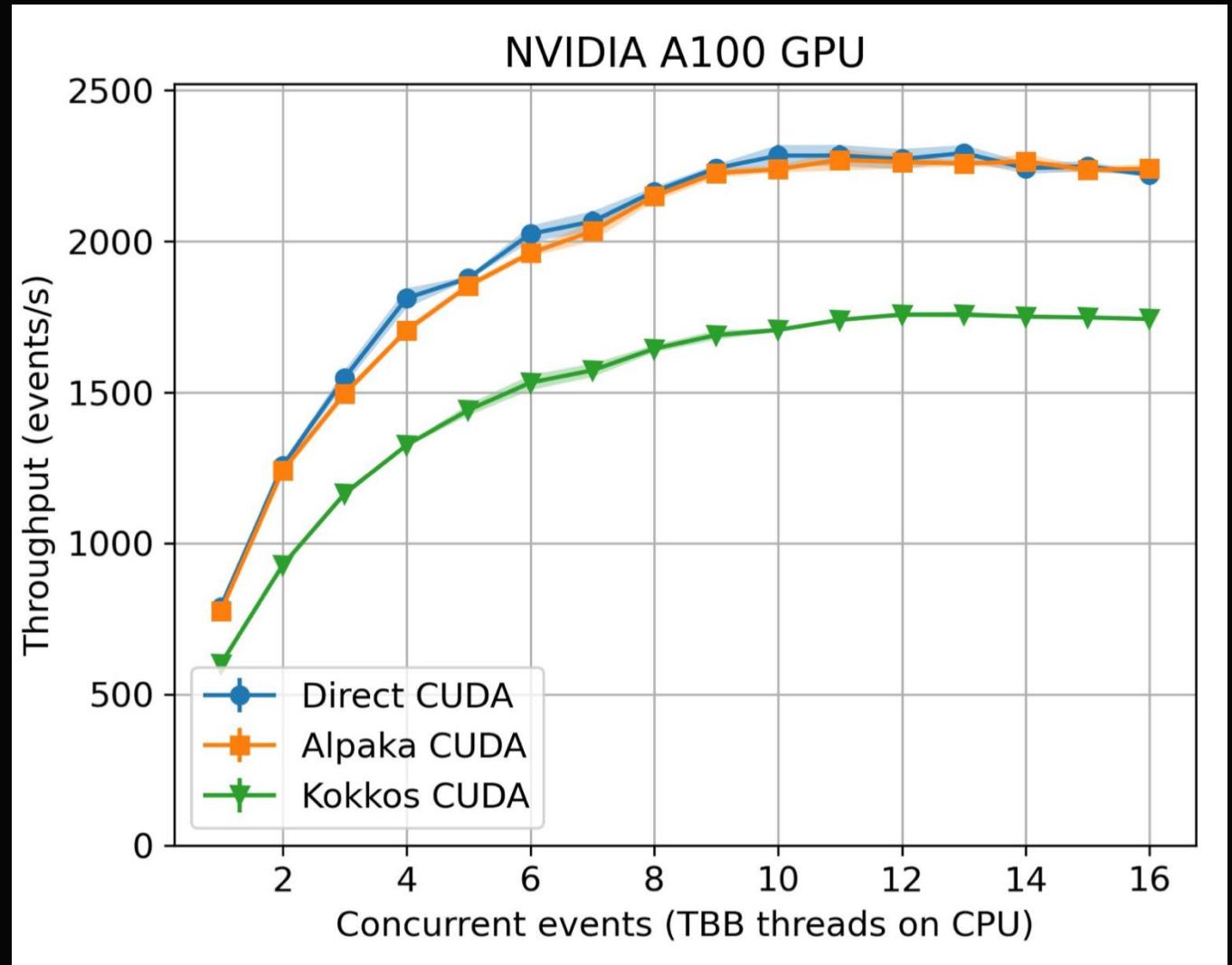
- The **alpaka** library is a header-only C++20 abstraction library for accelerator development. <https://alpaka-group.github.io/alpaka/>

Accelerator Back-ends

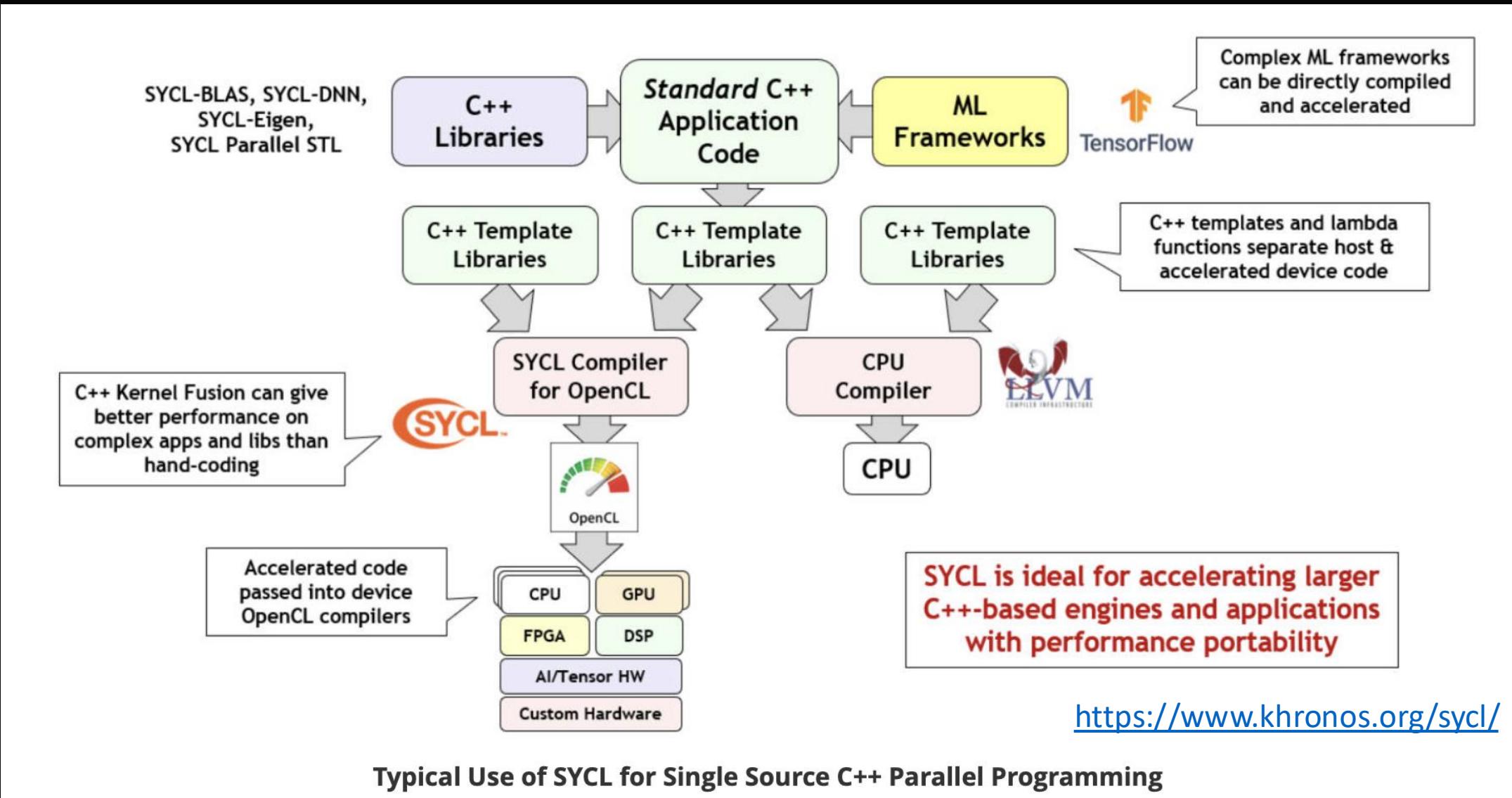
Accelerator Back-end	Lib/API	Devices	Execution strategy grid-blocks	Execution strategy block-threads
Serial	n/a	Host CPU (single core)	sequential	sequential (only 1 thread per block)
OpenMP 2.0+ blocks	OpenMP 2.0+	Host CPU (multi core)	parallel (preemptive multitasking)	sequential (only 1 thread per block)
OpenMP 2.0+ threads	OpenMP 2.0+	Host CPU (multi core)	sequential	parallel (preemptive multitasking)
std::thread	std::thread	Host CPU (multi core)	sequential	parallel (preemptive multitasking)
TBB	TBB 2.2+	Host CPU (multi core)	parallel (preemptive multitasking)	sequential (only 1 thread per block)
CUDA	CUDA 12.0+	NVIDIA GPUs	parallel (undefined)	parallel (lock-step within warps)
HIP(clang)	HIP 6.0+	AMD GPUs	parallel (undefined)	parallel (lock-step within warps)
SYCL(oneAPI)	oneAPI 2024.2+	CPUs, Intel GPUs and FPGAs	parallel (undefined)	parallel (lock-step within warps)

Alpaka performance examples

- CMS algorithm evaluation on a CPU comparing a CUDA implementation with an Alpaka implementation running on NVIDIA A100
- Essentially no loss of throughput seen (while enabling other backends with the same code)



Example: SYCL



SYCL code example..

```
#include<CL/syd.hpp>
#include<iostream>

int main() {
    using namespace d::syd;
    int data[1024]; // Allocates data to be worked on

    // Include all the SYCL work in a {} block to ensure all
    // SYCL tasks are completed before exiting the block.
    {
        // Create a queue to enqueue work to
        queue myQueue;

        // Wrap the data array variable in a buffer.
        buffer<size_t, 1> resultBuf { data, range<1>{ 1024 } };

        // Create a command group to
        // issue commands to the queue.
        myQueue.submit([&](handler & cgh)
        {
            // Request access to the buffer
            auto writeResult = resultBuf.get_access<access::mode::write>(cgh);

            // Enqueue a parallel_for task.
            cgh.parallel_for<class simple_test>(range<1>{ 1024 }, [=](id<1> idx)
            {
                writeResult[idx] = idx[0];
            }); // End of the kernel function
        }); // End of the queue commands
    } // End of scope, so wait for the queued work to complete

    // Print result
    for (int i = 0; i < 1024; i++) {
        std::cout << "data[" << i << "] = " << data[i] << std::endl;
    }
    return 0;
}
```

SYCL programs must include the `<CL/syd.hpp>` header file to provide all of the SYCL features.

Namespace

All SYCL names are defined in the `d::syd` namespace.

Queue

See queue class functions [4.6.5] on page 2 of this reference guide.

Buffer

See buffer class functions [4.7.2] on page 2 of this reference guide.

Accessor

See accessor class functions [4.7.6] on page 3 of this reference guide.

Handler

See handler class functions [4.8.3] on page 5 of this reference guide.

Scopes

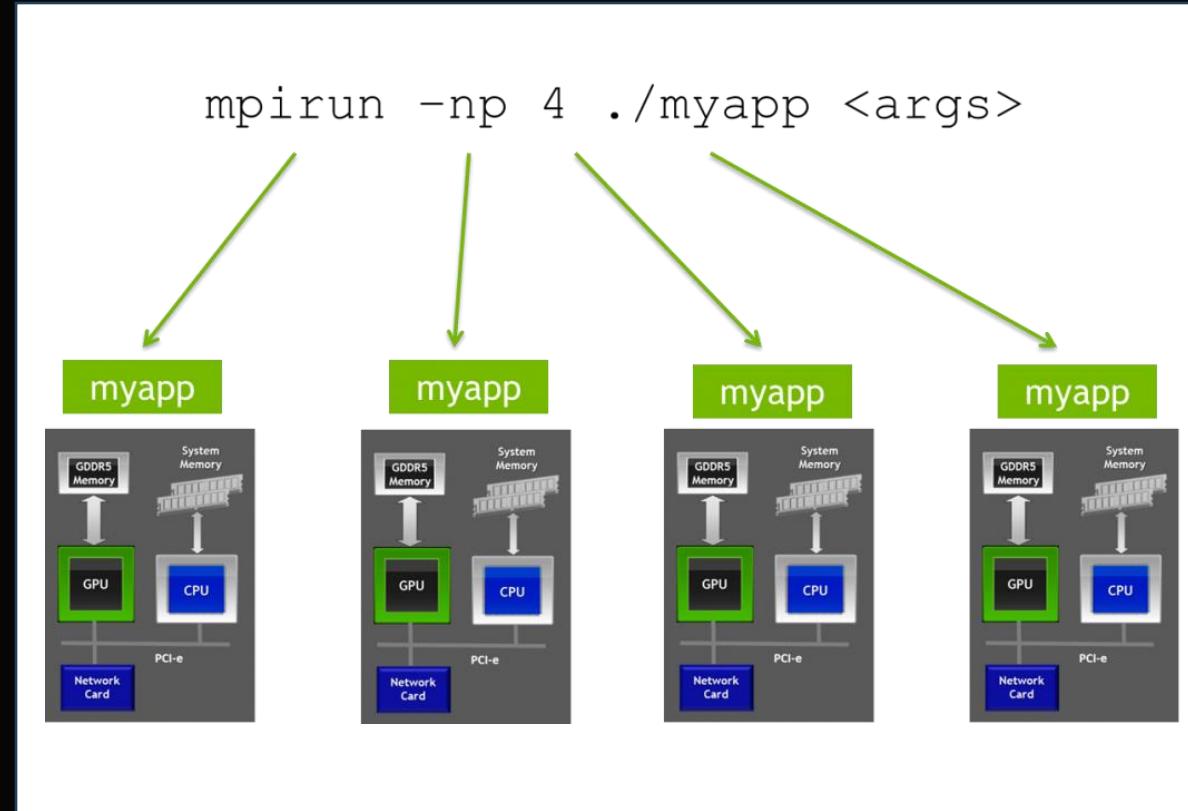
The **kernel scope** specifies a single kernel function that will be, or has been, compiled by a device compiler and executed on a device. See Invoking Kernels [4.8.5] in the spec.

The **command group scope** specifies a unit of work which is comprised of a kernel function and accessors.

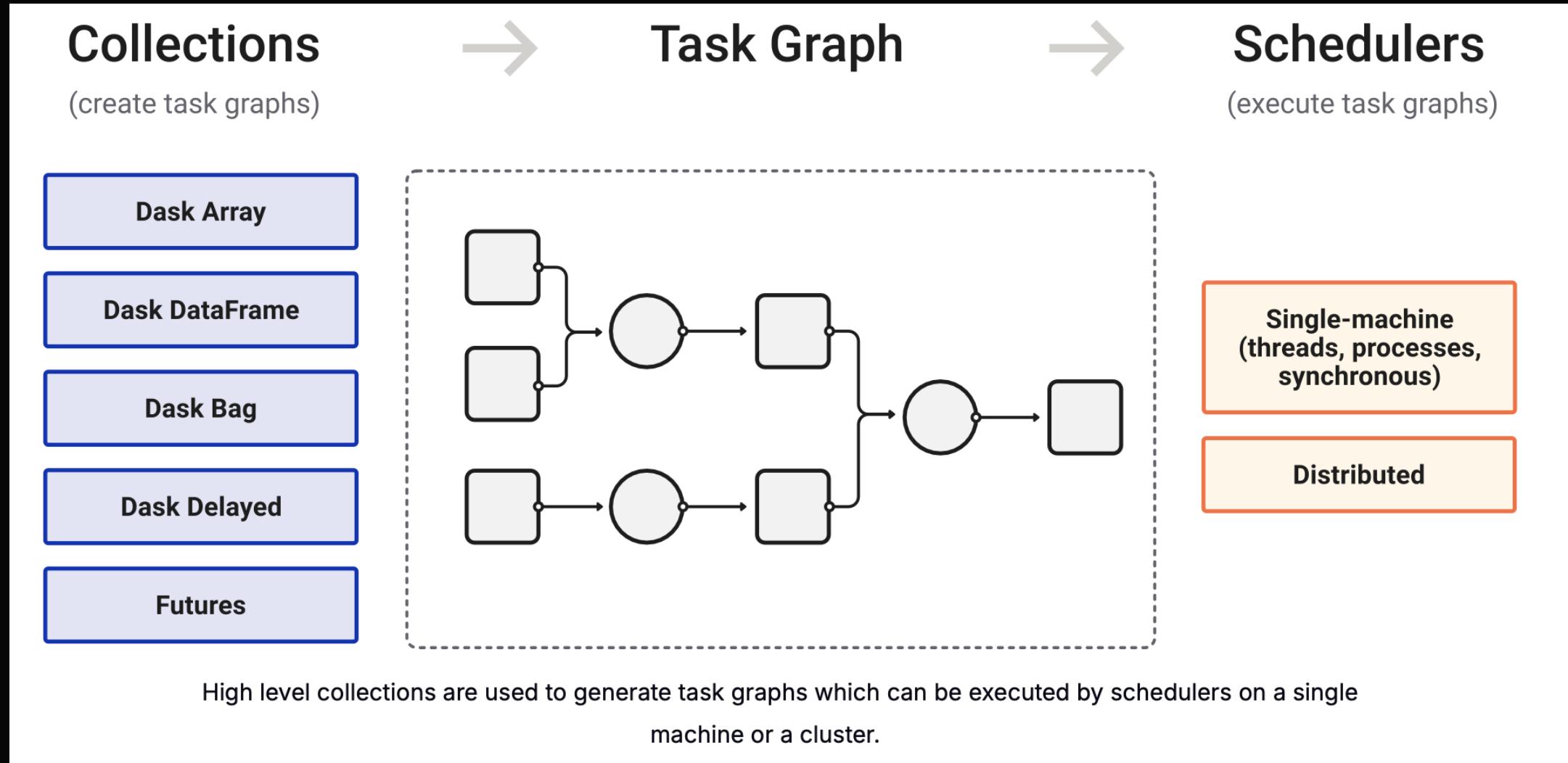
The **application scope** specifies all other code outside of a command group scope.

Libraries for scaling codes to large problems: MPI

- Message Passing Interface (MPI) is interface defining a model of parallel computing where each parallel process has its own local memory, and data must be explicitly shared by passing messages between processes.
- Using MPI allows programs to scale beyond the processors and shared memory of a single compute server, to the *distributed memory* and processors of multiple compute servers combined together.



Another library for horizontal scaling but for Python: Dask



Scientific Python on GPUs

Python and GPU

Our introduction to array-oriented programming started by addressing Python's slowness (or the slowness of interactive environments in general).

But GPUs are also best addressed in an array-at-a-time fashion: the memory layout explicitly favors arrays over objects addressed by pointers or references (even in C++).

It should not be a surprise, then, that Python's GPU libraries are array libraries.

Python and ML

- Most users of GPUs in Python are using them for Machine Learning (ML).
- Since ML libraries are focused on ML, rather than programming in general, they provide an interface that is identical for CPUs and GPUs, usually by a device flag in the array object.
 - Switching from `device="cpu"` to `device="cuda"` simply speeds up the calculation by an order of magnitude with no other changes.
- In this lesson, we'll focus on GPU programming in which the differences from CPU programming are visible.

Most libraries specifically address CUDA, Nvidia's brand of GPU programming

- This reflects the very tight monopoly Nvidia has on GPUs for general-purpose programming.
 - CuPy has *experimental* support for AMD but Numba deprecated support for AMD due to lack of community interest.
 - There are Python projects that support OpenCL, Vulkan, Kokkos, and SYCL.
- Unfortunately, this excludes many laptops, including all MacOS: CUDA is something you'll more likely use on a data center than a laptop.

Installing a working version of CUDA can be hard, as it depends strongly on the operating system and GPU hardware

- The GPU driver must be installed in the operating system (Windows or Linux) and not with any Python package installers (pip, conda, uv, pixi). You might get it directly from Nvidia or through the operating system's package installer (such as apt for Ubuntu).
- The CUDA compiler and its libraries are distinct from the device driver, and they can be installed with conda. The story is more complicated for CUDA 11, but for CUDA 12, just include `cuda-version=12` in the list of packages to install.
- Nvidia publishes a [table of version compatibility](#) between GPU hardware, drivers, and CUDA packages.

Lets have a look at our setup

Two useful shell commands:

- nvidia-smi is Nvidia's tool to find out what hardware is installed and what it is doing
- Numba comes with a numba -s command line tool to print the versions of everything: incompatibilities between driver and CUDA package can be found here.

Programming model

An algorithm intended for a CPU, like this:

```
void some_function(float* array, int index) {
    array[index] = ...;
}

for (int i = 0; i < 1000000; i++) {
    some_function(array, i);
}
```

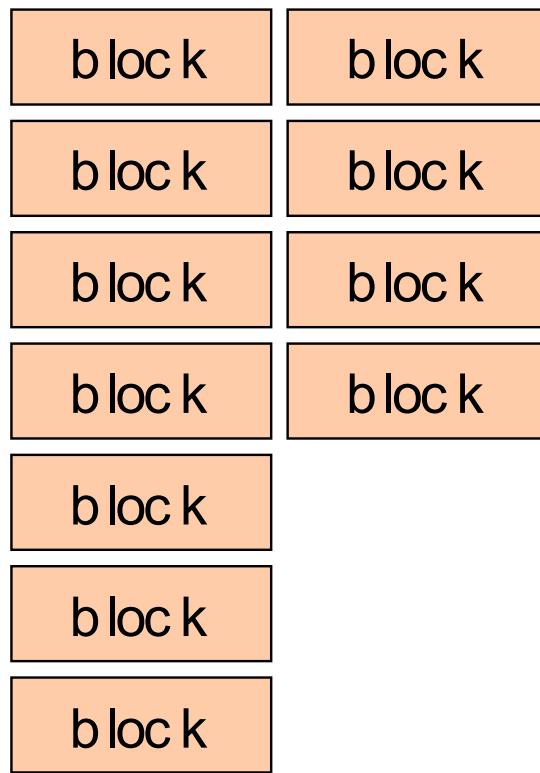
would get translated into something like this for a GPU:

```
__global__ some_function(float* array) {
    int index = blockIdx.x * blockDim.x + threadIdx.x;
    array[index] = ...;
}

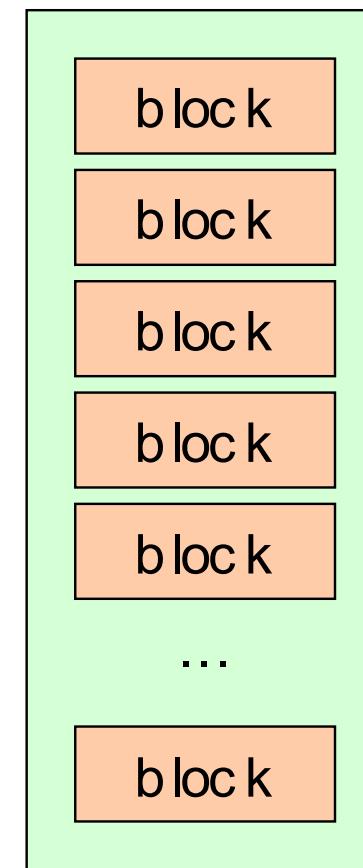
some_function<<<blocks_per_grid, threads_per_block>>>(array);
```

- The explicit iteration over elements of an array on the CPU is replaced by a GPU "kernel" function that applies to just one array element
- Each kernel-execution "thread" might run at any time.
- Execution is organized into "blocks": threads within a block run at the same time on shared resources, and the GPU schedules block execution in a way that keeps its streaming multiprocessors busy.
- The array must already be in the GPU's global memory.

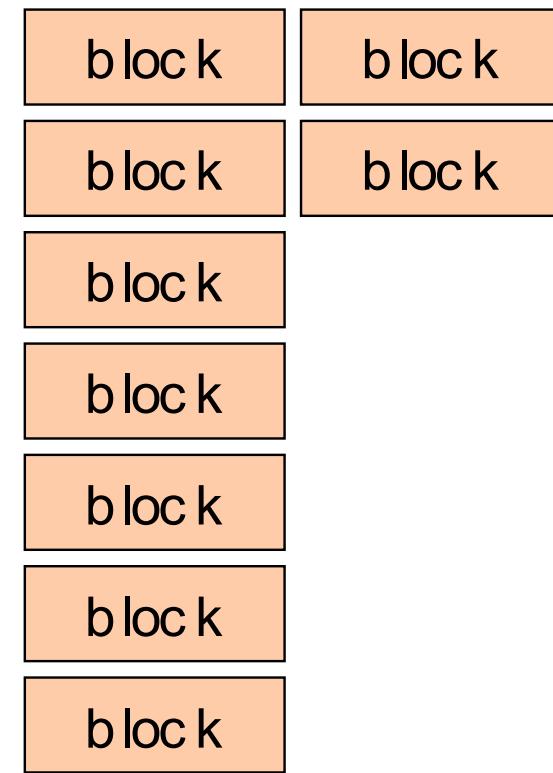
waiting to start



in progress



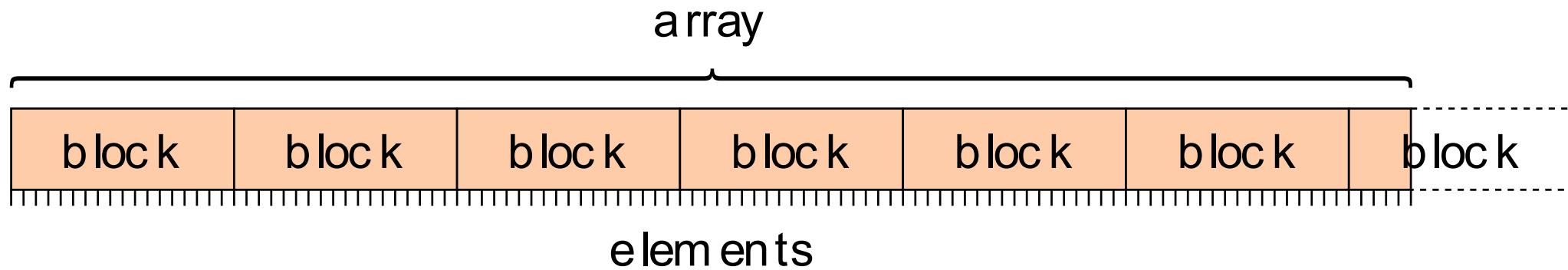
completed



- The **blocks_per_grid** and **threads_per_block** are units of work—they control how computations are scheduled on the GPU—not units of array allocation. However, for 1-dimensional data, it's often best to define them in terms of the array size:

```
threads_per_block = 1024;    // maximum possible on most GPUs  
blocks_per_grid = int(ceil(1.0 * size_of_array / threads_per_block));
```

- This minimizes the number of blocks needed to make sure that each array element is updated by exactly one thread.



If the `size_of_array` doesn't fit neatly into an integer number of blocks, a few threads will be wasted. To make sure that they don't update uninitialized data, you usually need to check for an excess in the kernel function:

```
__global__ some_function(float* array) {
    int index = blockIdx.x * blockDim.x + threadIdx.x;
    if (index < size_of_array) {
        array[index] = ....;
    }
}
```

- Now for more examples to help get you started with Python libraries on GPUs

Wrap up

Parallel programming libraries exist for any language and some abstract away the underlying heterogeneous hardware

BUT: you can not avoid formulating your algorithms to be efficient running in parallel – be it as independent tasks or data parallel

- Use libraries where you can. They will result in more efficient than code you (or I) can write

Types of parallelism

