

Parallel Programming with Python



Tim's backyard around sunset

Tim Mattson, University of Bristol and Merly.ai

We all love python ... but what about performance

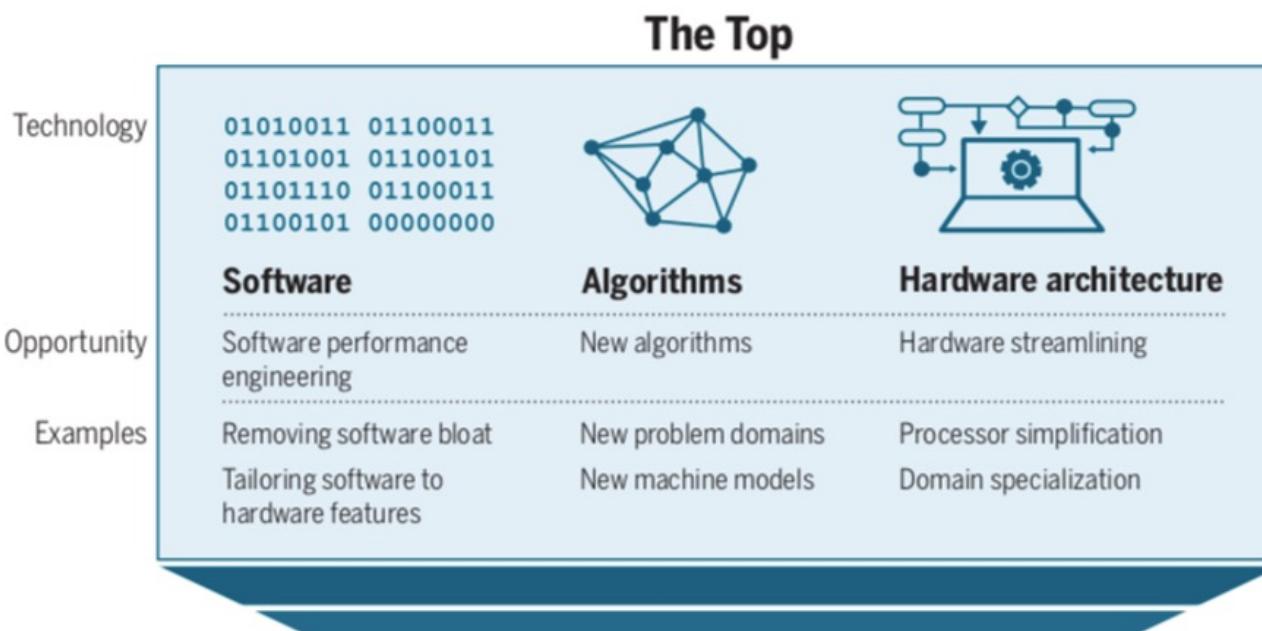
Software vs. Hardware and the nature of Performance

Up until ~2005,
performance came
from semiconductor
technology

* It's because of the end of
Dennard Scaling ...
Moore's law has nothing to
do with it

There's plenty of room at the Top: What will drive computer performance after Moore's law?*

Charles E. Leiserson¹, Neil C. Thompson^{1,2*}, Joel S. Emer^{1,3}, Bradley C. Kuszmaul^{1†},
Butler W. Lampson^{1,4}, Daniel Sanchez¹, Tao B. Schardl¹
Leiserson *et al.*, *Science* **368**, eaam9744 (2020) 5 June 2020



The Bottom
for example, semiconductor technology

Since ~2005
performance comes
from
“the top”

Better software Tech.
Better algorithms
Better HW architecture[#]

#HW architecture matters,
but dramatically LESS than
software and algorithms

The view of Python from an HPC perspective

(from the "Room at the top" paper).

```
for i in range(4096):
    for j in range(4096):
        for k in range (4096):
            C[i][j] += A[i][k]*B[k][j]
```

A proxy for computing over nested loops ...
yes, they know you should use optimized library code for DGEMM

Table 1. Speedups from performance engineering a program that multiplies two 4096-by-4096 matrices. Each version represents a successive refinement of the original Python code. “Running time” is the running time of the version. “GFLOPS” is the billions of 64-bit floating-point operations per second that the version executes. “Absolute speedup” is time relative to Python, and “relative speedup,” which we show with an additional digit of precision, is time relative to the preceding line. “Fraction of peak” is GFLOPS relative to the computer’s peak 835 GFLOPS. See Methods for more details.

| Version | Implementation | Running time (s) | GFLOPS | Absolute speedup | Relative speedup | Fraction of peak (%) |
|---------|-----------------------------|------------------|---------|------------------|------------------|----------------------|
| 1 | Python | 25,552.48 | 0.005 | 1 | — | 0.00 |
| 2 | Java | 2,372.68 | 0.058 | 11 | 10.8 | 0.01 |
| 3 | C | 542.67 | 0.253 | 47 | 4.4 | 0.03 |
| 4 | Parallel loops | 69.80 | 1.969 | 366 | 7.8 | 0.24 |
| 5 | Parallel divide and conquer | 3.80 | 36.180 | 6,727 | 18.4 | 4.33 |
| 6 | plus vectorization | 1.10 | 124.914 | 23,224 | 3.5 | 14.96 |
| 7 | plus AVX intrinsics | 0.41 | 337.812 | 62,806 | 2.7 | 40.45 |

The view of Python from an HPC perspective

(from the "Room at the top" paper).

```
for I in range(4096):  
    for j in range(4096):  
        for k in range (4096):  
            C[i][j] += A[i][k]*B[k][j]
```

A proxy for computing over nested loops ... yes, they know you should use optimized library code for DGEMM

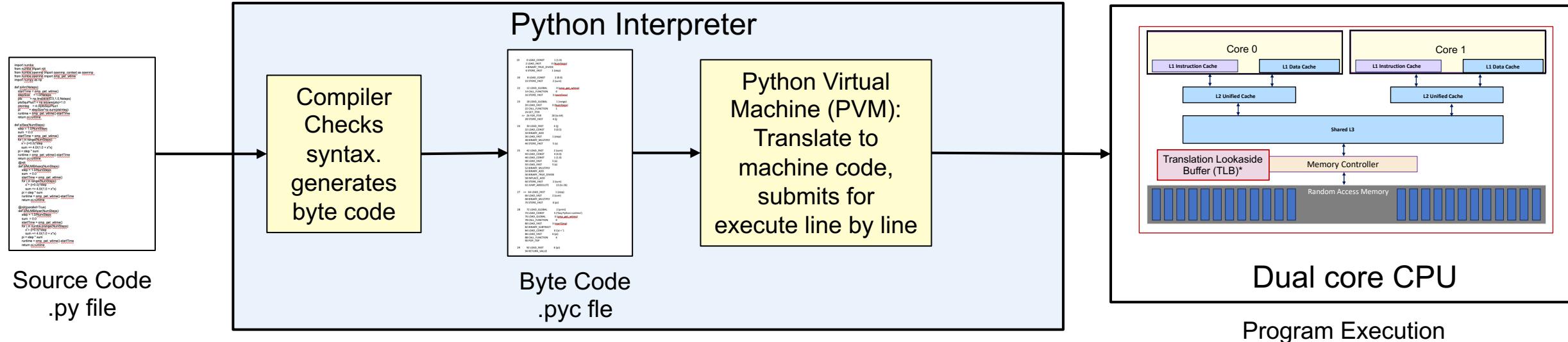
This demonstrates a common attitude in the HPC community

Python is great for productivity, algorithm development, and combining functions from high-level modules in new ways to solve problems. If getting a high fraction of peak performance is a goal ... recode in C.

| Version | Implementation | Running time (s) | GFLOPS | Absolute speedup | Relative speedup | Fraction of peak (%) |
|---------|-----------------------------|------------------|---------|------------------|------------------|----------------------|
| 1 | Python | 25,552.48 | 0.005 | 1 | - | 0.00 |
| 2 | Java | 2,372.68 | 0.058 | 11 | 10.8 | 0.01 |
| 3 | C | 542.67 | 0.253 | 47 | 4.4 | 0.03 |
| 4 | Parallel loops | 69.80 | 1.969 | 366 | 7.8 | 0.24 |
| 5 | Parallel divide and conquer | 3.80 | 36.180 | 6,727 | 18.4 | 4.33 |
| 6 | plus vectorization | 1.10 | 124.914 | 23,224 | 3.5 | 14.96 |
| 7 | plus AVX intrinsics | 0.41 | 337.812 | 62,806 | 2.7 | 40.45 |

Why is Python so slow?

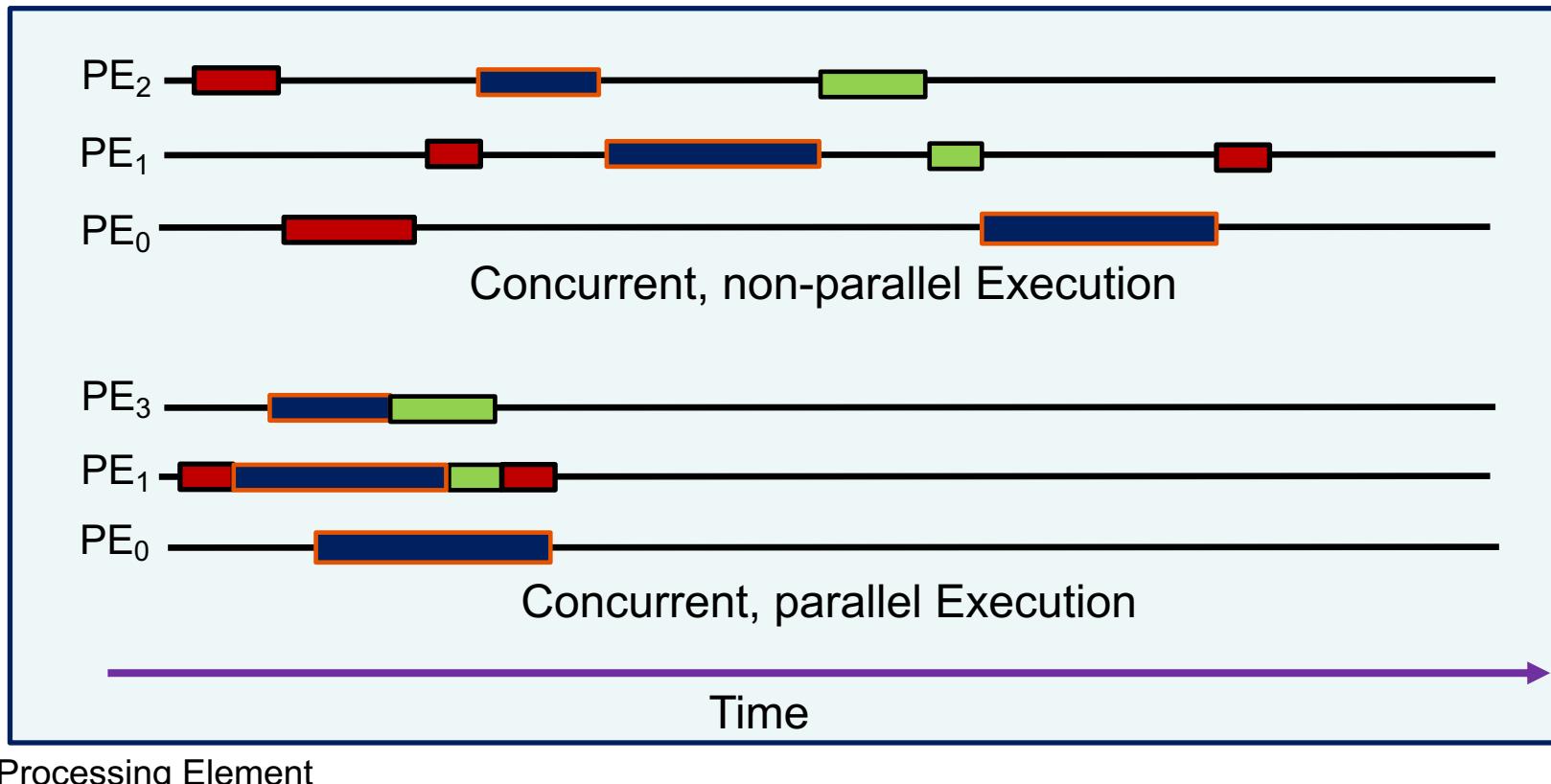
- Python is interpreted ... not compiled



- What if I want my Python program to run in parallel. Does that work?
- Not really. Python has a **Global Interpreter lock (GIL)**. This is a mutex (mutual exclusion lock) so only one thread at a time can make forward progress.

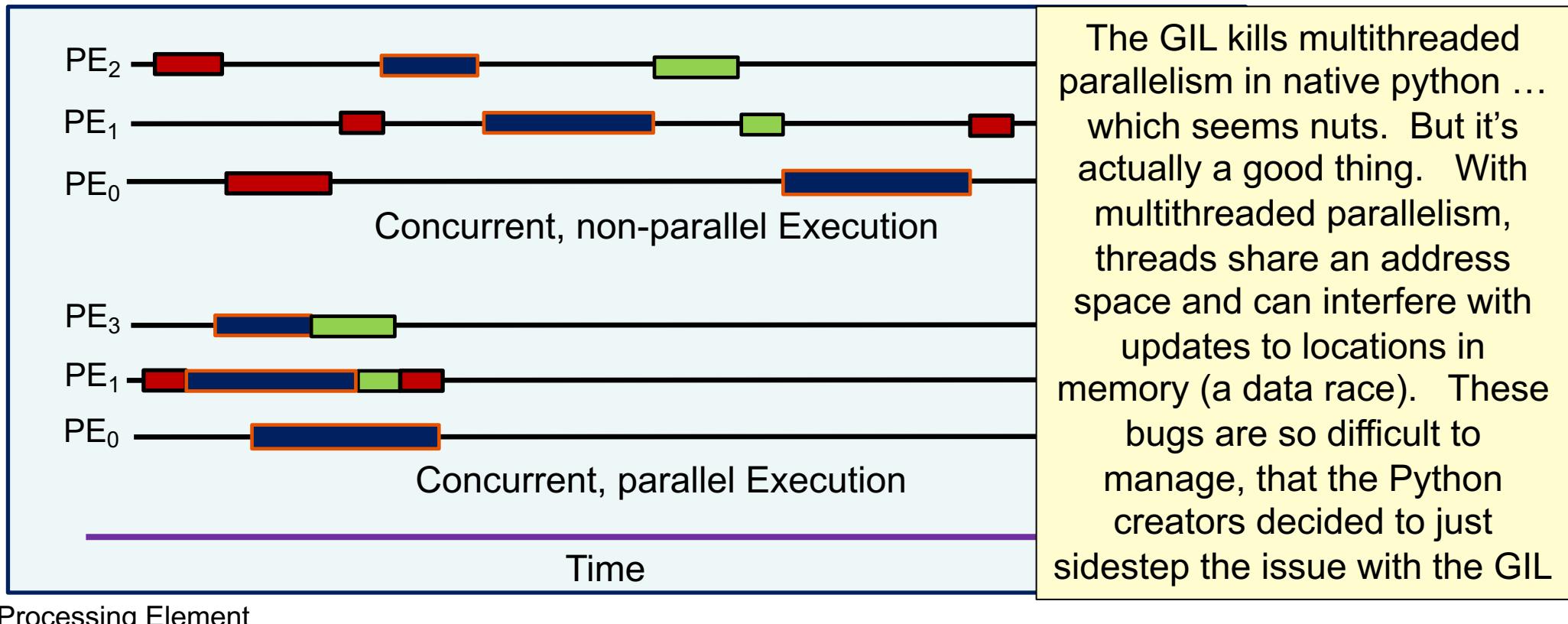
Concurrency vs. Parallelism

- Two important definitions:
 - Concurrency: A condition of a system in which multiple tasks are active and unordered. If **scheduled fairly**, they can be described as logically making **forward progress** at the same time.
 - Parallelism: A condition of a system in which multiple tasks are actually making **forward progress** at the same time.



Concurrency vs. Parallelism

- Two important definitions:
 - Concurrency: A condition of a system in which multiple tasks are active and unordered. If **scheduled fairly**, they can be described as logically making **forward progress** at the same time.
 - Parallelism: A condition of a system in which multiple tasks are actually making **forward progress** at the same time.



Exploiting parallel hardware requires parallel software.

Before looking at Parallel Programming in Python, lets pause for a moment and look at the history of parallel programming languages

Consider the state of programming models from the early days of parallel computing.

Parallel programming environments in the 90's

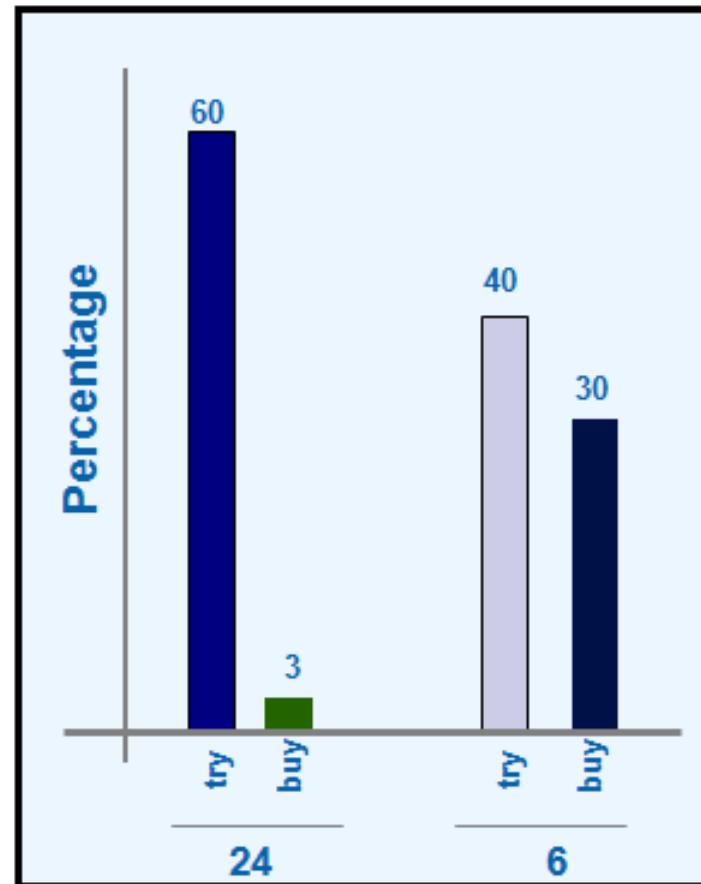
| | | | | | | | | |
|-----------------|-----------------|-------------|----------------|------------|--------------|--------------|-------------|--------------|
| ABCPL | C4 | DOLIB | HAsL. | P4-Linda | Nano-Threads | Parallel-C++ | QPC++ | Sthreads |
| ACE | CC++ | DOME | Haskell | Glenda | NESL | Parallaxis | PVM | Strand. |
| ACT++ | Chu | DOSMOS. | HPC++ | POSYBL | NetClasses++ | ParC | PSI | SUIF. |
| Active messages | Charlotte | DRL | JAVAR. | Objective- | Nexus | ParLib++ | PSDM | Synergy |
| Adl | Charm | DSM-Threads | HORUS | Linda | Nimrod | ParLin | Quake | Telegrphos |
| Adsmith | Charm++ | Ease . | HPC | LiPS | NOW | Parmacs | Quark | SuperPascal |
| ADDAP | Cid | ECO | IMPACT | Locust | Objective | Parti | Quick | TCGMSG. |
| AFAPI | Cilk | Eiffel | ISIS. | Lparx | Linda | pC | Threads | Threads.h++. |
| ALWAN | CM-Fortran | Eilean | JAVAR | Lucid | Occam | pC++ | Sage++ | TreadMarks |
| AM | Converse | Emerald | JADE | Maisie | Omega | PCN | SCANDAL | TRAPPER |
| AMDC | Code | EPL | Java RMI | Manifold | OpenMP | PCP: | SAM | uC++ |
| AppLeS | COOL | Excalibur | javaPG | Mentat | Orca | PH | pC++ | UNITY |
| Amoeba | CORRELATE | Express | JavaSpace | Legion | OOF90 | PEACE | SCHEDULE | UC |
| ARTS | CPS | Falcon | JIDL | Meta Chaos | P++ | PCU | SciTL | V |
| Athapascan-0b | CRL | Filaments | Joyce | Midway | P3L | PET | POET | ViC* |
| Aurora | CSP | FM | Khoros | Millipede | p4-Linda | PETSc | SDDA. | Visifold V- |
| Automap | Cthreads | FLASH | Karma | CparPar | Pablo | PENNY | SHMEM | NUS |
| bb_threads | CUMULVS | The FORCE | KOAN/Fortran-S | Mirage | PADE | Phosphorus | SIMPLE | VPE |
| Blaze | DAGGER | Fork | LAM | MpC | PADRE | POET. | Sina | Win32 |
| BSP | DAPPLE | Fortran-M | Lilac | MOSIX | Panda | Polaris | SISAL. | threads |
| BlockComm | Data Parallel C | FX | Linda | Modula-P | Papers | POOMA | distributed | WinPar |
| C*. | DC++ | GA | JADA | Modula-2* | AFAPI. | POOL-T | smalltalk | WWWinda |
| "C* in C | DCE++ | GAMMA | WWWinda | Multipol | Para++ | PRESTO | SMI. | XENOOPS |
| C** | DDD | Glenda | ISETL-Linda | MPI | Paradigm | P-RIO | SONiC | XPC |
| CarlOS | DICE. | GLU | ParLin | MPC++ | Parafrase2 | Prospero | Split-C. | Zounds |
| Cashmere | DIPC | GUARD | Eilean | Munin | Paralation | Proteus | SR | ZPL |

Third party names are the property of their owners.

**A warning I've
been making for
the last 25 years**

Is it bad to have so many languages? Too many options can hurt you

- The Draeger Grocery Store experiment consumer choice:
 - Two Jam-displays with coupon's for purchase discount.
 - 24 different Jam's
 - 6 different Jam's
 - How many stopped by to try samples at the display?
 - Of those who "tried", how many bought jam?



Programmers don't need a glut of options ... just give us something that works OK on every platform we care about. Give us a decent standard and we'll do the rest

The findings from this study show that an extensive array of options can at first seem highly appealing to consumers, yet can reduce their subsequent motivation to purchase the product.

Iyengar, Sheena S., & Lepper, Mark (2000). When choice is demotivating: Can one desire too much of a good thing? *Journal of Personality and Social Psychology*, 76, 995-1006.

Parallel programming environments: post-90s

- The application community (with leadership from the Accelerated Strategic Computing Initiative) pushed for convergence around a small number of programming languages:
 - For clusters and massively parallel computers: **MPI**
 - For shared memory systems: **OpenMP**
- With only two languages, vendors could focus on engineering high quality solutions ... rather than chasing the latest fad.
- All was good until ~2006 when fully programmable GPUs came along. We are still sorting out what will become the converged solution ...
 - **Cuda, OpenCL, Sycl, OpenACC, OpenMP** ← I know, this is a bit of a mess. Sorry about that

How about Parallel programming with Python

| | | |
|------------|-------------------|---|
| dispy | Dask | pyPastSet |
| Delegate | Deap | pypvm |
| forkmap | disco | pynpvm |
| forkfun | dispy | Pyro |
| Jobibppmap | DistributedPYthon | Ray |
| POSH | exec_proxy | Rthread |
| pp | execnet | ScientificPython.BSP |
| pprocess | iPython | Scientific.DistrubedComputing.MasterSlave |
| processing | job_stream jug | Scientific.MPI |
| PyCSP | mi4py | SCOOP |
| PyMP | NetWorkSpaces | seppo |
| Ray | PaPy | PySpark |
| remoteD | papyrus | Star-P |
| torcp | PyCOMPSSs | superrpy |
| VecPy | PyLinda | torcpy |
| batchlib | pyMPI | StarCluster |
| Celery | pypyar | dpctl |
| Charm4py | multiprocessing | arkouda |
| PyCUDA | PyOpenCL | PyOMP |
| Ramba | | dnpn |

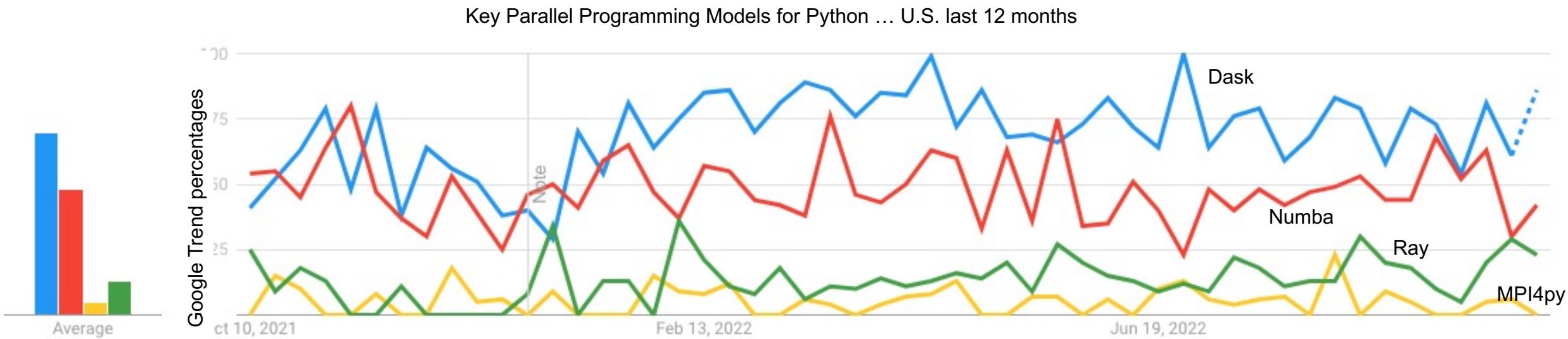
How about Parallel programming with Python

| | | |
|------------|--|-------------|
| dispy | Dask | pyPastSet |
| Delegate | Deap | pypvm |
| forkmap | disco | pynpvm |
| forkfun | dispy | Pyro |
| Jobibppmap | DistributedPYthon | Ray |
| POSH | exec proxy | Rthread |
| pp | We are still early (compared to HPC) in the evolution of parallel programming models for Python. | |
| pprocess | Hopefully, soon the python application community will come together and help us narrow down to a handful of systems to focus on. | |
| processing | That would allow vendors to carry out HW/SW optimization and focus on quality over "chasing fads". | |
| PyCSP | PyLinda | StarCluster |
| PyMP | puMPI | dpctl |
| Ray | pypar | arkouda |
| remoteD | multiprocessing | PyOMP |
| torcp | PyOpenCL | dppnp |
| VecPy | | |
| batchlib | | |
| Celery | | |
| Charm4py | | |
| PyCUDA | | |
| Ramba | | |

Popular python parallel Programming models

We compared many python parallel programming models with google-trends (which tracks web searches)

These four systems are popular and (in our opinion) are the key systems to consider

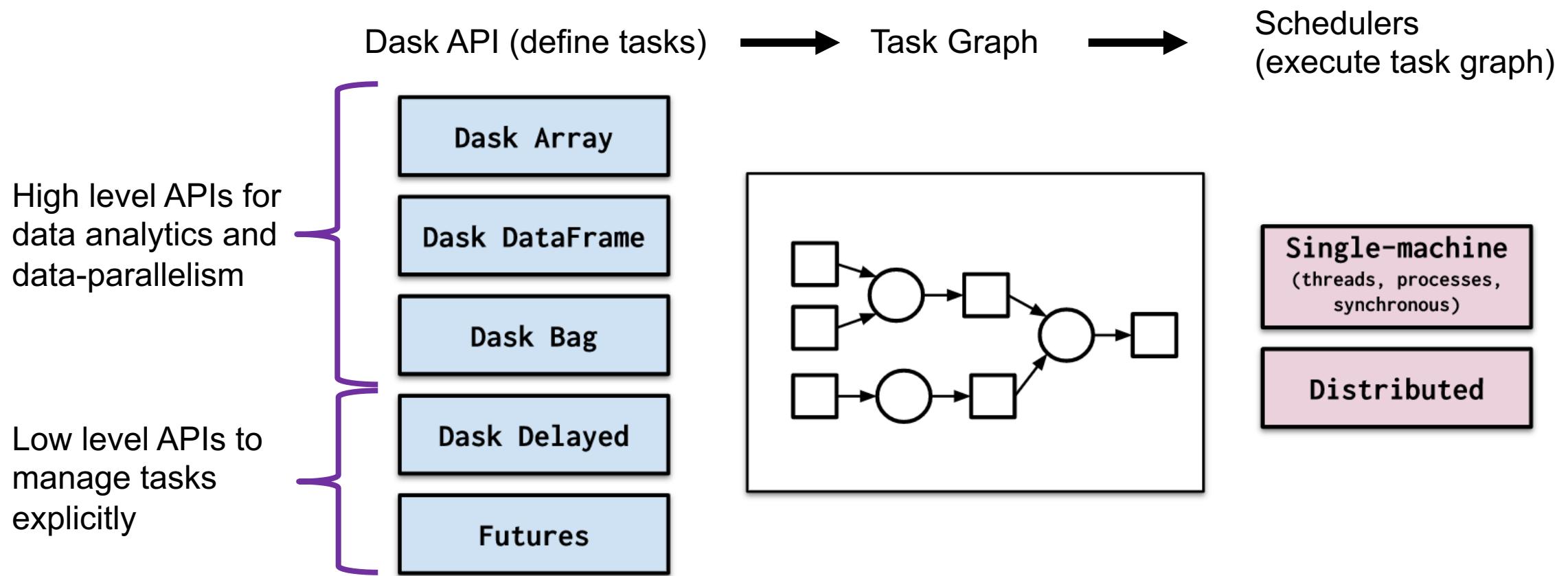


Most people writing parallel code in Python use Dask.

**That's because Dask was designed around the needs
of data scientists, and lets face, many python
programmers work on data science problems.**

DASK

- Parallel and distributed computing library for Python
- Client / driver submits tasks to Dask cluster (set of worker processes on one or more physical nodes)



Dask Delayed – lazy, remote functions

- Define a remote function:

```
@dask.delayed  
def add_one(i):  
    time.sleep(1)  
    return i+1
```

Decorator turns normal Python function into Dask lazy function

- Calling remote function, getting results:

```
futurevalue = add_one(7)  
v = futurevalue.compute()
```

Returns immediately after creating task in task graph

Triggers execution of task graph
Returns value 8 after about 1 second when task completes

Dask – parallel and chaining calls

- Parallel execution:

```
fv = [add_one(i) for i in range(5)]  
v = sum(fv)  
v = v.compute()
```

Returns immediately with a list of “futures”

Standard Python sum function;
Returns immediately with a future

- Chained execution:

```
v = 2  
for x in range(5):  
    v = add_one(v)  
v = v.compute()
```

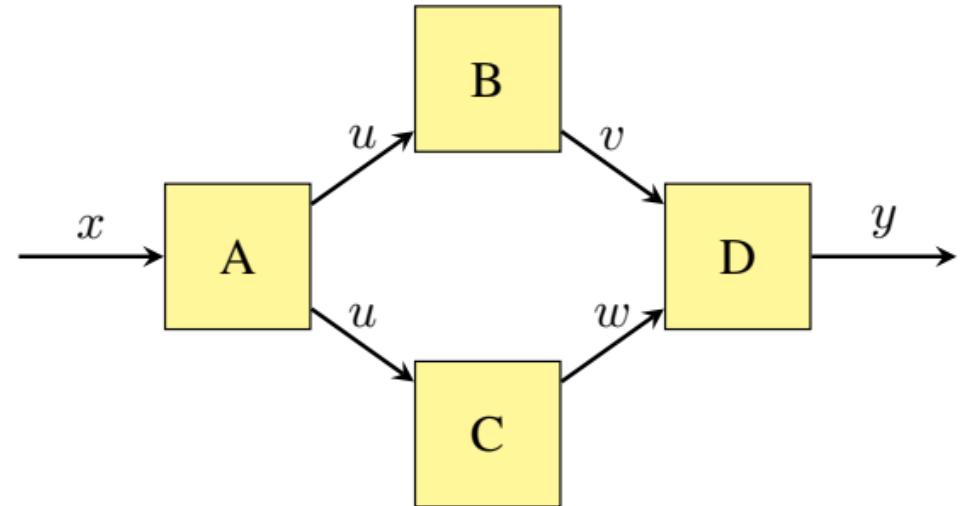
Returns value 15 after about 1 second

These return immediately

Returns value 7 after about 5 seconds

Chaining forms DAGs of Tasks

```
# A, B, C, and D are delayed functions  
u = A(x)  
v = B(u)  
w = C(u)  
y = D(v, w)  
y = y.compute()
```



Dask Futures

- Same concept, but eager asynchronous execution, different syntax

```
def add_one(i):  
    time.sleep(1)  
    return i+1
```

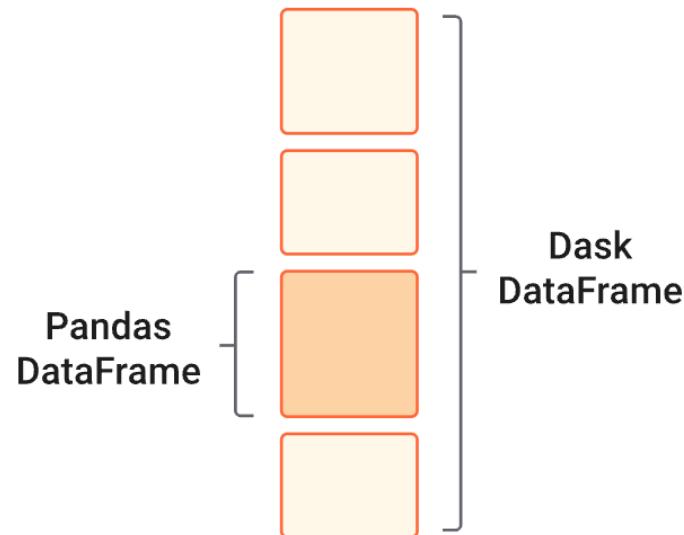
```
future = client.submit(add_one, 3) ← Submits add_one(3) for  
result = future.result()           distributed execution
```

← Returns value 4 in about 1 second

- Note: no decorator, explicit job submission
- Can pass futures as parameters to chain functions/construct DAGs

Dask DataFrames

- The Dask users I know use it to migrate single node programs using Pandas into distributed programs using Dask DataFrames.



- The Dask DataFrame coordinates pandas DataFrames running on individual cluster nodes.
- Distribute rows by index (primary key)

```
>>> import pandas as pd  
  
>>> df = pd.read_parquet('s3://mybucket/myfile.parquet')  
>>> df.head()  
0 1 a  
1 2 b  
2 3 c
```



```
>>> import dask.dataframe as dd  
  
>>> df = dd.read_parquet('s3://mybucket/myfile.*.parquet')  
>>> df.head()  
0 1 a  
1 2 b  
2 3 c
```

- Data Frames supports huge data sets ... embarrassingly parallel combination of disk bandwidth across a cluster

Dask DataFrames

- API matches pandas hence, if you know pandas, you pretty much know how to user Dask DataFrames

```
import pandas as pd
```

```
>>> df = df[df.value >= 0]
>>> joined = df.merge(other, on="account")
>>> result = joined.groupby("account").value.mean()

>>> result
alice 123
bob 456
```

```
import dask.dataframe as dd
```

```
>>> df = df[df.value >= 0]
>>> joined = df.merge(other, on="account")
>>> result = joined.groupby("account").value.mean()

>>> result.compute()
alice 123
bob 456
```

- you trigger computation by calling the `.compute()` method or persist data in distributed memory with the `.persist()` method.

Dask Array

- High-level API provides distributed, Numpy-like array interface
- Arrays partitioned into chunks – serves as unit of storage and computation
- Arrays can be disk-backed, and thus larger than memory
- Array operations are lazy, internally constructing DAG of operations
- Explicit triggering of execution using `compute()` method
 - Parallel execution of relevant portions of task graph on Dask cluster
 - Computation at chunk granularity
 - Only necessary chunks computed for requested result

Dask Array example

```
import dask.array as da
```

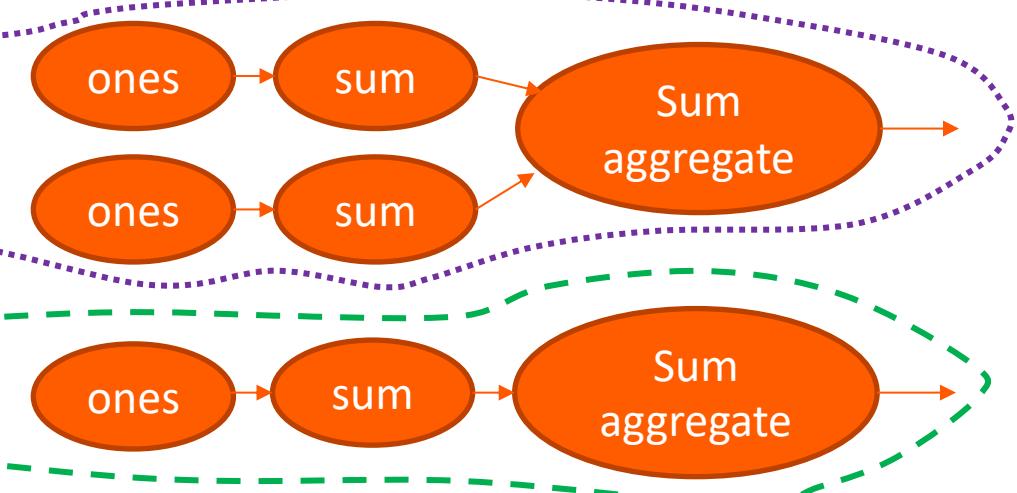
- `A = da.ones((1000,1000),chunks=(1000,500))`
 - Constructs 1000x1000 array, with two chunks of size 1000x500

- `B = da.sum(A, axis=0)`
 - Sum along axis 0 → should produce a 1000 element array

- `B.compute()`
 - Triggers computation of DAG:
 - Parallel execution on chunks

- `B[0].compute()`
 - Only compute chunks needed for `B[0]`

- Typically, Dask will not materialize a derived array
 - Keeps the DAG that describes how to compute it
 - May need to recompute (but may cache results as well)
 - Optimized for computations on disk-based data that won't fit in memory
- `Persist()` method to force computation, materialization of an array



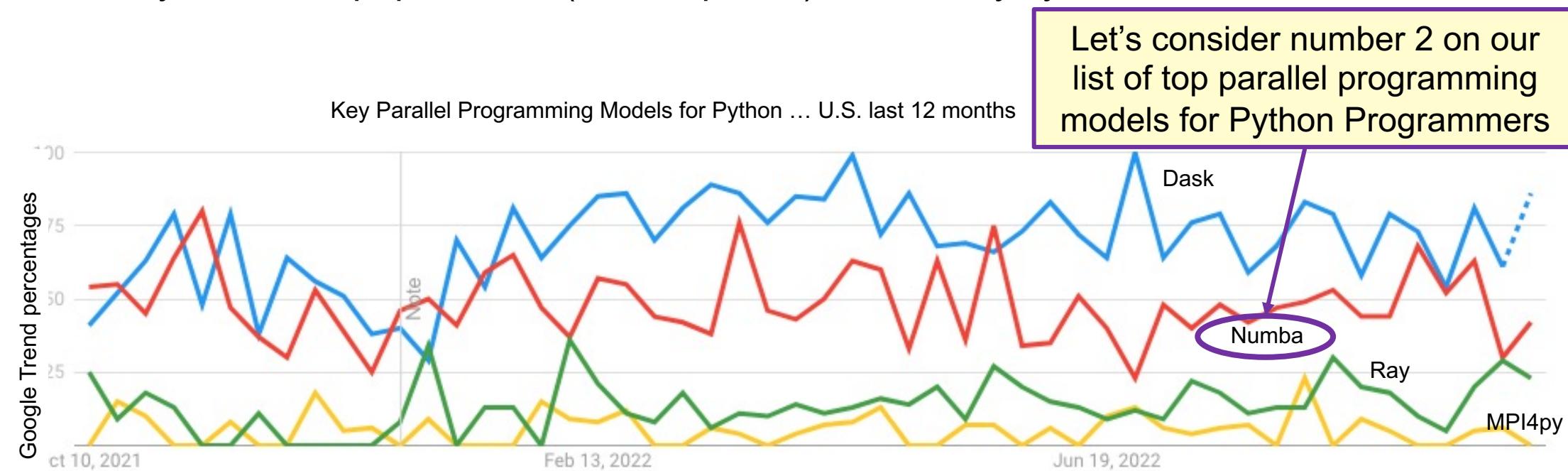
Dask is great ... and if you are interested in Distributed Data Science problems and a pandas API, it is the right choice.

But it doesn't handle a wide range of parallel systems and has too much overhead to support a wide range of algorithms

Popular python parallel Programming models

We compared many python parallel programming models with google-trends (which tracks web searches)

These four systems are popular and (in our opinion) are the key systems to consider



Numba with ParallelAccelerator

Numba ... C-like performance from Python code



- Numba is a JIT compiler. Maps a subset of python with numpy arrays onto LLVM
 - Once code is JIT'ed into LLVM, all performance enhancements exposed at the level of LLVM are directly available ... result is performance that approaches that from raw C or Fortran
 - Source code is pure python for maximum portability
-
- Just add the `@jit` decorator to enable numba for a function.

```
from numba import jit

@jit
def addit(A,B):
    return (A+B)
```

Numba jit compiler applied the first time a function is encountered. Caches the code so subsequent calls to the function don't run the jit step.

Numba defines elementwise functions called **ufuncs**

This generates the LLVM code and calls the addition **ufunc** to do an elementwise add of A and B

- Numerous options in numba ... we are barely scratching the surface
 - `@jit(nopython = true)` tells the system to NOT use any python objects in the generated code. Can be much faster
 - `@jit(parallel = true)` invoke parallel accelerator

Parallel Accelerator

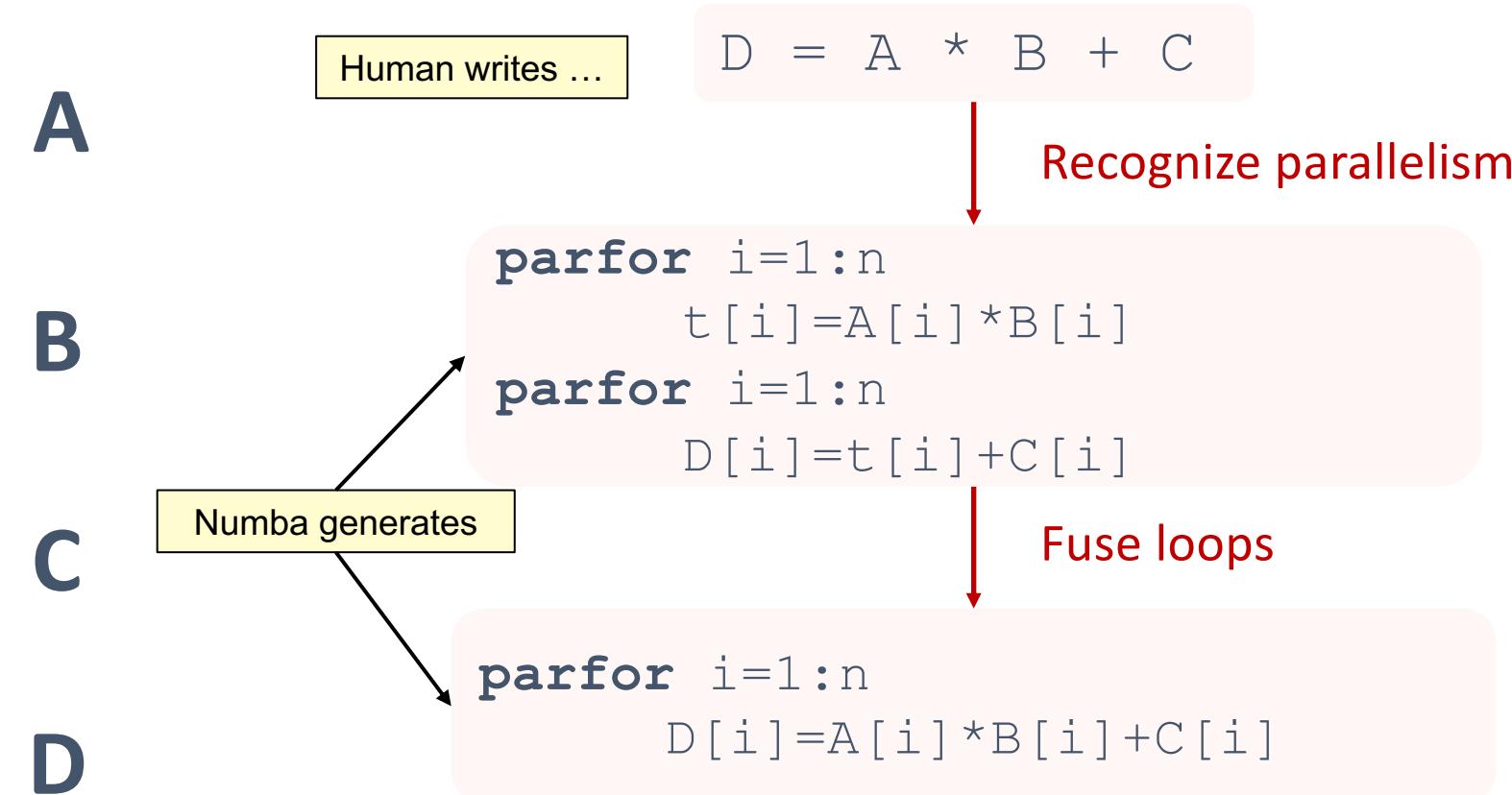
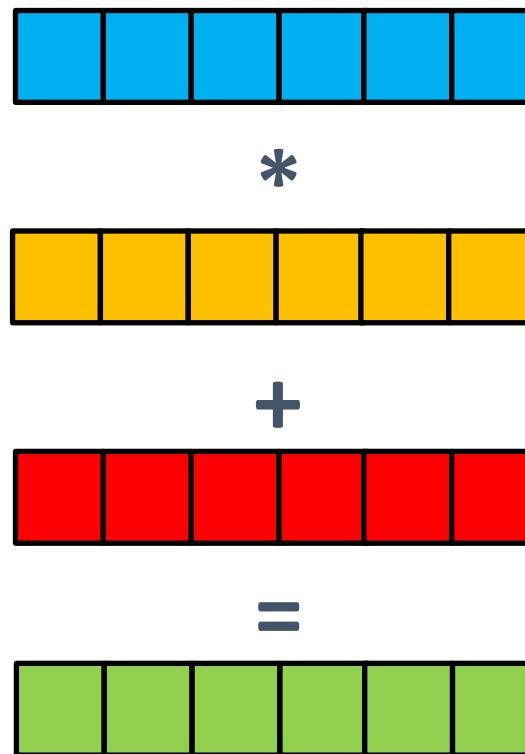


Works with numba to JIT code that executes in parallel. It does the following:

1. Recognize parallelism.
 - Pattern recognition of operations with concurrent semantics.
2. Represent parallelism.
 - Numba's parfor node* – represents a strictly nested set of for loops known to have no cross-iteration dependencies.
3. Optimizations.
 - Fusion – combine compatible parfors together. Eliminates unnecessary temporary arrays and traverses arrays only once for better cache utilization.
4. Run in parallel.
 - Improves performance by leveraging multiple cores and vector instructions.

*Note: a parfor node refers to the intermediate code generated by Numba. This is NOT a human callable API.

Transformation carried out for array-based data parallelism



We will talk about writing explicit parallel loops later, you don't put parfor's in your python code

ParallelAccelerator – Softmax program



```
import numba

@numba.njit(parallel=True)
def sigArr(A):
    Amax = np.max(A)
    Ashift = A - Amax
    expAshift = np.exp(Ashift)
    Normalization = np.sum(expAshift)
    reciNorm = 1/Normalization
    sigma = expAshift*reciNorm
    return sigma
```

- Same as the NumPy version.
- np.max executed in one parallel region.
- Subtraction, exp, and sum fused into one parallel region.
- Ashift temporary eliminated.
- expAshift * reciNorm the final parallel region.

Recognizing Parallelism



The following patterns are recognized by ParallelAccelerator for parallel execution:

1. Implicit

- Element-wise operations: unary(+,-,~,), binary(+,-,*/,//?,%,|,>>,^,&,**,//), comparison(==,!=<,<=,>,>=), NumPy ufuncs, user-defined DUFunc.
- NumPy reductions: sum, prod, min, max, argmin, argmax, mean, var, std.
- Array creation: zeros, ones, arrange, linspace, and random array create for all available distributions.
- NumPy dot: matrix/vector or vector/vector.
- Array assignment.
- Functools.reduce.
- Stencil decorator.

2. Explicit

- prange, pndindex

```
Import numba

@numba.njit(parallel=True)
def fun(A):
    for i in numba.prange(aa.size):
        A[i]=A[i]*2
    return A
```

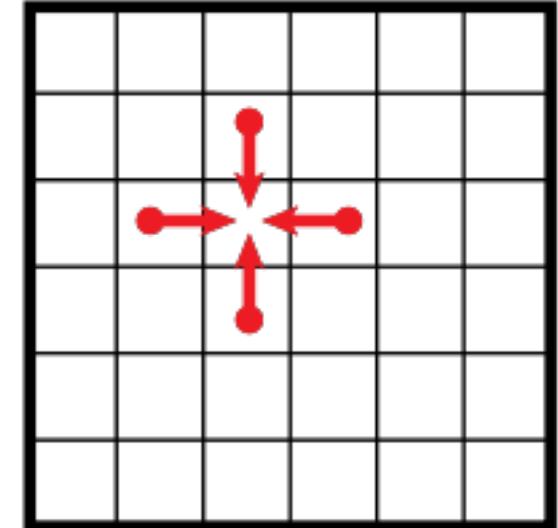
Other ParallelAccelerator Technology



- Stencils are very common in scientific computing.
- ParallelAccelerator provides a productive stencil abstraction with automatic parallelization.

```
@stencil
def jacobi_kernel(a):
    return 0.25 * (a[0,1] + a[0,-1] + a[-1,0] + a[1,0])
```

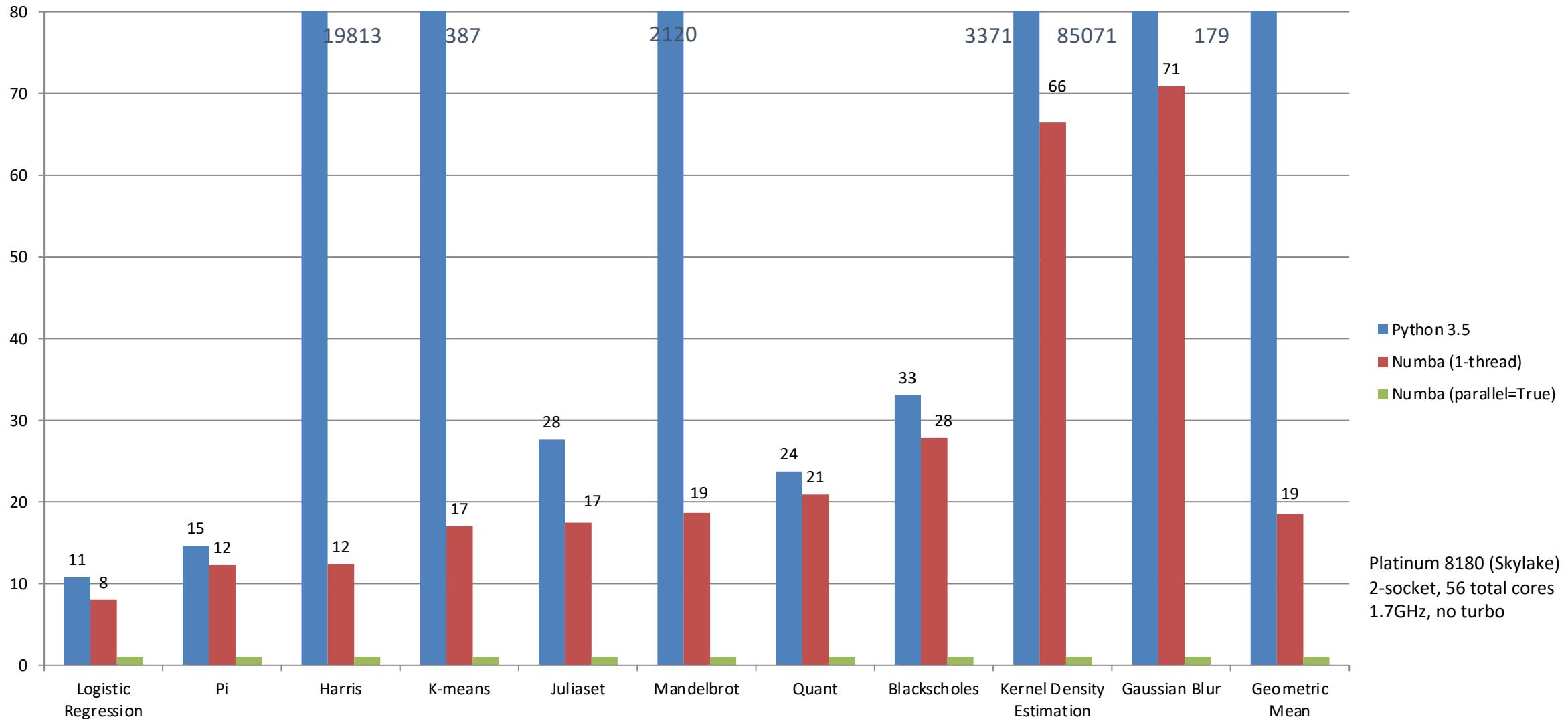
```
@numba.njit(parallel=True)
def run_jacobi(a):
    return jacobi_kernel(a)
```



Performance



Kernel Times Relative to Numba with parallel=True (lower is better)



**Dask and Numba parallel accelerator only address
a small fraction of parallel algorithms.**

**We need a more general parallel programming API
to use in Python.**

**The ubiquitous standard for multithreaded
programming on CPUs is OpenMP**

OpenMP* Overview

OpenMP: An API for Writing Multithreaded Applications

- A set of compiler directives and library routines for parallel application programmers
- Greatly simplifies writing multi-threaded (MT) programs in Fortran, C and C++
- Standardizes established SMP practice + vectorization and heterogeneous device programming

C\$OMP PARALLEL COPYIN(/blk/)

C\$OMP DO lastprivate(XX)

Nthrds = OMP_GET_NUM_PROCS()

omp_set_lock(lck)

PyOMP: OpenMP projected into Python

- A parallel multithreaded “hello world” program with PyOMP

```
from numba import njit
from numba.openmp import openmp_context as openmp

@njit
def hello():
    with openmp("parallel"):
        print("hello")
        print("world")

    hello()
    print("DONE")
```

PyOMP: OpenMP projected into Python

- A parallel multithreaded “hello world” program with PyOMP

Numba Just In Time (JIT) compiler compiles the Python code into LLVM.

Compiled code cached for later use.

```
from numba import njit
from numba.openmp import openmp_context as openmp

@njit
def hello():
    with openmp("parallel"):
        print("hello")
        print("world")
    hello()
    print("DONE")
```

“parallel” creates a team of threads

The code inside the with context manager is packaged into a function and executed by each thread

OpenMP managed through the *with* context manager.

- Numba Just In Time (JIT) compiler compiles the Python code into LLVM thereby bypassing the GIL. Hence, the threads execute in parallel.
- The string in the with openmp context manager is identical to the constructs in OpenMP. If you know OpenMP for C/C++/Fortran, then you know it for Python

PyOMP: OpenMP projected into Python

When I run this program,
here is the output.

- A parallel multithreaded “hello world” program with PyOMP

```
from numba import njit
from numba.openmp import openmp_context as openmp

@njit
def hello():
    with openmp("parallel"):
        print("hello")
        print("world")

    hello()
    print("DONE")
```

hello
world
hello
hello
hello
world
hello
world
hello
world
hello
world
hello
world
world
world
hello
world
DONE

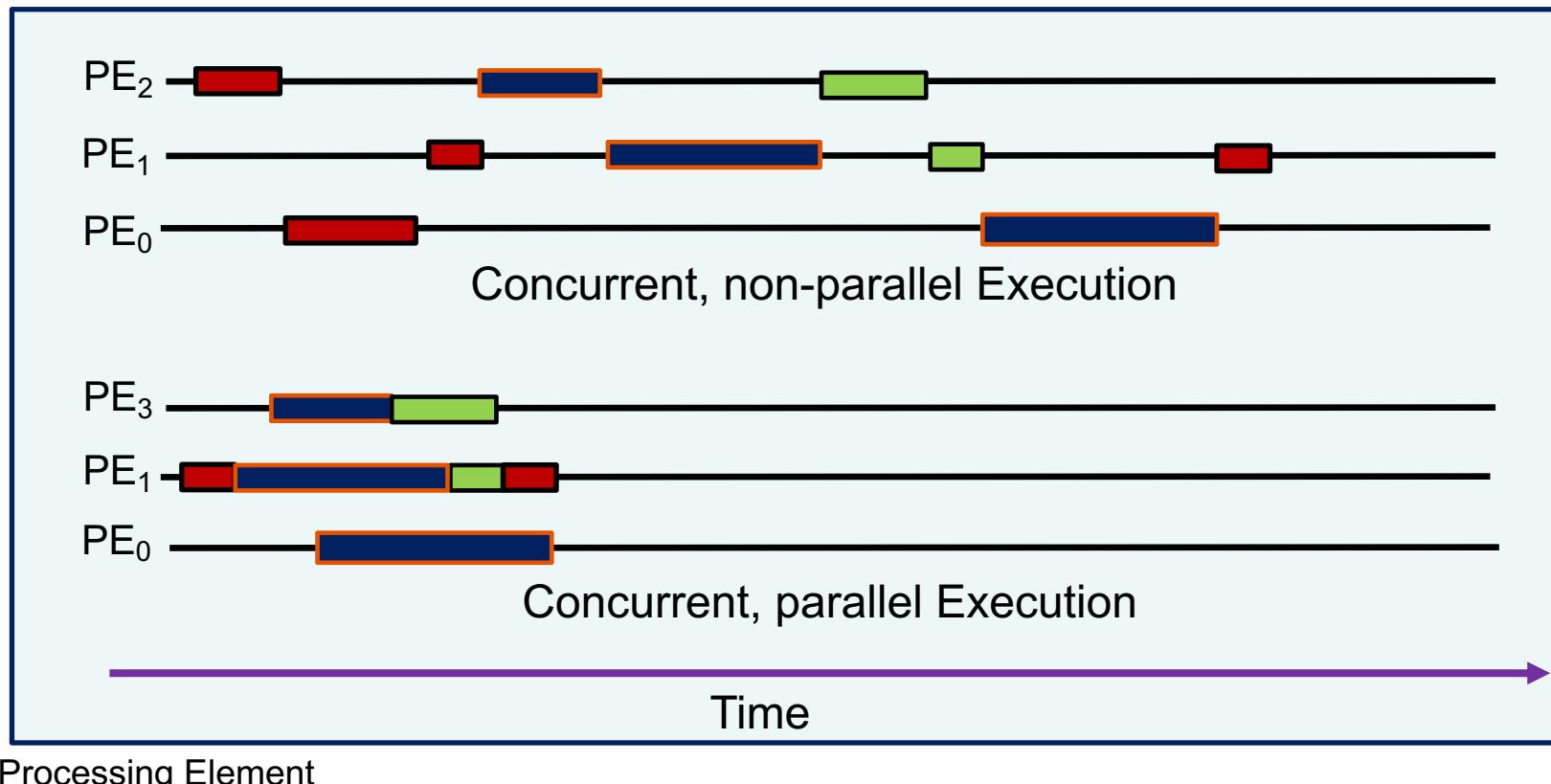
The interleaved print
output is different each
time I run the program

Why is the output from our hello world program so weird?

To answer that question, we must digress briefly and settle on a few key definitions

Concurrency vs. Parallelism

- Two important definitions:
 - Concurrency: A condition of a system in which multiple tasks are active and unordered. If **scheduled fairly**, they can be described as logically making **forward progress** at the same time.
 - Parallelism: A condition of a system in which multiple tasks are actually making **forward progress** at the same time.



PyOMP: OpenMP projected into Python

When I run this program,
here is the output.

- A parallel multithreaded “hello world” program with PyOMP

```
from numba import njit
from numba.openmp import openmp_context as openmp

@njit
def hello():
    with openmp("parallel"):
        print("hello")
        print("world")

    hello()
    print("DONE")
```

hello
world
hello
hello
hello
world
hello
world
hello
world
hello
world
hello
world
world
world
hello
world
DONE

The challenge for programmers writing multithreaded code is to make sure every semantically allowed way statements can interleave results in correct code.

How did we implement PyOMP?

We used the “magic” of Numba

PyOMP Implementation in Numba

- PyOMP changes to Numba:
 - Adds an OpenMP context manager
 - Provides the ability to call all the OpenMP runtime functions from both Python and Numba JITed code.
- Exception handling disabled in OpenMP regions since Numba exception mechanism breaks OpenMP single-entry/single-exit requirement.
- Variables not listed in a data clauses are SHARED if used before or after OpenMP region, PRIVATE otherwise*.
- Supports most OpenMP 3.5 and much of OpenMP 4.5. Supported directives and clauses can be found at <https://pyomp.readthedocs.io/en/latest/>.
- Note that one can use `@jit(cache=True)` Numba decorator to compile the function once and store the result on disk to avoid recompilation each time the program is restarted.

* The scope of shared/private variables exposes subtle issues in how the rules for an OpenMP data environment interacts with how Numba manages the visibility of variables. This is a topic that is still evolving, though in practice it hasn't impacted the usability of PyOMP .

How do you install PyOMP on your own system?

PyOMP installation

- Preferred installation method is through conda. Running python from 3.8 to 3.10:
 - `conda install -c python-for-hpc -c conda-forge pyomp`
- We currently support PyOMP on four systems
 - linux-ppc6le
 - linux-64 (x86_64)
 - osx-arm64 (mac)
 - linux-arm64
- We also have a working (free) JupyterLab under binder for OpenMP CPU at:
 - <https://mybinder.org/v2/gh/Python-for-HPC/binder/HEAD>

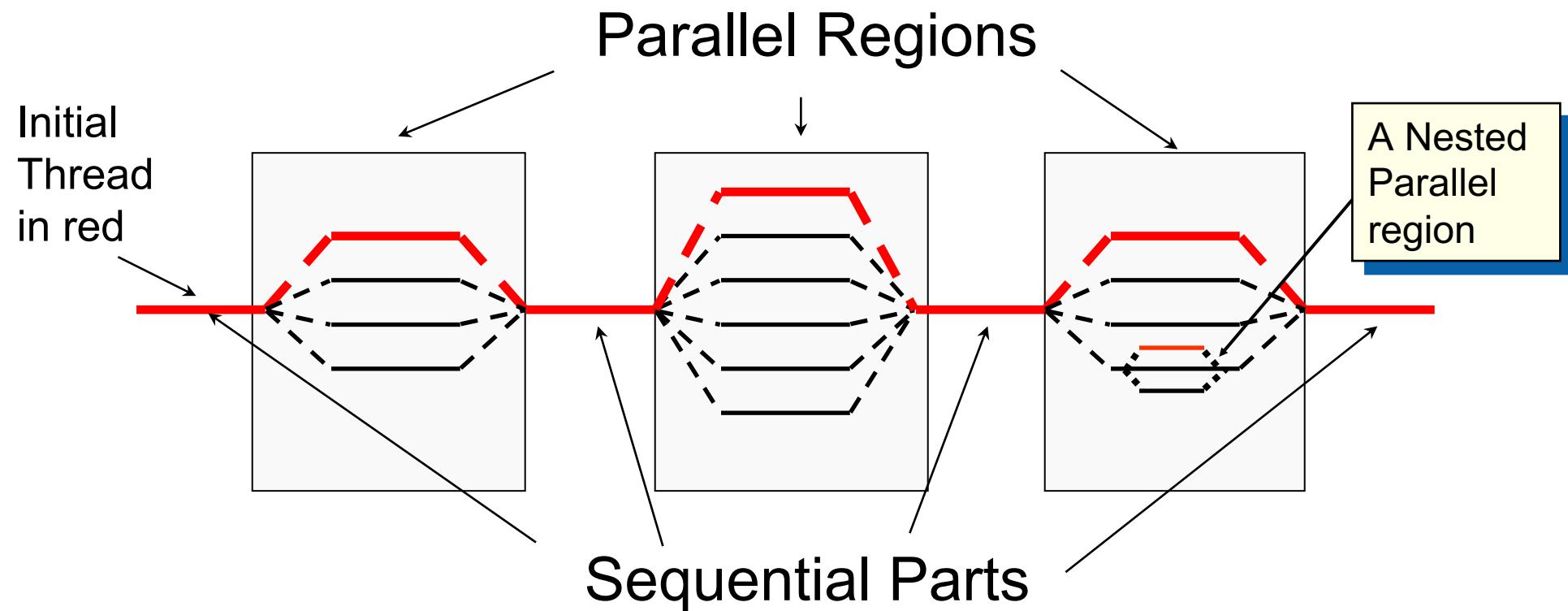
<https://github.com/Python-for-HPC/PyOMP>

**Lets dive into the details of
multithreading and how they are most
commonly used in an application**

OpenMP Execution Model

Fork-Join Parallelism:

- Initial thread **forks** a team of threads as needed.
- They execute in a shared address space ... All reads read/write a common set of the variables.
- When the team is finished, the threads **join** together and the initial thread continues
- Parallelism added incrementally until performance goals are met, i.e., the sequential program evolves into a parallel program.



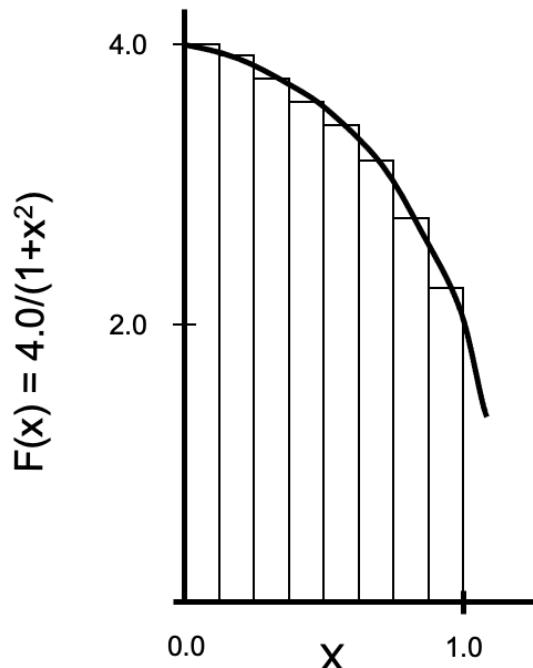
Understanding OpenMP

We will explain the key elements of OpenMP as we explore the three fundamental design patterns of OpenMP (**Loop parallelism**, SPMD, and **divide and conquer**) applied to the following problem

Numerical Integration (the *hello world* program of parallel computing)

Mathematically, we know that:

$$\int_0^1 \frac{4.0}{(1+x^2)} dx = \pi$$



We can approximate the integral as a sum of rectangles:

$$\sum_{i=0}^N F(x_i) \Delta x \approx \pi$$

Each rectangle: width Δx , height $F(x_i)$ at i^{th} interval midpoint.

```
def piFunc(NumSteps):
    step=1.0/NumSteps
    pisum = 0.0
    x = 0.5
    for i in range(NumSteps):
        x+=step
        pisum += 4.0/(1.0+x*x)
    pi=step*pisum
    return pi
```

The Loop-level parallelism design pattern

- Parallelism defined in terms of parallel loops ... that is, loops where iterations can safely execute when divided between a collection of threads.
- Key elements:
 - identify compute intensive loops in a program
 - Expose concurrency by removing or managing loop carried dependencies
 - Exploit concurrency for parallel execution usually using a parallel loop construct/directive.

The Loop-level parallelism design pattern

- Parallelism defined in terms of parallel loops ... that is, loops where iterations can safely execute when divided between a collection of threads.
- Key elements:
 - identify compute intensive loops in a program
 - Expose concurrency by removing or managing loop carried dependencies
 - Exploit concurrency for parallel execution usually using a parallel loop construct/directive.

```
def piFunc(NumSteps):  
    step=1.0/NumSteps  
    pisum = 0.0  
    x = 0.5  
    for i in range(NumSteps):  
        x+=step  
        pisum += 4.0/(1.0+x*x)  
    pi=step*pisum  
    return pi
```

The Loop-level parallelism design pattern

- Parallelism defined in terms of parallel loops ... that is, loops where iterations can safely execute when divided between a collection of threads.
- Key elements:
 - identify compute intensive loops in a program
 - Expose concurrency by removing or managing loop carried dependencies
 - Exploit concurrency for parallel execution usually using a parallel loop construct/directive.

```
def piFunc(NumSteps):  
    step=1.0/NumSteps  
    pisum = 0.0  
    x = 0.5  
    for i in range(NumSteps):  
        x+=step  
        pisum += 4.0/(1.0+x*x)  
    pi=step*pisum  
    return pi
```

A loop carried dependency

The Loop-level parallelism design pattern

- Parallelism defined in terms of parallel loops ... that is, loops where iterations can safely execute when divided between a collection of threads.
- Key elements:
 - identify compute intensive loops in a program
 - Expose concurrency by removing or managing loop carried dependencies
 - Exploit concurrency for parallel execution usually using a parallel loop construct/directive.

```
def piFunc(NumSteps):
    step=1.0/NumSteps
    pisum = 0.0
    x = 0.5
    for i in range(NumSteps):
        x+=step
        pisum += 4.0/(1.0+x*x)
        pi=step*pisum
    return pi
```

A loop carried dependency → Recast to compute from i

```
def piFunc(NumSteps):
    step=1.0/NumSteps
    pisum = 0.0

    for i in range(NumSteps):
        x=(i+0.5)*step
        pisum += 4.0/(1.0+x*x)
        pi=step*pisum
    return pi
```

The Loop-level parallelism design pattern

- Parallelism defined in terms of parallel loops ... that is, loops where iterations can safely execute when divided between a collection of threads.
- Key elements:
 - identify compute intensive loops in a program
 - Expose concurrency by removing or managing loop carried dependencies
 - Exploit concurrency for parallel execution usually using a parallel loop construct/directive.

```
def piFunc(NumSteps):  
    step=1.0/NumSteps  
    pisum = 0.0  
    x = 0.5  
    for i in range(NumSteps):  
        x+=step  
        pisum += 4.0/(1.0+x*x)  
    pi=step*pisum  
    return pi
```

A loop carried dependency

Recast to compute from i

```
def piFunc(NumSteps):  
    step=1.0/NumSteps  
    pisum = 0.0  
  
    for i in range(NumSteps):  
        x=(i+0.5)*step  
        pisum += 4.0/(1.0+x*x)  
    pi=step*pisum  
    return pi
```

This dependency is more complicated. It's called a reduction

Loop Parallelism code

```
from numba import njit
from numba.openmp import openmp_context as openmp      OpenMP managed through the with context manager.

@njit                                                 Numba Just In Time (JIT) compiler compiles the Python code into
def piFunc(NumSteps):                                LLVM thereby bypassing the GIL. Compiled code cached for
    step = 1.0/NumSteps                            later use.

    pisum = 0.0

with openmp ("parallel for private(x) reduction(+:pisum)": Pass the OpenMP directive into the OpenMP context
    for i in range(NumSteps):                      manager as a string
        x = (i+0.5)*step
        pisum += 4.0/(1.0 + x*x)

    pi = step*pisum
    return pi

pi = piFunc(100000000)
```

- **parallel**: creates a team of threads
- **for**: maps loop iterations onto threads.
- **private(x)**: each threads gets its own x
- Loop control index of a parallel for (**i**) is private to each thread.
- **reduction(+:sum)**: combine sum from each thread using +

Reduction

- OpenMP reduction clause added to a **parallel for**:
reduction (op : list)
- Inside the **parallel for**:
 - Each thread gets a private copy of each variable in **list** ... initialized depending on the “op”
(e.g., 0 for “+”).
 - Updates to the reduction variable from each thread happens to its private copy.
 - The private copies from each thread are combined into a single value ... and then combined with the original global value ... all using the **op** from the reduction clause.
- The variables in the “list” must be shared in the enclosing parallel region.

```
from numba import njit
from numba.openmp import openmp_context as openmp

@njit
def piFunc(NumSteps):
    step = 1.0/NumSteps
    pisum = 0.0

    with openmp ("parallel for private(x) reduction(+:pisum)"):
        for i in range(NumSteps):
            x = (i+0.5)*step
            pisum += 4.0/(1.0 + x*x)

    pi = step*pisum
    return pi

pi = piFunc(100000000)
```

Numerical Integration results in seconds ... lower is better

| Threads | PyOMP | | C | |
|---------|--------|--|--------|--|
| | Loop | | Loop | |
| 1 | 0.447 | | 0.444 | |
| 2 | 0.252 | | 0.245 | |
| 4 | 0.160 | | 0.149 | |
| 8 | 0.0890 | | 0.0827 | |
| 16 | 0.0520 | | 0.0451 | |

10^8 steps

Intel® Xeon® E5-2699 v3 CPU with 18 cores running at 2.30 GHz.

For the C programs we used Intel®icc compiler version 19.1.3.304 as `icc -qnextgen -O3 –fopenmp`

Ran each case 5 times and kept the minimum time. **JIT time is not included** for PyOMP (it was about 1.5 seconds)

**Parallel Loop are great ... but sometimes
you want more control over individual
threads**

Understanding OpenMP

We will explain the key elements of OpenMP as we explore the three fundamental design patterns of OpenMP (**Loop parallelism**, **SPMD**, and **divide and conquer**) applied to the following problem

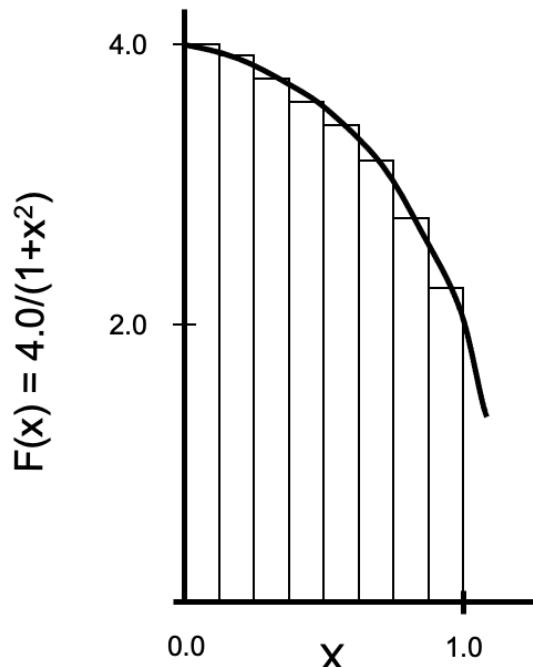
Numerical Integration (the *hello world* program of parallel computing)

Mathematically, we know that:

$$\int_0^1 \frac{4.0}{(1+x^2)} dx = \pi$$

We can approximate the integral as a sum of rectangles:

$$\sum_{i=0}^N F(x_i) \Delta x \approx \pi$$

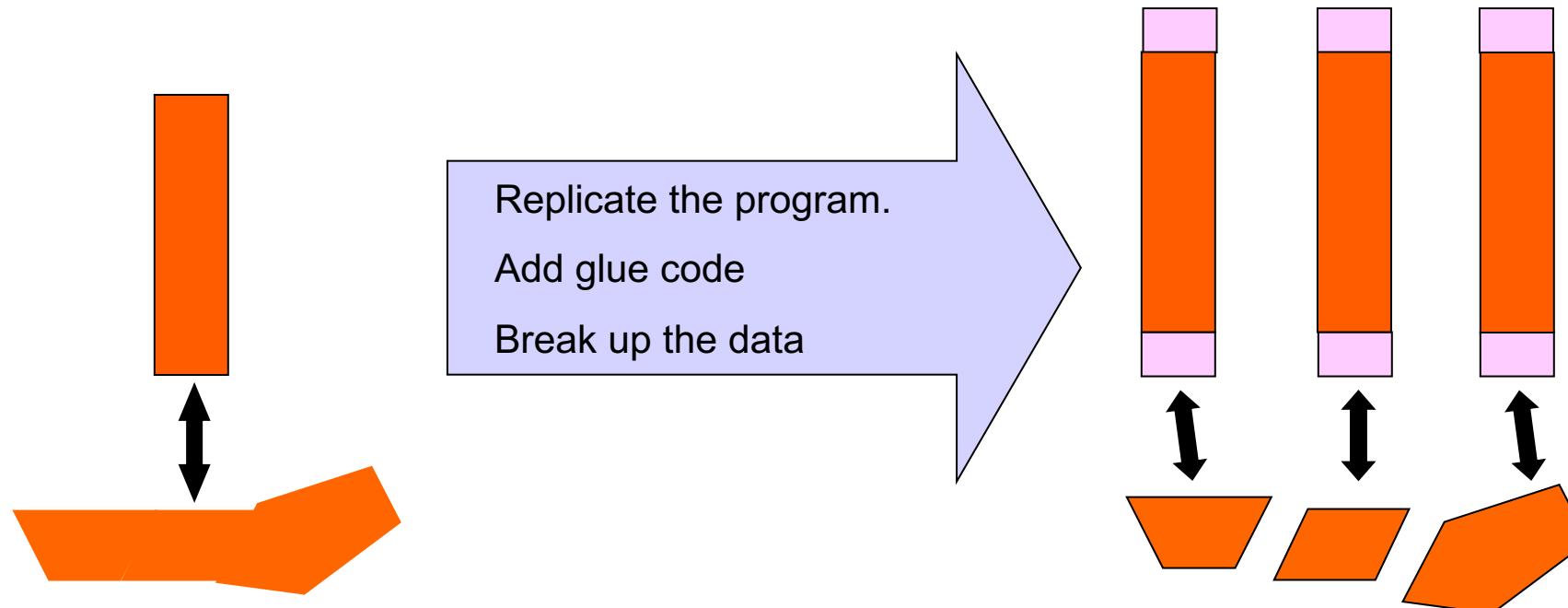


Each rectangle: width Δx , height $F(x_i)$ at i^{th} interval midpoint.

```
def piFunc(NumSteps):
    step=1.0/NumSteps
    pisum = 0.0
    x = 0.5
    for i in range(NumSteps):
        x+=step
        pisum += 4.0/(1.0+x*x)
    pi=step*pisum
    return pi
```

SPMD (Single Program Multiple Data) design pattern

- Run the same program on P processing elements where P can be arbitrarily large.
- Use the rank ... an ID ranging from 0 to $(P-1)$... to select between a set of tasks and to manage any shared data structures.



This pattern is very general and has been used to support most (if not all) the algorithm strategy patterns.

MPI programs almost always use this pattern ... it is probably the most commonly used pattern in the history of parallel programming.

Single Program Multiple Data (SPMD)

```
from numba import njit
import numpy as np
from numba.openmp import openmp_context as openmp
from numba.openmp import omp_get_thread_num, omp_get_num_threads
MaxTHREADS = 32
@njit
def piFunc(NumSteps):
    step = 1.0/NumSteps
    partialSums = np.zeros(MaxTHREADS)
    with openmp("parallel shared(partialSums,numThrds) private(threadID,i,x,localSum)"):
        threadID = omp_get_thread_num()
        with openmp("single"):
            numThrds = omp_get_num_threads()
            localSum = 0.0
            for i in range(threadID, NumSteps, numThrds):
                x = (i+0.5)*step
                localSum = localSum + 4.0/(1.0 + x*x)
            partialSums[threadID] = localSum
    return step*np.sum(partialSums)
```

```
pi = piFunc(100000000)
```

- **omp_get_num_threads()**: get N=number of threads
- **omp_get_thread_num()**: thread rank = 0...(N-1)
- **single**: One thread does the work, others wait
- **private(x)**: each threads gets its own x
- **shared(x)**: all threads see the same x

Deal out loop iterations as if a deck of cards (a cyclic distribution)
... each threads starts with the Iteration = ID, incremented by the
number of threads, until the whole “deck” is dealt out.

The data environment seen by OpenMP threads

- The data environment is the collection of variables visible to the threads in a team.
- Variables can be **shared** or **private**.
 - **Shared variable**: A variable that is visible (i.e. can be read or written) to all threads in a team.
 - **Private variable**: A variable that is only visible to an individual thread.
- All the code associated with an OpenMP directive (such as **parallel** or **for**), including the code in functions called inside that code, is called a **region**. A directive plus code in the immediate block associated with it, is called a **construct**
- Rules for defining a variable as shared or private:
 - A variable is **shared** if it is used before or after an OpenMP construct, otherwise it is **private**.
 - Variables can be made shared or private through clauses included with a directive.

```
from numba import njit
from numba.openmp import openmp_context as openmp

@njit
def piFunc(NumSteps):
    step = 1.0/NumSteps
    pisum = 0.0
    with openmp ("parallel for reduction(+:pisum)"):
        for i in range(NumSteps):
            x = (i+0.5)*step
            pisum += 4.0/(1.0 + x*x)
    pi = step*pisum
    return pi

pi = piFunc(100000000)
```

x first used inside the OpenMP construct ... it is private.

Numerical Integration results in seconds ... lower is better

| Threads | PyOMP | | C | |
|---------|--------|--------|--------|--------|
| | Loop | SPMD | Loop | SPMD |
| 1 | 0.447 | 0.450 | 0.444 | 0.448 |
| 2 | 0.252 | 0.255 | 0.245 | 0.242 |
| 4 | 0.160 | 0.164 | 0.149 | 0.149 |
| 8 | 0.0890 | 0.0890 | 0.0827 | 0.0826 |
| 16 | 0.0520 | 0.0503 | 0.0451 | 0.0451 |

10^8 steps

Intel® Xeon® E5-2699 v3 CPU with 18 cores running at 2.30 GHz.

For the C programs we used Intel®icc compiler version 19.1.3.304 as `icc -qnextgen -O3 -fopenmp`

Ran each case 5 times and kept the minimum time. **JIT time is not included** for PyOMP (it was about 1.5 seconds)

How do we handle problems without such regular structure or with complex load balancing problems?

We do this in OpenMP with explicit tasks

Explicit tasks in PyOMP

- A task is a sequence of statements and an associated data environment. Lots of flexibility in how those tasks are created, so handles irregular parallelism, recursive parallelism, and many other control structures.
- A common pattern ... one thread creates explicit tasks and puts them in a queue. All the threads work together to execute them. The single construct causes one thread to execute statements while the other threads wait at a barrier at the end of the single. It's perfect for task level parallelism.

```
from numba import njit
from numba.openmp import openmp_context as openmp
```

```
@njit
def irregularPar():
    with openmp("parallel"):
        with openmp("single"):
            StateVal = 1
            while (StateVal > 0):
                with openmp("task firstprivate(StateVal)"):
                    BigComp(StateVal)
                    StateVal = ExitYet()
    return
irregularPar()
```

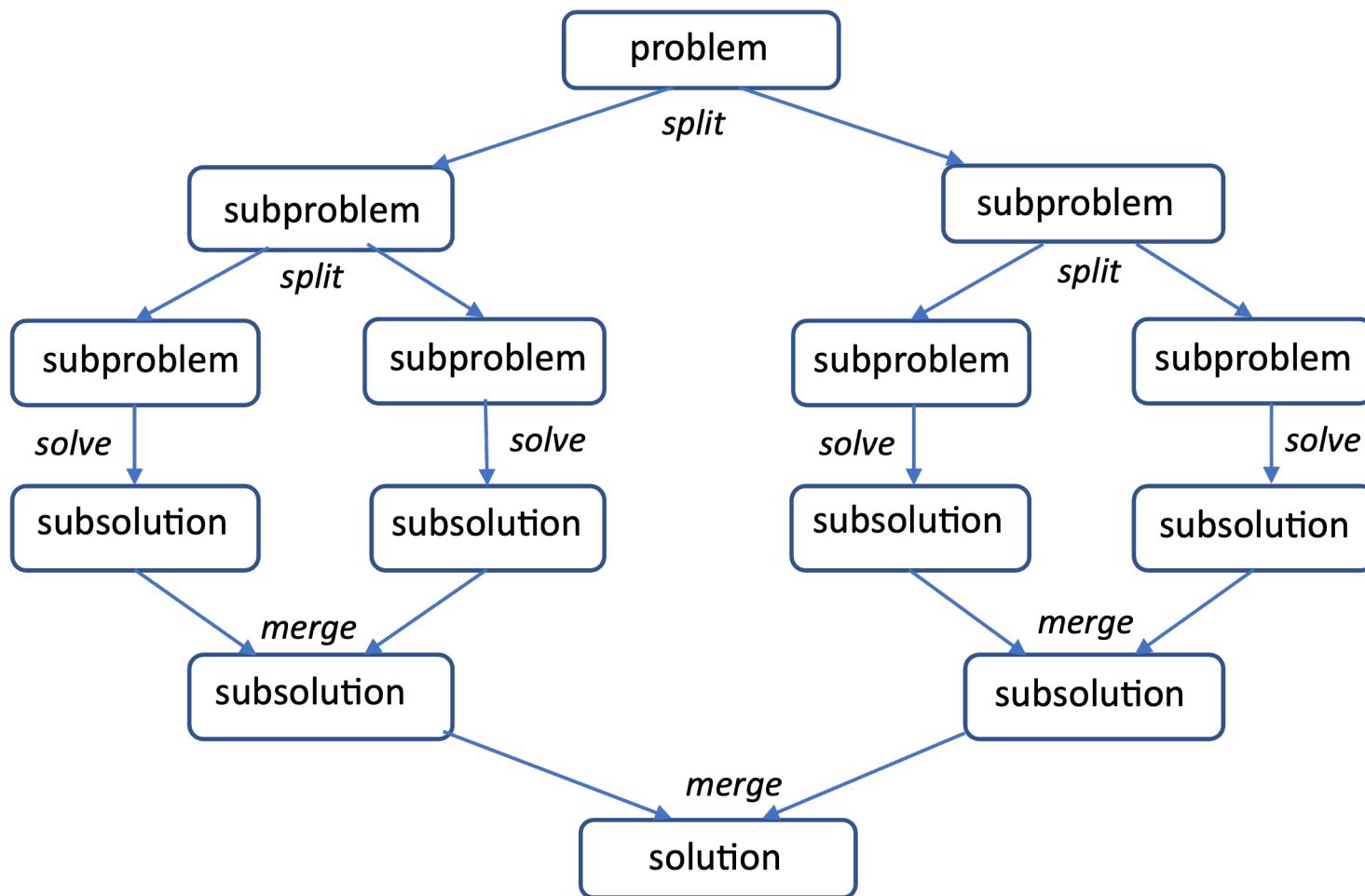
Single: one thread does the work while the other threads wait (and execute tasks) at the barrier implied at the end of single

An explicit task ... captures value of the variable StateVal and calls BigComp.

Returns a negative value at some point (function not shown)

Divide and conquer design pattern

- Split the problem into smaller sub-problems; continue until the sub-problems can be solved directly



3 Options for parallelism:

- Do work as you split into sub-problems
- Do work at the leaves
- Do work as you recombine

Divide and conquer (with explicit tasks)

```
from numba import njit
from numba.openmp import openmp_context as openmp
from numba.openmp import omp_get_num_threads, omp_set_num_threads
MIN_BLK = 1024*256
@njit
def piComp(Nstart, Nfinish, step):
    iblk = Nfinish-Nstart
    if(iblk<MIN_BLK):
        pisum = 0.0
        for i in range(Nstart,Nfinish):
            x= (i+0.5)*step
            pisum += 4.0/(1.0 + x*x)
    else:
        sum1 = 0.0
        sum2 = 0.0
        with openmp ("task shared(sum1)"):
            sum1 = piComp(Nstart, Nfinish-iblk/2,step)
        with openmp ("task shared(sum2)"):
            sum2 = piComp(Nfinish-iblk/2,Nfinish,step)
        with openmp ("taskwait"):
            pisum = sum1 + sum2
return pisum
```

Solve

Split

Merge

```
@njit
```

```
def piFunc(NumSteps):
```

```
    step = 1.0/NumSteps
```

```
    sum = 0.0
```

```
    startTime = omp_get_wtime()
```

```
    with openmp ("parallel"):
```

```
        with openmp ("single"):
```

```
            pisum = piComp(0,NumSteps,step)
```

Fork threads
and launch the
computation

```
pi = step*pisum
```

```
return pi
```

```
pi = piFunc(100000000)
```

- **single**: One thread does the work, others wait
- **task**: code block enqueued for execution
- **taskwait**: wait until task in the code block finish

Numerical Integration results in seconds ... lower is better

| Threads | PyOMP | | | C | | |
|---------|--------|--------|--------|--------|--------|--------|
| | Loop | SPMD | Task | Loop | SPMD | Task |
| 1 | 0.447 | 0.450 | 0.453 | 0.444 | 0.448 | 0.445 |
| 2 | 0.252 | 0.255 | 0.245 | 0.245 | 0.242 | 0.222 |
| 4 | 0.160 | 0.164 | 0.146 | 0.149 | 0.149 | 0.131 |
| 8 | 0.0890 | 0.0890 | 0.0898 | 0.0827 | 0.0826 | 0.0720 |
| 16 | 0.0520 | 0.0503 | 0.0517 | 0.0451 | 0.0451 | 0.0431 |

10^8 steps

Intel® Xeon® E5-2699 v3 CPU with 18 cores running at 2.30 GHz.

For the C programs we used Intel®icc compiler version 19.1.3.304 as `icc -qnextgen -O3 -fopenmp`

Ran each case 5 times and kept the minimum time. **JIT time is not included** for PyOMP (it was about 1.5 seconds)

**There is more But this is enough
to get you started with CPU
programming in PyOMP**

**So let's wrap up our discussion of
CPU programming**

PyOMP subset of OpenMP for CPU programming

| | |
|---|---|
| <code>with openmp("parallel"):</code> | Create a team of threads. Execute a parallel region |
| <code>with openmp("for"):</code> | Use inside a parallel region. Split up a loop across the team. |
| <code>with openmp("parallel for"):</code> | A combined construct. Same a <code>parallel</code> followed by a <code>for</code> . |
| <code>with openmp ("single"):</code> | One thread does the work. Others wait for it to finish |
| <code>with openmp("task"):</code> | Create an explicit task for work within the construct. |
| <code>with openmp("taskwait"):</code> | Wait for all tasks in the current task to complete. |
| <code>with openmp("barrier"):</code> | All threads arrive at a barrier before any proceed. |
| <code>with openmp("critical"):</code> | Mutual exclusion. One thread at a time executes code |
| <code>schedule(static [,chunk])</code> | Map blocks of loop iterations across the team. Use with <code>for</code> . |
| <code>reduction(op:list)</code> | Combine values with op across the team. Used with <code>for</code> |
| <code>private(list)</code> | Make a local copy of variables for each thread. Use with <code>parallel</code> , <code>for</code> or <code>task</code> . |
| <code>firstprivate(list)</code> | <code>private</code> , but initialize with original value. Use with <code>parallel</code> , <code>for</code> or <code>task</code> |
| <code>shared(list)</code> | Variables shared between threads. Use with <code>parallel</code> , <code>for</code> or <code>task</code> . |
| <code>default(None)</code> | Force definition of variables as <code>private</code> or <code>shared</code> . |
| <code>omp_get_num_threads()</code> | Return the number of threads in a team |
| <code>omp_get_thread_num()</code> | Return an ID from 0 to the number of threads minus one |
| <code>omp_set_num_threads(int)</code> | Set the number of threads to request for parallel regions |
| <code>omp_get_wtime()</code> | Return a snapshot of the wall clock time. |
| <code>OMP_NUM_THREADS=N</code> | Environment variable to set the default number of threads |

PyOMP subset of OpenMP for CPU programming

| | | |
|---|---|-------------------|
| <code>with openmp("parallel"):</code> | Create a team of threads. Execute a parallel region | Fork threads |
| <code>with openmp("for"):</code> | Use inside a parallel region. Split up a loop across the team. | |
| <code>with openmp("parallel for"):</code> | A combined construct. Same a <code>parallel</code> followed by a <code>for</code> . | |
| <code>with openmp ("single"):</code> | One thread does the work. Others wait for it to finish | Work sharing |
| <code>with openmp("task"):</code> | Create an explicit task for work within the construct. | |
| <code>with openmp("taskwait"):</code> | Wait for all tasks in the current task to complete. | |
| <code>with openmp("barrier"):</code> | All threads arrive at a barrier before any proceed. | Synchronization |
| <code>with openmp("critical"):</code> | Mutual exclusion. One thread at a time executes code | |
| <code>schedule(static [,chunk])</code> | Map blocks of loop iterations across the team. Use with <code>for</code> . | |
| <code>reduction(op:list)</code> | Combine values with op across the team. Used with <code>for</code> | Par. Loop support |
| <code>private(list)</code> | Make a local copy of variables for each thread. Use with <code>parallel</code> , <code>for</code> or <code>task</code> . | |
| <code>firstprivate(list)</code> | <code>private</code> , but initialize with original value. Use with <code>parallel</code> , <code>for</code> or <code>task</code> | |
| <code>shared(list)</code> | Variables shared between threads. Use with <code>parallel</code> , <code>for</code> or <code>task</code> | Data Environment |
| <code>default(None)</code> | Force definition of variables as <code>private</code> or <code>shared</code> . | |
| <code>omp_get_num_threads()</code> | Return the number of threads in a team | |
| <code>omp_get_thread_num()</code> | Return an ID from 0 to the number of threads minus one | |
| <code>omp_set_num_threads(int)</code> | Set the number of threads to request for parallel regions | |
| <code>omp_get_wtime()</code> | Return a snapshot of the wall clock time. | runtime libraries |
| <code>OMP_NUM_THREADS=N</code> | Environment variable to set the default number of threads | Environment |

The view of Python from an HPC perspective

```
for I in range(4096):  
    for j in range(4096):  
        for k in range (4096):  
            C[i][j] += A[i][k]*B[k][j]
```

We know better ...
the IKJ order is more
cache friendly

And we picked a
smaller problem

```
for I in range(1000):  
    for k in range(1000):  
        for j in range (1000):  
            C[i][j] += A[i][k]*B[k][j]
```

Table 1. Speedups from performance engineering a program that multiplies two 4096-by-4096 matrices. Each version represents a successive refinement of the original Python code. “Running time” is the running time of the version. “GFLOPS” is the billions of 64-bit floating-point operations per second that the version executes. “Absolute speedup” is time relative to Python, and “relative speedup,” which we show with an additional digit of precision, is time relative to the preceding line. “Fraction of peak” is GFLOPS relative to the computer’s peak 835 GFLOPS. See Methods for more details.

| Version | Implementation | Running time (s) | GFLOPS | Absolute speedup | Relative speedup | Fraction of peak (%) |
|---------|-----------------------------|------------------|---------|------------------|------------------|----------------------|
| 1 | Python | 25,552.48 | 0.005 | 1 | — | 0.00 |
| 2 | Java | 2,372.68 | 0.058 | 11 | 10.8 | 0.01 |
| 3 | C | 542.67 | 0.253 | 47 | 4.4 | 0.03 |
| 4 | Parallel loops | 69.80 | 1.969 | 366 | 7.8 | 0.24 |
| 5 | Parallel divide and conquer | 3.80 | 36.180 | 6,727 | 18.4 | 4.33 |
| 6 | plus vectorization | 1.10 | 124.914 | 23,224 | 3.5 | 14.96 |
| 7 | plus AVX intrinsics | 0.41 | 337.812 | 62,806 | 2.7 | 40.45 |

PyOMP DGEMM (Mat-Mul with double precision numbers)

```
from numba import njit
import numpy as np
from numba.openmp import openmp_context as openmp
from numba.openmp import omp_get_wtime

@njit(fastmath=True)
def dgemm(iterations,order):

    # allocate and initialize arrays
    A = np.zeros((order,order))
    B = np.zeros((order,order))
    C = np.zeros((order,order))

    # Assign values to A and B such that
    # the product matrix has a known value.
    for i in range(order):
        A[:,i] = float(i)
        B[:,i] = float(i)
```

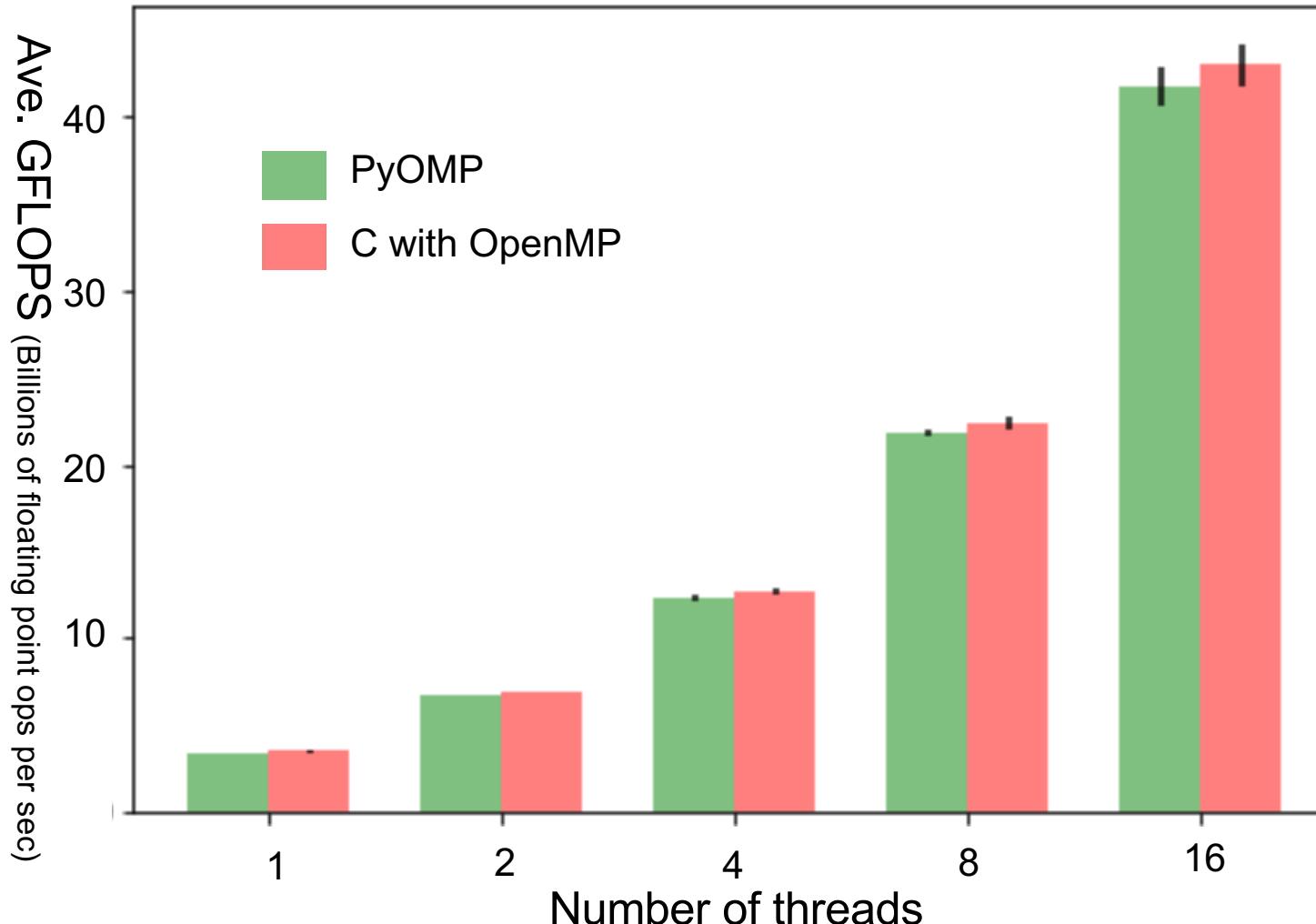
```
tInit = omp_get_wtime()
with openmp("parallel for private(j,k)"):
    for i in range(order):
        for k in range(order):
            for j in range(order):
                C[i][j] += A[i][k] * B[k][j]

dgemmTime = omp_get_wtime() - tInit

# Check result
checksum = 0.0;
for i in range(order):
    for j in range(order):
        checksum += C[i][j]
ref_checksum = order*order*order
ref_checksum *= 0.25*(order-1.0)*(order-1.0)
eps=1.e-8
if abs((checksum - ref_checksum)/ref_checksum) < eps:
    print('Solution validates')
nflops = 2.0*order*order*order
print('Rate (MF/s): ',1.e-6*nflops/dgemmTime)
```

DGEMM PyOMP vs C-OpenMP

Matrix Multiplication, double precision, order = 1000, with error bars (std dev)



250 runs for order
1000 matrices

PyOMP times
DO NOT include
the one-time JIT
cost of ~2
seconds.

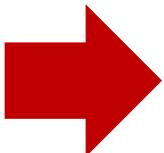
... but remember,
the JIT'ed code
can be cached for
future use. It's
straightforward to
hide the JIT cost.

And we can use PyOMP for GPU programming

The “BIG idea” Behind GPU programming

Traditional Loop based vector addition (vadd)

```
int main() {  
    int N = . . . ;  
    float *a, *b, *c;  
  
    a* =(float *) malloc(N * sizeof(float));  
  
    // ... allocate other arrays (b and c)  
    // and fill with data  
  
    for (int i=0;i<N; i++)  
        c[i] = a[i] + b[i];  
}
```



Data Parallel vadd with CUDA

```
// Compute sum of length-N vectors: C = A + B  
void __global__  
vecAdd (float* a, float* b, float* c, int N) {  
    int i = blockIdx.x * blockDim.x + threadIdx.x;  
    if (i < N) c[i] = a[i] + b[i];  
}  
  
int main () {  
    int N = . . . ;  
    float *a, *b, *c;  
    cudaMalloc (&a, sizeof(float) * N);  
    // ... allocate other arrays (b and c)  
    // and fill with data  
  
    // Use thread blocks with 256 threads each  
    vecAdd <<< (N+255)/256, 256 >>> (a, b, c, N);  
}
```

Assume a GPU with unified shared memory
... allocate on host, visible on device too

How do we execute code on a GPU: The SIMT model (Single Instruction Multiple Thread)

1. Turn source code into a scalar work-item

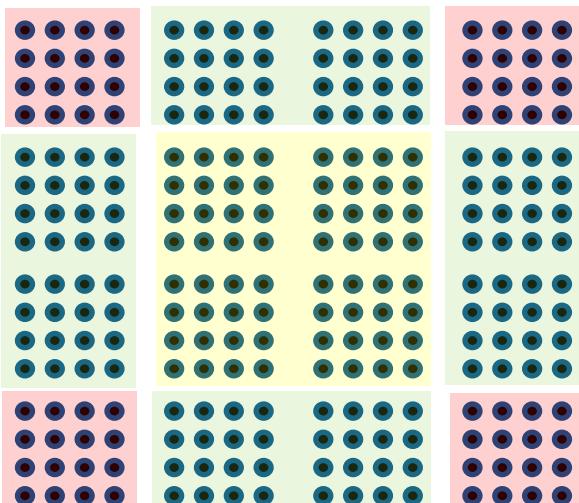
```
// Compute sum of length-N vectors: C = A + B
void __global__
vecAdd (float* a, float* b, float* c, int N) {
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    if (i < N) c[i] = a[i] + b[i];
}

int main () {
    int N = ... ;
    float *a, *b, *c;
    cudaMalloc (&a, sizeof(float) * N);
    // ... allocate other arrays (b and c)
    // and fill with data

    // Use thread blocks with 256 threads each
    vecAdd <<< (N+255)/256, 256 >>> (a, b, c, N);
}
```

This is CUDA code ... the sort of code the OpenMP compiler generates on your behalf

2. Map work-items onto an N dim index space.



4. Run on hardware designed around the same SIMT execution model

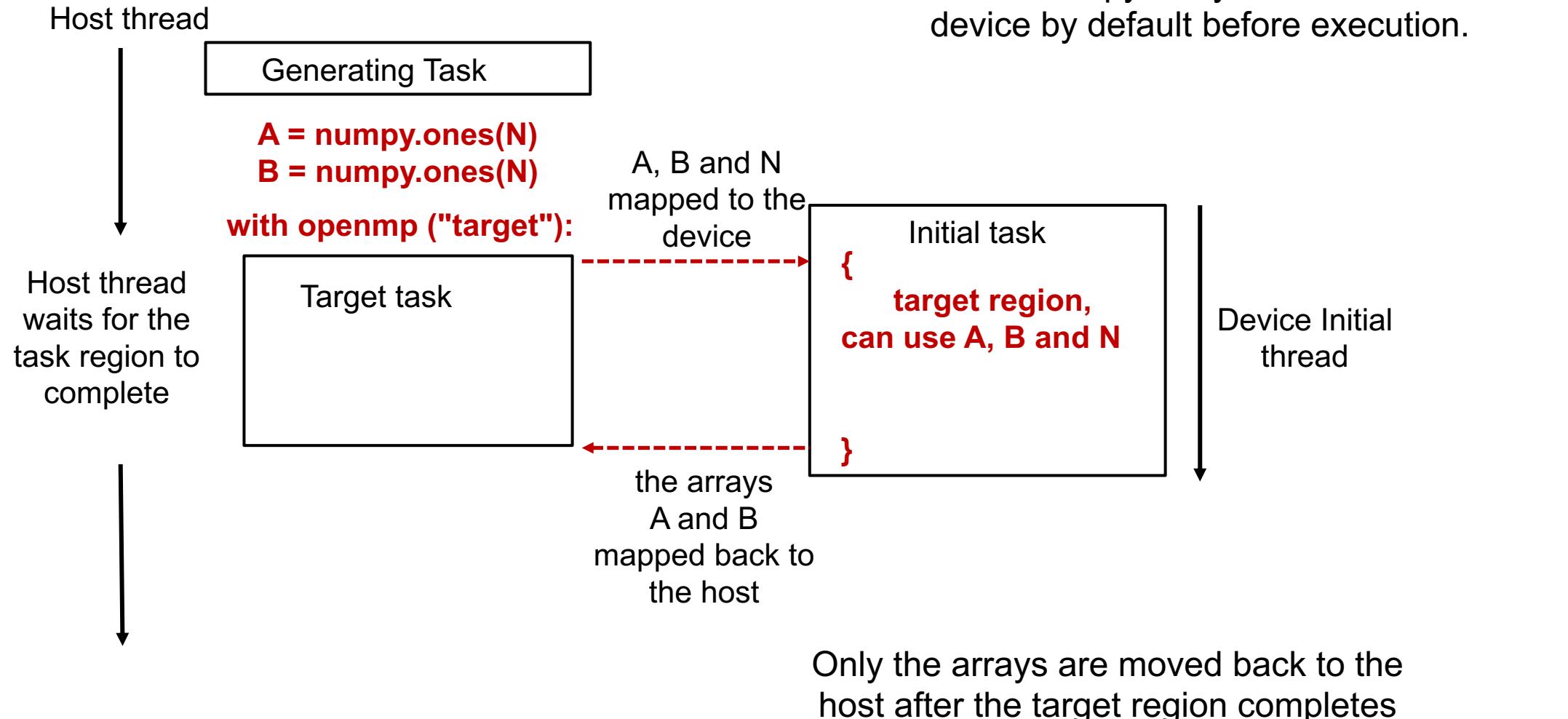


3. Map data structures onto the same index space

Note: The CUDA code defines a 1D grid. I show a 2D grid on this slide to make kernel execution and its relation to data more clear.

**How do we map a loop onto the
GPU execution model in PyOMP?**

Step 1: move code and data onto the GPU: The target construct and default data movement



Step 2: Map loop iterations onto the GPU's SIMD lanes

```
@njit  
def main():  
    N = 1024  
    A = numpy.ones(N)  
    B = numpy.ones(N)  
  
    with openmp ("target "):  
        with openmp ("loop"):  
            for i in range(N):  
                A[i] += B[i]
```

The loop construct tells the compiler:
"this loop will execute correctly if the loop iterations run in any order. You can safely run them concurrently. And the loop-body doesn't contain any OpenMP constructs. So do whatever you can to make the code run fast"

The loop construct is a declarative construct. You tell the compiler what you want done but you DO NOT tell it how to "do it". This is new for OpenMP

Step 2: Map loop iterations onto the GPU's SIMD lanes

```
@njit  
def main():  
    N = 1024  
    A = numpy.ones(N)  
    B = numpy.ones(N)
```

with openmp ("target"):

with openmp ("loop"):

for i in range(N):

A[i] += B[i]

1. Variables created in host memory.

2. Scalar **N** and arrays **A** and **B** are copied to device memory. Execution transferred to device.

3. For-loop index variables (such as **i**) are **private** in openmp regions

4. Loop iterations define the index space, work-items, and work-groups.

5. After the OpenMP construct, arrays **A** and **B** are copied *from* device memory back to the host. Host resumes execution.

Difference from OpenMP/C: PyOMP only has NumPy arrays, which carry size information. So, PyOMP arrays sent in full by default ... as it is with C static-arrays.

Loop Parallelism code naturally maps onto the CPU

```
from numba import njit
import numpy as np
from numba.openmp import openmp_context as openmp

@njit(fastmath=True)
def dgemm(iterations,N):

    # allocate and initialize numpy arrays
    # A, B and C of size N by N. <<< code not shown>>>

with openmp("parallel for private(j,k)"):
    for i in range(N):
        for k in range(N):
            for j in range(N):
                C[i][j] += A[i][k] * B[k][j]
```

OpenMP constructs managed through the *with* context manager.

Create a team of threads. Map loop iterations onto them

- **parallel**: creates a team of threads
- **for**: maps loop iterations onto threads.
- **private(j,k)**: each threads gets its own j and k variables
- Loop control index of a parallel for (**i**) is private to each thread.

Loop Parallelism code naturally maps onto the CPU

```
from numba import njit
import numpy as np
from numba.openmp import openmp_context as openmp

@njit(fastmath=True)
def dgemm(iterations,N):

    # allocate and initialize numpy arrays
    # A, B and C of size N by N. <<< code not shown>>>

with openmp("target teams loop collapse(2) private(j)"):
    for i in range(N):
        for k in range(N):
            for j in range(N):
                C[i][j] += A[i][k] * B[k][j]
```

OpenMP constructs managed through the *with* context manager.

Map the loop onto a 2D index space ... the loop body defines the kernel function

- **target**: map execution from the host onto the device
- **teams loop**: Map kernel instances onto PEs inside the compute units
- **collapse(2)**: combine following two loops into a single iteration space.
- **private(j)**: each threads gets its own j variable
- Indices of parallelized loops (**i,k**) are private to each thread.

**Implicit data movement covers a small subset of
the cases you need in a real program.**

**To be more general ... we need to manage data
movement explicitly**

Implicit data movement

- Previously, we described the rules for *implicit* data movement ... **N**, **A** and **B** moved to the GPU on entry to the target construct. **A** and **B** moved to the CPU on exit from the target construct.
- Notice that in this case, **B** is not changed on the GPU ... moving it is a waste of resources

```
@njit
def main():
    N = 1024
    A = numpy.ones(N)
    B = numpy.ones(N)
```

```
with openmp ("target"):
    for i in range(N):
        A[i] += B[i]
```

Controlling data movement with the map clause

```
@njit
```

```
def main():
```

```
    N = 1024
```

```
    A = numpy.ones(N)
```

```
    B = numpy.ones(N)
```

```
    with openmp ("target map(tofrom: A) map(to: B)"):
```

```
        for i in range(N):
```

```
            A[i] += B[i]
```

map(tofrom: A) Map data at the start and end of **target** region.

map(to: B) map data at the start of **target** region but NOT at the end.

We use the term “map” since depending on the detailed memory architecture of the CPU and the GPU, data may be in a shared address space so copying may not be needed.

PyOMP array notation

- When mapping data arrays, if you only give the array name then PyOMP transfers the entire array (using the NumPy array metadata to determine the size)
- To transfer less than the full array, the array section syntax can be used
 - **array_name[begin:end]**
 - This follows Python/NumPy slicing syntax where **begin** is inclusive but **end** is exclusive.
 $A[N:M]$. In set notation implies elements $[N:M)$
 - Multi-dimensional arrays work as expected when transferred in full. Currently PyOmp doesn't support array-section syntax for multi-dimensional arrays.

C Difference: In C, arrays are usually dynamically allocated and referenced through a pointer. You must use array-section syntax to move data. In C, array-syntax is “(initial-offset: number-of-items)”. Fortran uses “begin:end” syntax (as Python does), but the ending index is inclusive (i.e., [begin:end]).

Controlling data movement: the map clause

- **map(to:list)**: On entering the region, variables in the list are initialized on the device using the original values from the host (host to device copy).
- **map(from:list)**: At the end of the target region, the values from variables in the list are copied into the original variables on the host (device to host copy). On entering the region, the initial value of the variables on the device is not initialized.
- **map(tofrom:list)**: the effect of both a map-to and a map-from (host to device copy at start of region, device to host copy at end).
- **map(alloc:list)**: On entering the region, data is allocated and uninitialized on the device.
- **map(list)**: equivalent to **map(tofrom:list)**.

Note: Data movement is defined from the perspective of the host.

```
@njit
def main():
    a = numpy.ones(N)
    b = numpy.ones(N)
    c = numpy.empty(N)
    with openmp ("target teams loop map(to: a,b) map(tofrom: c)") :
        for i in range(N):
            c[i] = a[i] + b[i]
```

When applied to an array, the mapping mode applies only to the array's data. Array metadata is always transferred as **to** and no operations which would change the metadata (e.g., resize) are permitted.

Going beyond simple vector addition ...

**Using OpenMP for GPU application
programming ... the heat diffusion problem**

5-point stencil: the heat program

- The heat equation models changes in temperature over time.

$$\frac{\partial u}{\partial t} - \alpha \nabla^2 u = 0$$

- We'll solve this numerically on a computer using an explicit **finite difference** discretisation.
- $u = u(t, x, y)$ is a function of space and time.
- Partial differentials are approximated using diamond difference formulae:

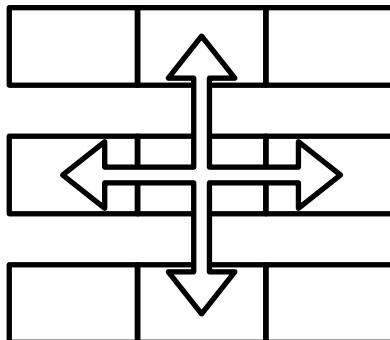
$$\frac{\partial u}{\partial t} \approx \frac{u(t+1, x, y) - u(t, x, y)}{dt}$$

$$\frac{\partial^2 u}{\partial x^2} \approx \frac{u(t, x+1, y) - 2u(t, x, y) + u(t, x-1, y)}{dx^2}$$

- Forward finite difference in time, central finite difference in space.

5-point stencil: the heat program

- Given an initial value of u , and any boundary conditions, we can calculate the value of u at time $t+1$ given the value at time t .
- Each update requires values from the north, south, east and west neighbours only:



- Computation is essentially a weighted average of each cell and its neighbouring cells.
- If on a boundary, look up a boundary condition instead.

Heat diffusion problem ...

```
# Loop over time steps
for _ in range(nsteps):
    # solve over spatial domain for step t
    solve(n, alpha, dx, dt, u, u_tmp)

    # Array swap to get ready for next step
    u, u_tmp = u_tmp, u
```

Array-swap on the host works. Why?

u and u_tmp are references to structs that hold NumPy metadata and a data pointer.

The OpenMP runtime creates a device struct at the target enter data construct and maintains a fixed association between host and device struct references.

Hence, as you swap the array variables, the references to the struct addresses in device memory are swapped.

5-point stencil: solve kernel

```
@njit
def solve(n, alpha, dx, dt, u, u_tmp):
    # Finite difference constant multiplier
    r = alpha * dt / (dx ** 2)
    r2 = 1 - 4 * r
    # Loop over the nxn grid
    for i in range(n):
        for j in range(n):
            # Update the 5-point stencil.
            # Using boundary conditions on the edges of the domain.
            # Boundaries are zero because the MMS solution is zero there.
            u_tmp[j, i] = (r2 * u[j, i] +
                           (u[j, i+1] if i < n-1 else 0.0) +
                           (u[j, i-1] if i > 0    else 0.0) +
                           (u[j+1, i] if j < n-1 else 0.0) +
                           (u[j-1, i] if j > 0    else 0.0))
```

25,000x25,000 grid for 10 time steps
* Xeon Platinum 8480+: 67.6 secs

Solution: parallel stencil (heat)

```
@njit

def solve(n, alpha, dx, dt, u, u_tmp):
    """Compute the next timestep, given the current timestep"""

    # Finite difference constant multiplier
    r = alpha * dt / (dx ** 2)
    r2 = 1 - 4 * r

    with openmp ("target loop collapse(2) map(tofrom: u, u_tmp)"):

        # Loop over the nxn grid
        for i in range(n):
            for j in range(n):
                u_tmp[j, i] = (r2 * u[j, i] +
                               (u[j, i+1] if i < n-1 else 0.0) +
                               (u[j, i-1] if i > 0    else 0.0) +
                               (u[j+1, i] if j < n-1 else 0.0) +
                               (u[j-1, i] if j > 0    else 0.0))
```

25,000x25,00 grid for 10 time steps

- Xeon Platinum 8480+: 67.6 secs
- Nvidia V100: 22.6 secs

Data Movement dominates...

25,000x25,00 grid for 10 time steps

- Xeon Platinum 8480+: 67.6 secs
- Nvidia V100: 22.6 secs

```
# Loop over time steps    Typically, many time steps!
for _ in range(nsteps):
    # solve over spatial domain for step t
    solve(n, alpha, dx, dt, u, u_tmp)
    # Array swap to get ready for next step
    u, u_tmp = u_tmp, u
```

solve() function uses this context:
with openmp ("target loop collapse(2) map(tofrom: u, u_tmp)"):

For each iteration, **copy from** device
 $(2*N^2)*\text{sizeof}(\text{TYPE})$ bytes

- We need to keep data resident on the device *between* target regions
- We need a way to manage the device data environment across iterations.

Target enter/exit data constructs

- The **target data** construct requires a *structured* block of code.
 - Often inconvenient in real codes.
- Can achieve similar behavior with two standalone directives:
with openmp ("target enter data map(..."):
with openmp ("target exit data map(..."):
- The **target enter data** maps variables to the device data environment.
- The **target exit data** unmaps variables from the device data environment.
- Future **target** regions inherit the existing data environment.

Solution: Reference swapping in action

```
with openmp ("target enter data map(to: u, u_tmp)") :  
    pass  
  
    Copy data to device  
    before iteration loop  
  
for _ in range(nsteps) :  
  
    solve(n, alpha, dx, dt, u, u_tmp);  
    Change solve() routine to remove map clauses:  
    with openmp ("target loop collapse(2)")  
  
    # Array swap to get ready for next step  
    u, u_tmp = u_tmp, u
```

```
with openmp ("target exit data map(from: u)") :  
    pass
```

Copy data from device
after iteration loop

- 25,000x25,00 grid for 10 time steps
- Xeon Platinum 8480+ default data movement: 67.6 secs
 - Nvidia V100 default data movement: 22.6 secs
 - Nvidia V100 target enter/exit: 1.2 secs

Target update directive

- You can update data between target regions with the **target update** directive.

Set up the data region ahead of time.

with openmp ("target data map(to: A, B) map(from: C)"):

with openmp ("target"):

{do lots of stuff with A, B and C}

map A on the device to A on the host.

with openmp ("target update from(A)"):

{do something on the host}

with openmp ("target update to(A)"):

pass

map A on the host to A on the device. Note: openmp context body cannot be empty so use "pass"

with openmp ("target"):

{do lots of stuff with A, B and C}

Note: update directive has the transfer direction as the clause: e.g. update to(...) Compare to map clause with direction inside: map(to: ...)

Data movement summary

- Data transfers between host/device occur at:
 - Beginning and end of **target** region
 - Beginning and end of **target data** region
 - At the **target enter data** construct
 - At the **target exit data** construct
 - At the **target update** construct
- Can use **target data** and **target enter/exit data** to reduce redundant transfers.
- Use the **target update** construct to transfer data on the fly within a **target data** region or between **target enter/exit data** directives.

Multi-tasking, Pi program with Dask

```
import numpy as np
import dask

@dask.delayed
def calc_pi(nstart, nstop, step):
    start = (nstart+0.5)*step
    stop = (nstop-0.5)*step
    nsteps = nstop-nstart
    X = np.linspace(start, stop, num=nsteps)
    Y = 4.0 / (1.0 + X*X)
    return np.sum(Y)

def piFunc(NumSteps, NumTasks):
    step = 1.0/NumSteps
    s = 0
    for i in range(NumTasks):
        nstart = (i*NumSteps)//NumTasks
        nstop = ((i+1)*NumSteps)//NumTasks
        s = s + calc_pi(nstart, nstop, step)
    s = s.compute()
    return step*s

if __name__=="__main__":
    from dask.distributed import Client
    client = Client()
    pi = piFunc(100000000, 100)
```

Calculate over part of the range;
Written in Numpy vector style
Faster than Python loops, but use
memory for the arrays X, Y, temps

Start NumTasks tasks,
construct DAG of operations
computing sum

Trigger execution, wait for
completion, get result

Initialize dask “cluster” on local
machine; can provide address
to connect to remote cluster

ParallelAccelerator: loop level parallelism



The Pi program

```
import numba

@numba.njit(parallel=True)
def pi():
    num_steps = 1000000
    step = 1.0 / num_steps
    the_sum = 0.0
    for i in numba.prange(num_steps):
        x = (0.5 + i) * step
        the_sum += 4.0 / (1.0 + x * x)
    pi = step * the_sum
    return pi

print(pi())
```

- ParallelAccelerator includes parallel loops for loop-level parallelism
- The **prange** construct causes equal portions of the iteration space from 0 to num_steps distributed to each core.
- The reduction (the_sum += ...) recognized and implemented safely and efficiently in parallel.

Pi program with PyOMP

```
from numba import njit
from numba.openmp import openmp_context as openmp
```

```
@njit
```

```
def piFunc(NumSteps):
    step = 1.0/NumSteps
    pisum = 0.0
```

```
with openmp ("parallel for private(x) reduction(:pisum)"):
```

```
    for i in range(NumSteps):
        x = (i+0.5)*step
        pisum += 4.0/(1.0 + x*x)
```

```
pi = step*pisum
```

```
return pi
```

```
pi = piFunc(100000000)
```

Pi program

- Single dual-socket server
 - 2x Intel® Xeon® E5-2699v3 @ 2.3Ghz (36 cores, 72 hypercores, total)
 - 128GB RAM
- Mean, stddev of 10 runs (unless stated otherwise), after 1 warmup so we do not time JIT overhead
- For multithreaded runs, we used the default number of threads.

| Num steps | 1e7 | 1e8 | 1e9 | 1e10 | |
|----------------------|-------------------------|-------------------------|------------------------|-----------------------|------------|
| Python loops | 0.92 (0.006) | | | | |
| Numpy | 0.135 (0.005) | 1.45 (0.0015) | | | |
| Numba | 0.039 (0.001) | 0.39 (0.001) | 3.92 (0.003) | | |
| Parallel Accelerator | | 0.019 (0.003) | 0.141 (0.002) | 1.48 (0.077) | |
| Dask | 0.133 (0.008) | 0.75 (0.04) | 6.9 (0.46) | | |
| PyOMP (loop) | 0.041 (0.005) 5 runs | 0.073 (0.005) 5 runs | 0.282 (0.02) 5 runs | 1.56 (0.02) 5 runs | ← Compiled |

All times in seconds

Various Pi programs in Python

```
import numba
from numba import njit
from numba.openmp import openmp_context as openmp
from numba.openmp import omp_get_wtime
import numpy as np

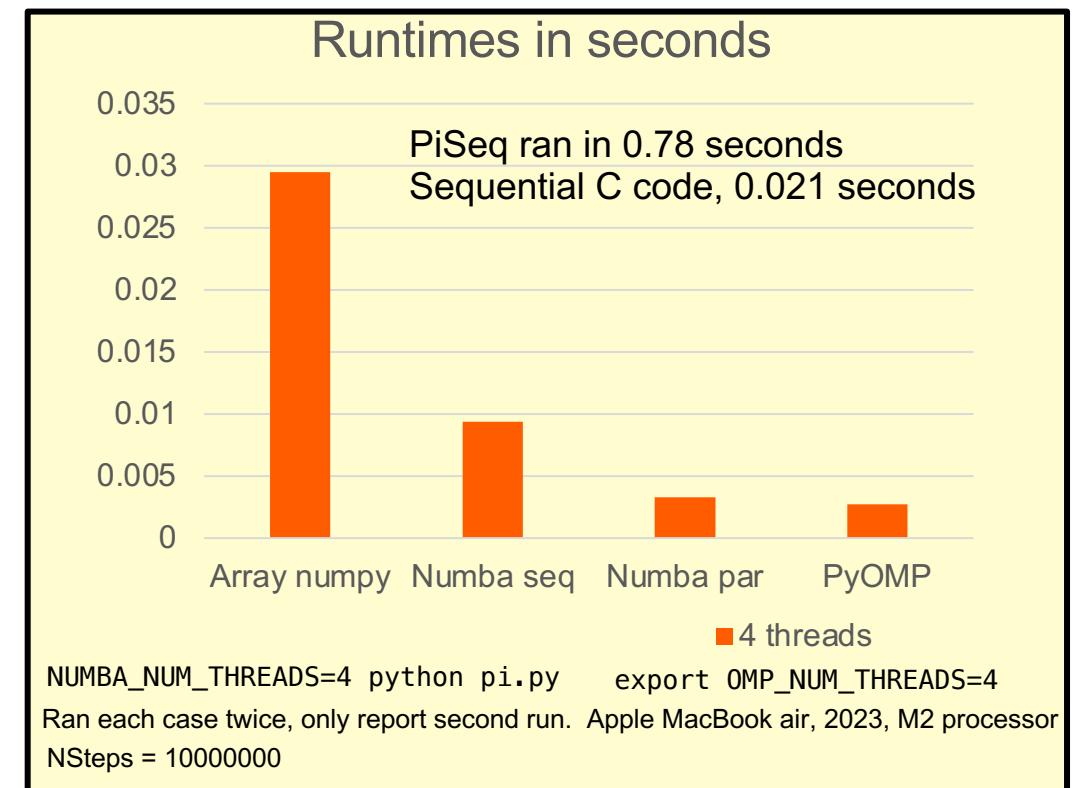
def piArr(Nsteps):
    startTime = omp_get_wtime()
    stepSize = 1.0/Nsteps
    pts = np.linspace(0.0,1.0,Nsteps)
    ptsSquPlus1 = np.square(pts)+1.0
    ptsInteg = 4.0/ptsSquPlus1
    pi = stepSize*np.sum(ptsInteg)
    runtime = omp_get_wtime()-startTime
    return pi,runtime

def piSeq(NumSteps):
    step = 1.0/NumSteps
    sum = 0.0
    startTime = omp_get_wtime()
    for i in range(NumSteps):
        x = (i+0.5)*step
        sum += 4.0/(1.0 + x*x)
    pi = step * sum
    runtime = omp_get_wtime()-startTime
    return pi,runtime
```

```
@njit
def piNUMBAseq(NumSteps):
    step = 1.0/NumSteps
    sum = 0.0
    startTime = omp_get_wtime()
    for i in range(NumSteps):
        x = (i+0.5)*step
        sum += 4.0/(1.0 + x*x)
    pi = step * sum
    runtime = omp_get_wtime()-startTime
    return pi,runtime

@njit(parallel=True)
def piNUMBAPar(NumSteps):
    step = 1.0/NumSteps
    sum = 0.0
    startTime = omp_get_wtime()
    for i in numba.prange(NumSteps):
        x = (i+0.5)*step
        sum += 4.0/(1.0 + x*x)
    pi = step * sum
    runtime = omp_get_wtime()-startTime
    return pi,runtime
```

```
@njit
def piOMP(NumSteps):
    step = 1.0/NumSteps
    sum = 0.0
    startTime = omp_get_wtime()
    with openmp("parallel for reduction(+:sum)"):
        for i in range(NumSteps):
            x = (i+0.5)*step
            sum += 4.0/(1.0 + x*x)
    pi = step * sum
    runtime = omp_get_wtime()-startTime
    return pi,runtime
```



Summary

- Parallel programming is here to stay. If you don't need it today, you will eventually. Fortunately, it's really fun.
- Software outlives hardware. Do not let a vendor lock you in to their platform. Portability must be non-negotiable.
- There are too many parallel programming models for python. Focus on the core principles and fundamental design patterns. Don't wear yourself out chasing the latest fad.



My Greenlandic skin-on-frame kayak in the middle of Budd Inlet during a negative tide

OpenMP Organizations

- OpenMP Architecture Review Board (ARB) URL, the “owner” of the OpenMP specification:

www.openmp.org

- OpenMP User’s Group (cOMPunity) URL:

www.community.org

Get involved, join the ARB and cOMPunity.

Help define the future of OpenMP

Resources

- www.openmp.org has a wealth of helpful resources

The screenshot shows the OpenMP website's "Specifications" page. At the top, the "Specifications" menu item is highlighted in orange. Below the navigation bar, the page title "Specifications" is displayed in large white text on a teal background. In the top right corner of the main content area, there is a breadcrumb trail: "Home > Specifications". The main content area contains two cards: one for the "OpenMP 5.2 Specification" and another for the "OpenMP 5.1 Specification". Each card features a document icon, the specification name, and a bulleted list of related links. A callout bubble from the bottom left points to the "OpenMP API 5.2 Examples" link in the 5.2 specification card.

The OpenMP API specification for parallel programming

Home Specifications Community ▾ Resources ▾ News & Events ▾ About ▾

Specifications

Home > Specifications

OpenMP 5.2 Specification

- OpenMP API 5.2 Specification – Nov 2021
 - Softcover Book on Amazon
- OpenMP API Additional Definitions 2.0 – Nov 2020
- OpenMP API 5.2 Reference Guide (English) (Japanese)
- OpenMP API 5.2 Supplementary Source Code
- OpenMP API 5.2 Examples – April 2022
 - Softcover Book on Amazon
- OpenMP API 5.2 Stack Overflow

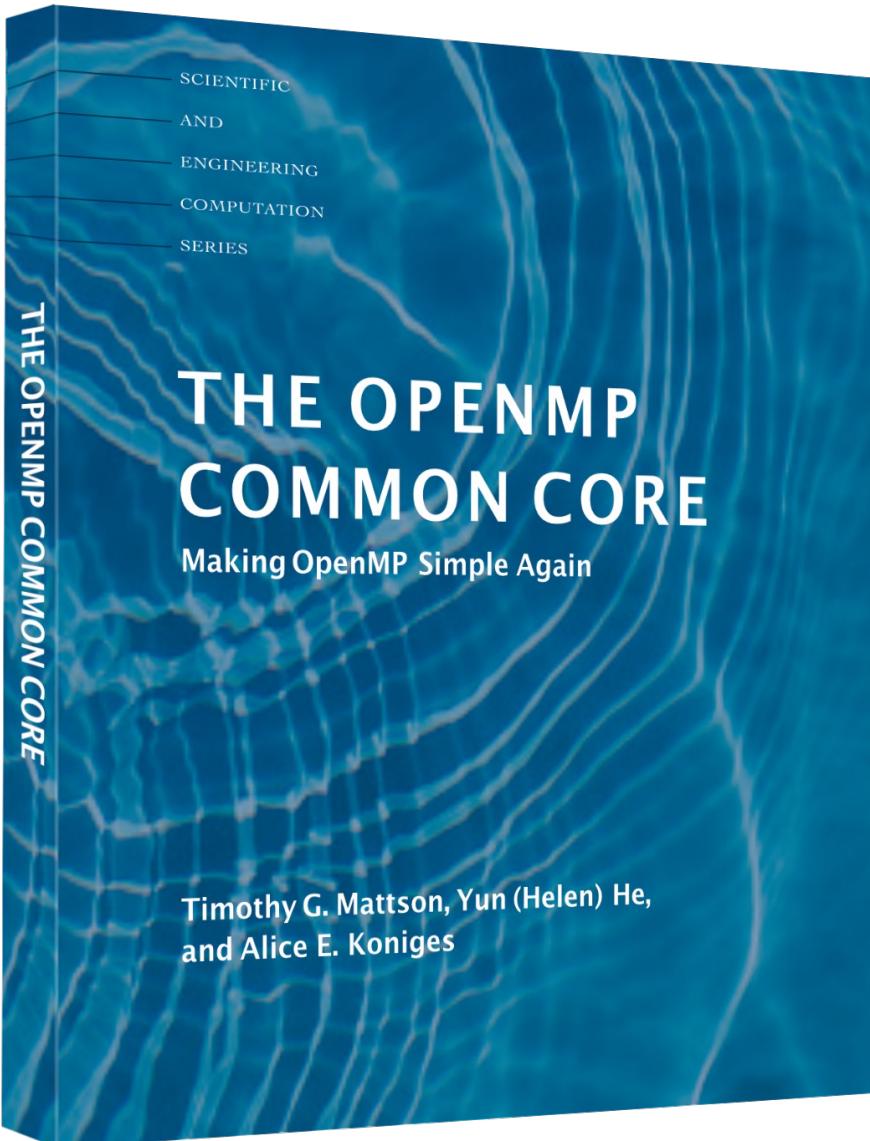
OpenMP 5.1 Specification

- OpenMP API 5.1 Specification – Nov 2020
 - HTML Version
 - Softcover Book on Amazon
- OpenMP API Additional Definitions 2.0 – Nov 2020
- OpenMP API 5.1 Reference Guide
- OpenMP API 5.1 Supplementary Source Code
- OpenMP API 5.1 Examples – August 2021
- OpenMP API 5.1 Stack Overflow

Including a comprehensive collection of examples of code using the OpenMP constructs

To learn OpenMP:

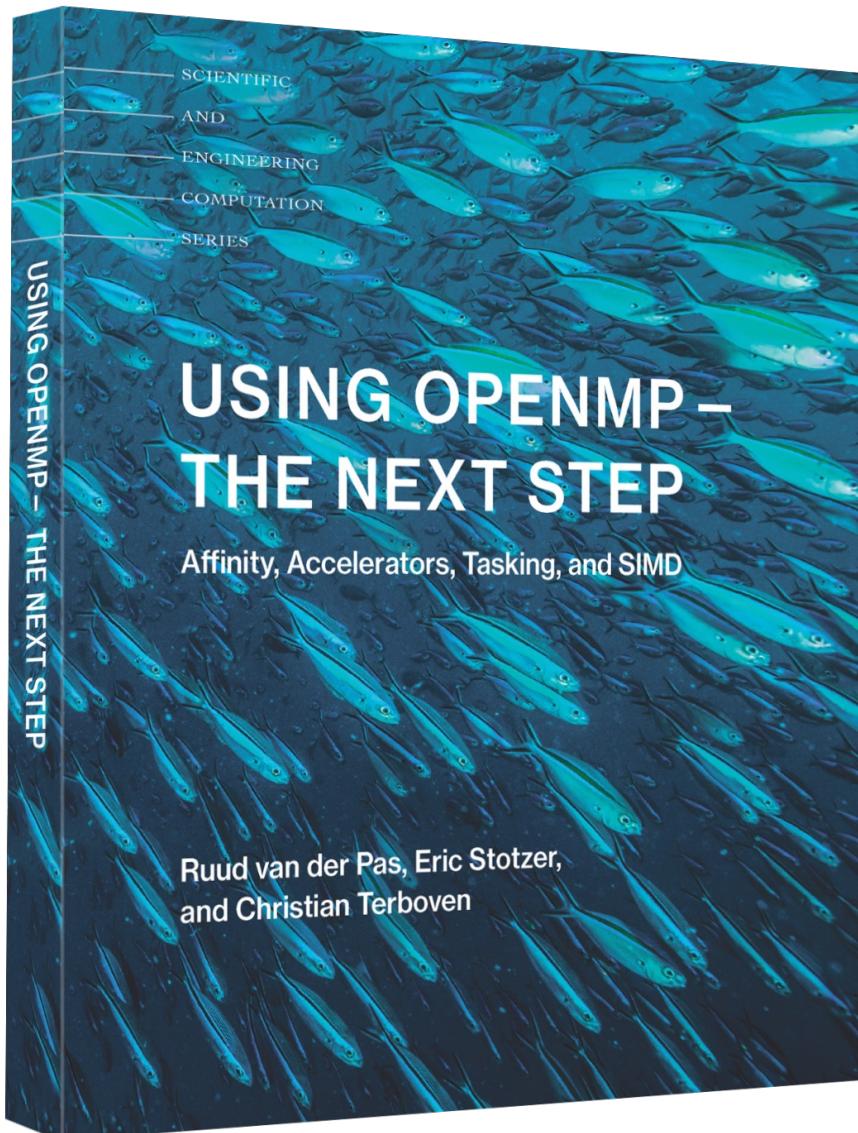
- An exciting new book that Covers the Common Core of OpenMP plus a few key features beyond the common core that people frequently use
- It's geared towards people learning OpenMP, but as one commentator put it ... **everyone at any skill level should read the memory model chapters.**
- Available from MIT Press



www.ompcore.com for code samples and the Fortran supplement

Books about OpenMP

A great book that covers
OpenMP features beyond
OpenMP 2.5



Books about OpenMP

The latest book on OpenMP ...

Now available at amazon.com and
MIT press.

A book about how to use OpenMP to
program a GPU.

